# Multi-Agent Reinforcement Learning for Problems with Combined Individual and Team Reward

Hassam Ullah Sheikh
*Department of Computer Science*
*University of Central Florida*
Orlando, USA
hassam.sheikh@knights.ucf.edu

Ladislau Bölöni
*Department of Computer Science*
*University of Central Florida*
Orlando, USA
lboloni@cs.ucf.edu

*Abstract*—**Many cooperative multi-agent problems require agents to learn individual tasks while contributing to the collective success of the group. This is a challenging task for current state-of-the-art multi-agent reinforcement algorithms that are designed to either maximize the global reward of the team or the individual local rewards. The problem is exacerbated when either of the rewards is sparse leading to unstable learning. To address this problem, we present** *Decomposed Multi-Agent Deep Deterministic Policy Gradient (DE-MADDPG)*: **a novel cooperative multi-agent reinforcement learning framework that simultaneously learns to maximize the global and local rewards. We evaluate our solution on the challenging defensive escort team problem and show that our solution achieves a significantly better and more stable performance than the direct adaptation of the MADDPG algorithm.**

*Index Terms*—**Multi-Agent Reinforcement Learning; Coordination and Collaboration; Dual-Reward Learning**

## I. INTRODUCTION

Cooperative multi-agent problems are prevalent in real-world settings such as strategic conflict resolution [1], coordination between autonomous vehicles [2] and collaboration of agents in defensive escort teams [3]. Such problems can be modelled as dual-interest: each agent is simultaneously working towards maximizing its own payoff (local reward) as well as the collective success of the team (global reward). For example, autonomous vehicles in double-lane merge conflicts must perform cooperative maneuvers without diverging from their destination-bound nominal trajectories. Similarly, in the case of a defensive escort team, each agent has to maintain a specific distance from the payload to avoid disrupting any social norms without sacrificing the security of the payload. Despite the recent success of multi-agent reinforcement learning (MARL) in multiplayer games like Dota 2 [4], Quake III Capture-the-Flag [5] and Starcraft II [6] or learning to use tools [7], learning multi-agent cooperation while simultaneously maximizing local rewards is still an open challenge. In this learning problem, to which we will refer as "**dual-reward MARL**", the agents are *explicitly* receiving two reward signals: the global team reward and the agent's individual local reward.

Current state-of-the-art MARL algorithms can be categorized in two types. For algorithms such as COMA [8] and QMIX [9], the goal is to maximize the global reward for the success of the group while algorithms such MADDPG [10] and M3DDPG [11] focus on optimizing local rewards without any explicit notion of coordination. As shown in [12] and in our findings in Section V, a direct adaptation of these algorithms to dual-reward problems often leads to poor performance and unstable learning. Generally, these adaptations happen in the reward function space where the local and the global reward signals are combined to form an entangled multi-objective reward function [13]. This coupling of reward functions leads to two problems. First, the entangled reward function becomes unfactorizable during training, causing the learning to oscillate between optimizing either the global or the local reward leading to a sub-optimal and unstable solution. This problem is exacerbated when either of the rewards is sparse, thus, leading to a bias towards the other. The second problem is that maximizing the entangled reward function does not correspond to maximizing the objective function.

To address these issues, we present *Decomposed Multi-Agent Deep Deterministic Policy Gradient (DE-MADDPG)*: a novel cooperative multi-agent reinforcement learning framework built on top of deterministic policy gradients that simultaneously learns to maximize the global and the local rewards without the need of creating an entangled multi-objective reward function. The core idea behind DE-MADDPG is to train two critics. The *global critic*, shared between all the cooperating agents takes as input the observations and actions of these agents and estimates the sum of the global expected reward. The *local critic* receives as input only the observation and action of the particular agent and estimates the sum of local expected reward. The advantage of training two critics is that the step of designing an entangled multi-objective reward function can be skipped altogether.

To summarize, our contributions in this paper are the following:

- We develop a dual-critic framework for multi-agent reinforcement learning that learns to simultaneously maximize the decomposed global and local rewards.
- Taking advantage of the decomposition, we treat the global critic as a single-agent critic. This allows us to

apply performance enhancement techniques such as Prioritized Experience Replay (PER) [14] and Twin Delayed Deep Deterministic Policy Gradients (TD3) [15] to tackle the overestimation bias problem in Q-functions. This was not previously feasible in the multi-agent RL setting.

- We evaluate our proposed solution on the defensive escort team problem [3], [16] (see Figure 1) and show that it achieves a significantly better and more stable performance than the direct adaptation of the MADDPG algorithm.

## II. RELATED WORK

Early theoretical work in MARL was limited to discrete state and action spaces [17], [18], [19]. Recent work have adopted techniques from single-agent deep RL to develop general algorithms for high-dimensional continuous space environments requiring complex agent interactions [1], [20], [10].

Cooperative multi-agent learning is important since many real-world problems can be formulated as distributed systems with decentralized agents that must coordinate to achieve shared objectives [21]. Similar to our work, [22] have shown that agents whose rewards depend on all agents' success perform better than agents who optimize for their own success. In the special case when all agents have a single goal and share a global reward, COMA [8] uses a counterfactual baseline. However, the centralized critic in these methods only focuses on optimizing the collective success of the group. When a global objective is the sum of agents' individual objectives, value-decomposition methods optimize a centralized Q-function while preserving scalable decentralized execution [9], but do not address credit assignment. While MADDPG [10] and M3DDPG [11] apply to agents with different reward functions, they do not specifically address the need for cooperation; in fact, they do not distinguish the problems of cooperation and competition, despite the fundamental difference.

To the best of our knowledge, dual-reward MARL was not explicitly addressed in the existing literature. Among the related problems, [23] explored the multi-goal problem and analyzed its convergence in a special networked setting restricted to fully-decentralized training. In contrast, we are conducting centralized training with decentralized execution. In contrast to multi-*task* MARL, which aims for generalization among *non-simultaneous* tasks [24], and in contrast to hierarchical methods with top-level managers that *sequentially* select subtasks [25], our decentralized agents must cooperate *in parallel* to successfully achieve the global and their respective local objectives.

## III. BACKGROUND

### A. Policy Gradients

Policy gradient methods have been shown to learn the optimal policy in a variety of reinforcement learning tasks. The main idea behind policy gradient methods is that instead of parameterizing the Q-function to extract the policy, we parameterize the policy using the parameters $\theta$ to maximize the objective represented as $J(\theta) = \mathbb{E}\left[\mathbb{R}^t\right]$ by taking a step in the direction

$$\nabla J(\theta) = \mathbb{E}\left[\nabla_\theta \log \pi_\theta\left(a|s\right) Q^\pi\left(s, a\right)\right]$$

Policy gradient methods are prone to the high variance problem. Several methods such as [26], [27] have been shown to reduce the variability by introducing a *critic*, a Q-function that tells about the goodness of a reward by working as a baseline. [28] has shown that it is possible to extend the policy gradient framework to deterministic policies *i.e.* $\pi_\theta : \mathcal{S} \rightarrow \mathcal{A}$. In particular we can write $\nabla J(\theta)$ as

$$\nabla J(\theta) = \mathbb{E}\left[\nabla_\theta \pi\left(a|s\right) \nabla_a Q^\pi\left(s, a\right)|_{a=\pi(s)}\right]$$

A variation of this model, Deep Deterministic Policy Gradients (DDPG) [29] is an off-policy algorithm that approximates the policy $\pi$ and the critic $Q^\pi$ with deep neural networks. DDPG uses an experience replay buffer alongside a target network to stabilize the training. Twin Delayed Deep Deterministic Policy Gradients (TD3) [15] improves on DDPG by addressing the overestimation bias of the Q-function, similarly to Double Q-learning [30]. They find that approximation errors of the neural network, combined with gradient descent make DDPG tend to overestimate the Q-values, leading to a slower convergence. TD3 addresses this by using two Q-networks $Q_{\psi_1}, Q_{\psi_2}$, along with two target networks. The Q-functions are updated with the target $y = r^t + \gamma \min_{1,2} Q_{\psi_i'}(s'^t, a'^t)$, while updating the policy with $Q_{\psi_1}$. Additionally, they introduce target policy smoothing by adding noise in the determination of the next action for the critic target $a'^t = \mu_{\theta_\pi'}(s') + \epsilon$, with $\epsilon$ being clipped Gaussian noise $\epsilon = \texttt{clip}(\mathcal{N}(0, \sigma), -c, c)$, where $c$ is a tunable parameter. Additionally, they use delayed policy updates and only update the policy $\pi$ and target network parameters once every $d$ critic updates.

Multi-agent deep deterministic policy gradients (MADDPG) [10] extends DDPG for the multi-agent setting where each agent has it's own policy. The gradient of each policy is written as

$$\nabla J(\theta_i) = \mathbb{E}\left[\nabla_{\theta_i} \pi_i\left(a_i|o_i\right) \nabla_{a_i} Q_i^\pi\left(s, a_1, \ldots, a_N\right)|_{a_i=\pi_i(o_i)}\right]$$

where $s = (o_1, \ldots, o_N)$ and $Q_i^\pi(s, a_1, \ldots, a_N)$ is a centralized action-value function that takes the actions of all the agents in addition to the state of the environment to estimate the Q-value for agent $i$. Since every agent has its own Q-function, the model allows the agents to have different action space and reward functions. The primary insight behind MADDPG is that knowing all the actions of other agents makes the environment stationary, even though their policy changes.

Very recently [31] proposed Multi-Agent TD3 (MATD3) extending MADDPG by replacing the deterministic policy gradients with twin delayed deterministic policy gradient to tackle the overestimation bias problem.

## IV. DECOMPOSED MULTI-AGENT DEEP DETERMINISTIC POLICY GRADIENT

We propose *Decomposed Multi-Agent Deep Deterministic Policy Gradient*: a multi-agent deep reinforcement learning
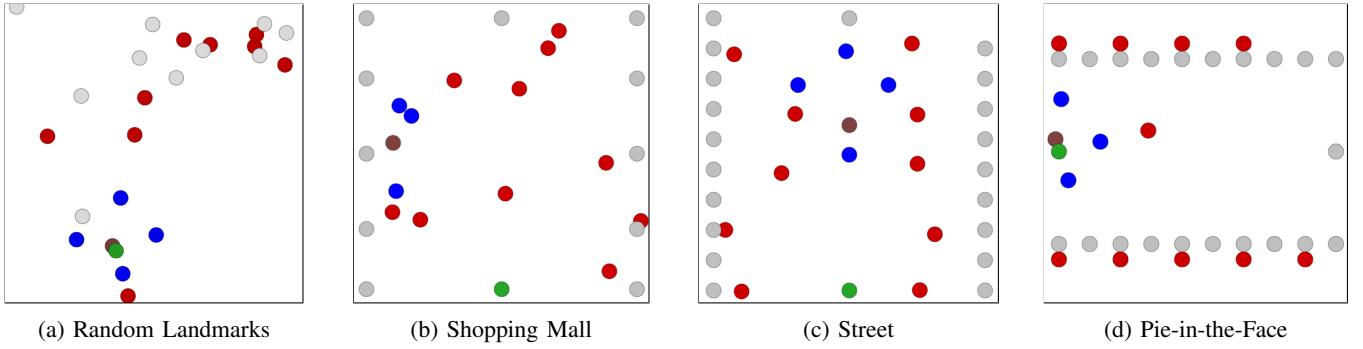
Fig. 1: Four environments for the defensive escort team problem [16]. The team of bodyguards (blue) need to protect the VIP (brown) in the environments, from left to right: "Random Landmarks", "Shopping Mall", "Street" and "Pie-in-the-Face".

(a) Random Landmarks    (b) Shopping Mall    (c) Street    (d) Pie-in-the-Face
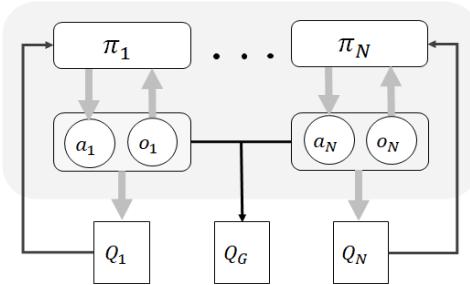


Fig. 2: An overview of the Decomposed Multi-Agent Deep Deterministic Policy Gradient architecture. In contrast to MADDPG which uses a single centralized critic, DE-MADDPG has a global centralized critic shared between all cooperating agents and a local critic specific to the agent.

algorithm that learns to simultaneously maximize the group's global reward and the agent's local rewards. Our approach uses a two critic approach to train policies and value functions that are optimal to maximize the global and local rewards respectively.

The main idea is to combine MADDPG (or MATD3) for maximizing global rewards with a standard single agent DDPG (or TD3). Intuitively, the goal is to move the policy in the direction that maximizes both the global and the local critic. The resulting learning paradigm is similar to the centralized training with decentralized execution during testing used by [10]. In this setting, additional information is provided for the agents during training that is not available during test time.

Concretely, we consider an environment with $N$ agents with policies $\pi = \{\pi_1, \ldots, \pi_N\}$ parameterized by $\theta = \{\theta_1, \ldots, \theta_N\}$. The *multi-agent deep deterministic policy gradient* for agent $i$ can written as

$$\nabla J(\theta_i) = \mathbb{E}\left[\nabla_{\theta_i} \pi_i \left(a_i | o_i\right) \nabla_{a_i} Q_i^\pi \left(s, a_1, \ldots, a_N\right)|_{a_i = \pi_i(o_i)}\right]$$

where $s = (o_1, \ldots, o_N)$ and $Q_i^\pi (s, a_1, \ldots, a_N)$ is a centralized action-value function parameterized by $\phi_i$ that takes the actions of all the agents in addition to the state of the environment to estimate the Q-value for agent $i$. We extend

the idea of MADDPG by introducing a local critic. Now the modified policy gradient for each agent $i$ can be written as

$$\nabla J(\theta_i) = \mathbb{E}_{s,a\sim\mathcal{D}} \overbrace{\left[\nabla_{\theta_i} \pi_i \left(a_i | o_i\right) \nabla_{a_i} Q_\psi^g \left(s, a_1, \ldots, a_N\right)\right]}^{MADDPG}$$
$$\left. + \mathbb{E}_{o_i,a_i\sim\mathcal{D}} \left[\nabla_{\theta_i} \pi_i \left(a_i | o_i\right) \nabla_{a_i} Q_i^\pi \left(o_i, a_i\right)\right] \right\} DDPG$$
$$(1)$$

where $a_i = \pi_i(o_i)$ is action from agent $i$ following policy $\pi_i$ and $\mathcal{D}$ is the experience replay buffer. The global critic is $Q_\psi^g$ is updated as:

$$\mathcal{L}(\psi) = \mathbb{E}_{s,a,r,s'}\left[\left(Q_\psi^g (s, a_1, \ldots, a_N) - y_g\right)^2\right]$$

where $y_g$ is defined as:

$$y_g = r_g + \gamma Q_{\psi'}^g (s', a_1', \ldots, a_N')|_{a_i' = \pi_i'(o_i')}$$

where $\pi' = \{\pi_1', \ldots, \pi_N'\}$ are target policies parameterized by $\theta' = \{\theta_1', \ldots, \theta_N'\}$. Similarly, The local critic is $Q_i^\pi$ is updated as:

$$\mathcal{L}(\phi_i) = \mathbb{E}_{o,a,r,o'}\left[\left(Q_i^\pi (o_i, a_i) - y_l\right)^2\right]$$

where $y_l$ is defined as:

$$y_l = r_l^i + \gamma Q_{\phi_i'}^{\pi'} (o_i', a_i')|_{a_i' = \pi_i'(o_i')}$$

Overestimation bias in Q-functions have been thoroughly studied in [15], [31]. This overestimation bias can be problemsome in multi-agent settings especially in real time autonomous systems. For example, in resolving the double lane merge conflict in autonomous vehicles, the vehicles might consider the current state to be near conflict resolution thus taking a dangerous turns. To solve this problem [15] have proposed a double critic approach to minimize the overestimation bias. Motivated from the results in [15], we replace the Multi-Agent Deterministic Policy Gradient of the global critic in Equation (1) with Twin Delayed Deterministic Policy Gradient. Therefore our updated policy gradient becomes
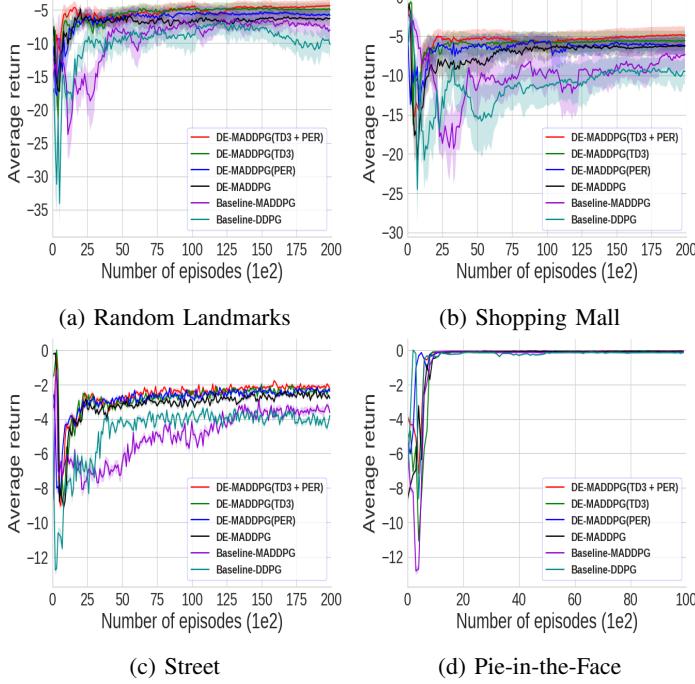
(a) Random Landmarks

(b) Shopping Mall

(c) Street

(d) Pie-in-the-Face

Fig. 3: Learning curves representing the average cumulative global reward. The higher reward represents higher protection to the VIP (payload).

$$\nabla J(\theta_i) = \mathbb{E}_{s,a\sim\mathcal{D}}\left[\nabla_{\theta_i}\pi_i\left(a_i|o_i\right)\nabla_{a_i}Q_{\psi_1}^{g_1}\left(s, a_1, \ldots, a_N\right)\right]$$
$$+ \mathbb{E}_{o_i,a_i\sim\mathcal{D}}\left[\nabla_{\theta_i}\pi_i\left(a_i|o_i\right)\nabla_{a_i}Q_i^{\pi}\left(o_i, a_i\right)\right]$$
(2)

The twin global critics are updated as

$$\mathcal{L}\left(\psi_i\right) = \mathbb{E}_{s,a,r,s'}\left[\left(Q_{\psi_i}^{g_i}\left(s, a_1, \ldots, a_N\right) - y_g\right)^2\right]$$

where $y_g$ is defined as:

$$y_g = r_g + \gamma \min_{i=1,2} Q_{\psi_i'}^{g_i}\left(s', a_1', \ldots, a_N'\right)\big|_{a_i'=\pi_i'(o_i')}$$

Similarly, the local critics can be updated using TD3 update style but for simplicity, we will use the standard DDPG to update the local critics. The overall algorithm to which we refer as *Decomposed Multi-Agent Deterministic Policy Gradient (DE-MADDPG)* is described in Algorithm 1. The overview of the architecture can be seen in Figure 2.

## V. EXPERIMENTS

### A. Environments

We perform our experiments using the defensive escort problem on four VIP protection environments [16], [3]. This is a medium-size collaborative problem where a defensive escort team of agents is learning to maintain an optimal formation

---

**Algorithm 1** Decomposed Multi-agent Deep Policy Gradient

1: Initialize main global critic networks $Q_{\psi}^{g_1}$ and $Q_{\psi}^{g_2}$.
2: Initialize target global critic networks $Q_{\psi'}^{g_1}$ and $Q_{\psi'}^{g_2}$.
3: Initialize each agents policy and critic networks.
4: **for** episode = 1 to $T$ **do**
5:     **for** t = 1 to episode–length **do**
6:         Get environment state $s^t$.
7:         For each agent $i$, select action $a_i^t = \pi_{\theta_i}\left(o_i^t\right)$
8:         Execute actions $\mathbf{a}^t = [a_1^t, \ldots, a_N^t]$
9:         Receive global $r_g^t$ and local rewards $\mathbf{r}_l^t$.
10:         Store $\left(s^t, \mathbf{a}^t, \mathbf{r}_l^t, r_g^t, s^{t+1}\right)$ in replay buffer.
11:     **end for**
12:     /* Train global critic*/
13:     Sample minibatch of size S $\left(\mathbf{s}^j, \mathbf{a}^j, \mathbf{r}_g^j, \mathbf{s}'^j\right)$ from buffer.
14:     $\left(a_1', \ldots, a_N'\right) := \left(\pi_{\theta_i}'(o_i'^j), \ldots, \pi_{\theta_N}'(o_N'^j)\right)$
15:     Set $y_g^j = r_g^j + \gamma \min_{i=1,2} Q_{\psi_i'}^{g_i}\left(s'^j, a_1', \ldots, a_N'\right)$
16:     Update global critics by minimizing

$$\frac{1}{S}\sum_j\left(y_g^j - Q_{\psi_i}^{g_i}\left(s^j, a_1^j, \ldots, a_N^j\right)\right)^2$$

17:     Update target network parameters

$$\psi_i' \leftarrow \tau\psi_i + (1-\tau)\psi_i'$$

18:     **if** episode mod $d$ **then**
19:         /* Train local critics and update agent policies*/
20:         **for** agent $i = 1$ to $N$ **do**
21:         Sample minibatch of size S $\left(\mathbf{s}^j, \mathbf{a}^j, \mathbf{r}_l^j, \mathbf{s}'^j\right)$
22:         Set $y^j = r_{i_l}^j + \gamma Q_{\phi_i'}^{\pi'}\left(o'^j, \pi_{\theta_i}'(o_i'^j)\right)$
23:         Update local critic by minimizing

$$\frac{1}{S}\sum_j\left(y^j - Q_{\phi_i}^{\pi}\left(o^j, a_i^j\right)\right)^2$$

24:

$$\theta_i = \theta_i + \frac{1}{S}\sum_j\nabla_{\theta_i}\pi_i\left(a_i|o_i^j\right)\nabla_{a_i}Q_{\psi}^{g_1}\left(s^j, a_1^j, \ldots, a_N^j\right)$$
$$+\nabla_{\theta_i}\pi_i\left(a_i|o_i^j\right)\nabla_{a_i}Q_{\phi_i}^{\pi}\left(o^j, a_i^j\right)$$

25:         **end for**
26:     Update target network parameters for each agent $i$

$$\theta_i' \leftarrow \tau\theta_i + (1-\tau)\theta_i'$$
$$\phi_i' \leftarrow \tau\phi_i + (1-\tau)\phi_i'$$

27:     **end if**
28: **end for**

---

around the VIP (payload). The objective of the defenders is to minimize the potential physical attacks as the VIP is moving in a variety of different real world scenarios and are implemented in the Multi-Agent Particle Environment [20]. An illustration of the environments can be seen in Figure 1.

(a) Random Landmarks
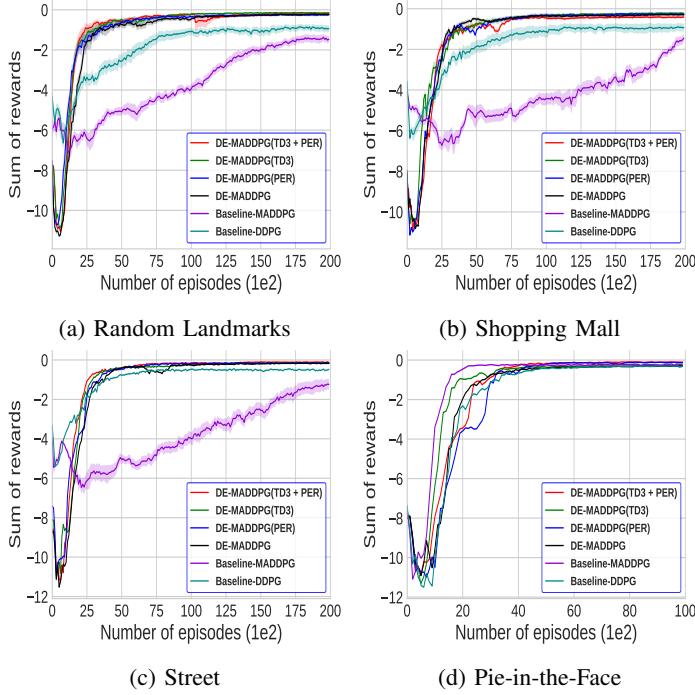
(b) Shopping Mall

(c) Street

(d) Pie-in-the-Face

Fig. 4: Learning curves of the experiments representing the sum of local rewards. Notice that similar to global rewards, the local reward learning of MADDPG and DDPG is also unstable and reaches only at sub-optimal performance.

The environment consists of the VIP (payload), defensive escort team and one or more classes of bystanders. The VIP (brown disk) starts from the starting point and moves towards the destination landmark (green disk). The goal of the defensive escort team is learn an optimal formation around the VIP to protect it from potential physical attacks. In order to maintain social norms, the defensive agents need to maintain a certain distance from the VIP. To closely simulate real world situations and providing substantial variability, four different scenarios were developed that differ in the number, arrangement of the landmarks and the behavior of the different classes of bystanders:

- **Random Landmark:** In this fully stochastic scenario, landmarks are placed randomly in the area. The starting point and destination for the VIP are randomly selected landmarks. The bystanders are performing random waypoint navigation: they choose a landmark at random, move towards it and repeat this process till the end of the episode.
- **Shopping Mall:** In this scenario, landmarks are placed at the edge of the environment emulating shops on a street or shopping mall. The bystanders are moving between randomly chosen shops.
- **Street:** This scenario aims to model a crowded sidewalk. The bystanders are moving in two different directions towards waypoints that are outside the current area. However, due to their proximity to each other, the position of

the other bystanders influence their movement described by laws of particle motion [32].

- **Pie-in-the-Face**: This scenario models a VIP walking a "red carpet", with bystanders standing behind designed lines. In this scenario, an unruly bystander breaks the line to approach the VIP (presumably, to throw a pie in his/her face).

The state of the environment is given by the locations of the landmarks, bystanders, VIP and the defensive team. To closely represent a real world bodyguard that has a limited range of perception, the observation of each agent is the relative physical state of the nearest $M$ bystanders, the VIP and the remaining members of the defensive escort team and represented as $o_i = [x_{j,...N+M}] \in \mathcal{O}_i$ where $x_j$ is the observation of the entity $j$ from the perspective of agent $i$. In our experiments, we used $M = 5$.

We chose the defensive escort team problem with these environments, because it is a representative case of a well structured dual-reward collaborative MARL problem, the environments are continuous with a relatively high dimensional state space. Furthermore, previous work using these environments provide strong baselines against which the proposed algorithms can be compared. Aligning our experimental setup with [3], [16], we chose *4* bodyguards and *10* bystanders.

*Decomposing the Reward Function*

In this section, we review the entangled multi-objective reward function defined in [13], [3], explain the problems with it and decompose it to be used by DE-MADDPG.

As mentioned in the the previous section that the goal of the defensive escort team is to learn an optimal formation around the VIP to minimize the physical threat while simultaneously maintain a certain distance from the VIP to follow the social norms. To achieve both of these objectives, the multi-objective reward function for each agent $i$ is defined as:

$$r_{total} = \alpha \overbrace{\left( -1 + \prod_{k=1}^{M} \left( 1 - RT\left( VIP, b_k, R \right) \right) \right)}^{r_{global}} + (1-\alpha) \underbrace{\left( \mathcal{D}\left( VIP, x_i \right) \right)}_{r_{local}} \quad (3)$$

where $r_{global}$ represents the reward that each agent receives given the formation of the team around the VIP at time-step $t$, therefore, represents the main objective of the team and $r_{local}$ represents the reward that the agent receives for maintaining a certain distance from the VIP and is defined as

$$\mathcal{D}\left( VIP, x_i \right) = \begin{cases} 0 & m \leq \|x_i - VIP\|_2 \leq d \\ -1 & \text{otherwise} \end{cases} \quad (4)$$

where $m$ is the minimum distance the agent has to maintain from VIP and $d$ is the safe distance. As per our understanding, $r_{total}$ has two problems. The first problem is the stability

issue since the policy oscillates between optimizing the global reward and the local reward. This stability problem exacerbates when either of the rewards are sparse and the other reward is dense. This phenomenon is also explained in Section V-C.

The second problem is the $\alpha$ hyper-parameter. The $\alpha$ hyper-parameter assigns weight to both rewards. Given the various difficulty settings of the simulations, security should be prioritized in some simulations as compared to the others. Moreover, as reinforcement learning experiments take substantial amount of time, finding an optimal $\alpha$ can be time consuming.

We solve this problem by having a dual critic architecture where the task of the global centralized critic is to approximate the cumulative global reward while each agent's local critic approximates its own local reward. In this particular scenario, $Q_\psi^g$ learns to approximate $r_{global}$ while $Q_{\theta_i}^\pi$ approximates $r_{local}$. The benefit of this decomposition is more stable learning and exclusion of the $\alpha$ hyper-parameter.

### B. Evaluations

To evaluate the efficacy of DE-MADDPG and its variants we compare our results with baseline MADDPG and DDPG on defensive escort team problem in the environments described above. We trained the global critic with two different approaches. In the first approach, we updated the policy by using the standard Multi-Agent Deep Deterministic Policy Gradient mentioned in Equation (1) while in the second approach, we updated the policy using the Twin Delayed Deep Deterministic Policy Gradient mentioned in Equation (2). Additionally, to mitigate the global sparse reward problem, we replace the standard replay buffer with prioritized experience replay buffer (PER). We performed our experiments on 8 different seeds. We used three layered neural networks for both the critic and actor networks. For each environment, we trained each approach for $20,000$ episodes except the *Pie-in-the-Face* environment which was trained for $10,000$ episodes.

Figure 3 shows the learning curves of our experiments. Figure 3a corresponds to the *Random Landmark* environment and it can be seen that DE-MADDPG based approaches outperforms the MADDPG and DDPG by a significant margin. Similar observation can be seen for the other two environments i.e., *Shopping Mall*, and *Street*. Finally, for the last environment *Pie-in-the-Face*, there is little difference between the performance of the different approaches. A possible reason for this is that this environment, focusing on a single attacking bystander, makes the positioning choice less complex.

Though, fig. 3 shows the learning curves that visually represents the performance of the different approaches, it does not quantitatively explain the improvement in performance across different approaches. To that end, we test our trained policies across all environments on 8 different seeds. Table I shows the average returns of the global reward over $1000$ episodes. We notice that in complex environments such as *Shopping Mall* and *Random Landmarks*, DE-MADDPG augmented with TD3 and PER were able to achieve 55% and 73% better performance than baseline MADDPG. Similarly, on the

slightly less complex environments *Street* and *Pie-in-the-Face*, the performance was about 61% and 100% better.

TABLE I: Average cumulative return of the global reward over 1000 episodes over 8 seeds. Maximum value for each task is bolded. ± corresponds to 95% confidence interval over seeds.

| Environment | DE-MADDPG (TD3+PER) | DE-MADDPG (TD3) | DE-MADDPG (PER) | DE-MADDPG | MADDPG | DDPG |
|---|---|---|---|---|---|---|
| Shopping Mall | **-4.87 ± 0.06** | -5.61 ± 0.08 | -6.17 ± 0.08 | -6.27 ± 0.07 | -7.59 ± 0.13 | -9.51 ± 0.12 |
| Rand. Landm. | **-4.43 ± 0.05** | -4.91 ± 0.06 | -5.75 ± 0.04 | -6.34 ± 0.08 | -7.67 ± 0.13 | -10.22 ± 0.16 |
| Street | **-2.13 ± 0.06** | -2.38 ± 0.07 | -2.36 ± 0.07 | -2.66 ± 0.08 | -3.43 ± 0.13 | -3.81 ± 0.11 |
| Pie-in-the-Face | -0.07 ± 0.002 | **-0.06 ± 0.002** | -0.11 ± 0.003 | -0.07 ± 0.004 | -0.12 ± 0.003 | -0.10 ± 0.006 |

Figure 4 and Table II shows the learning curves and the test results of the sum of the local rewards. Similar to the global reward, it can be seen in Figure 4 and and Table II that DE-MADDPG based approaches outperforms MADDPG and DDPG results. One point to be noted here is that maintaining a certain distance from a moving payload is fairly an easy problem for standard reinforcement learning algorithms as the reward is dense and the state space is relatively easy but adding an additional objective such as global reward maximization not only had catastrophic affect on the global reward maximization but also negatively effected the learning of a trivial task.

TABLE II: Average cumulative return of the local reward over 1000 episodes over 8 seeds. Maximum value for each task is bolded. ± corresponds to 95% confidence interval over seeds.

| Environment | DE-MADDPG (TD3+PER) | DE-MADDPG (TD3) | DE-MADDPG (PER) | DE-MADDPG | MADDPG | DDPG |
|---|---|---|---|---|---|---|
| Shopping Mall | -0.41 ± 0.02 | **-0.23 ± 0.03** | -0.28 ± 0.03 | -0.30 ± 0.07 | -1.44 ± 0.07 | -0.92 ± 0.05 |
| Rand. Landm. | **-0.15 ± 0.02** | -0.19 ± 0.01 | -0.25 ± 0.02 | -0.22 ± 0.03 | -1.46 ± 0.06 | -0.94 ± 0.05 |
| Street | **-0.12 ± 0.01** | -0.15 ± 0.02 | -0.16 ± 0.02 | -0.18 ± 0.08 | -1.23 ± 0.08 | -0.47 ± 0.04 |
| Pie-in-the-Face | **-0.11 ± 0.01** | -0.28 ± 0.01 | -0.12 ± 0.01 | -0.31 ± 0.01 | -0.20 ± 0.01 | -0.31 ± 0.01 |

A common pattern can be seen in Figure 3 and Figure 4 that DE-MADDPG based approaches not only achieve higher global and local rewards, they achieve it significantly faster as compared to MADDPG and DDPG. For example, in complex environments such as *Random Landmarks* and *Shopping Mall*, DE-MADDPG (TD3 + PER) reaches the maximum performance before 2500 episodes while the baseline MADDPG reaches its best performance for *Random Landmarks* and *Shopping Mall* environment at around $12,000$ episodes and $20,000$ episodes respectively.

### C. Stability

One of the focal point of this study was to find a stable solution that solves the dual-reward MARL problem. In this section, we analyse the stability of DE-MADDPG based solutions. Though, Figure 3 shows that DE-MADDPG when augmented with TD3 outperforms MADDPG in maximizing global reward, the average of results across all the seeds make it difficult to analyze the stability. For that purpose, we chose fixed single seed runs of all the solutions on the *Shopping Mall* environment and plotted their learning curves. It can be seen in Figure 5 that all DE-MADDPG based solutions are stable in learning to maximize global reward as compared to baseline MADDPG. This behavior is common across all seeds and environments except *Pie-in-the-Face* as it is does not provide any complexity.
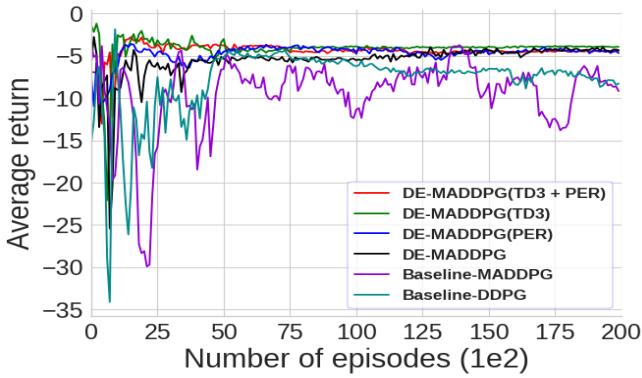
Fig. 5: Visualizing the single seed learning graph of shopping mall environment. Notice that DE-MADDPG(TD3) almost becomes flat after 10,000 episodes.
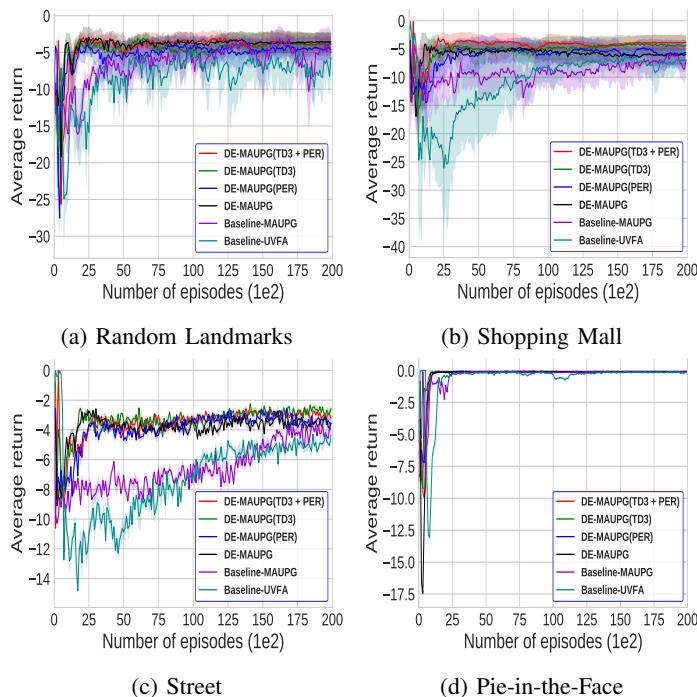


(a) Random Landmarks

(b) Shopping Mall

(c) Street

(d) Pie-in-the-Face

Fig. 6: Learning curves representing the average cumulative global reward of multi-scenario learning experiments..

TABLE III: Average cumulative return of the global reward of multi-scenario learning over 1000 episodes over 8 seeds. Maximum value for each task is bolded. ± corresponds to 95% confidence interval over seeds.

| Environment | DE-MAUPG (TD3+PER) | DE-MAUPG (TD3) | DE-MAUPG (PER) | DE-MAUPG | MAUPG | UVFA |
|---|---|---|---|---|---|---|
| Shopping Mall | **-3.85 ± 0.07** | -4.34 ± 0.08 | -5.70 ± 0.09 | -5.98 ± 0.10 | -6.97± 0.11 | -6.94± 0.11 |
| Rand. Landm. | -3.91 ± 0.07 | -3.86 ± 0.07 | -4.65 ± 0.08 | **-3.54 ± 0.08** | -5.02± 0.14 | -5.62± 0.11 |
| Street | **-2.54± 0.08** | -3.22 ± 0.09 | -3.58 ± 0.11 | -3.51 ± 0.10 | -3.97± 0.14 | -4.34± 0.14 |
| Pie-in-the-Face | **-0.07 ± 0.003** | -0.13± 0.008 | -0.08 ± 0.003 | **-0.07 ± 0.002** | -0.12± 0.003 | -0.19± 0.008 |

## D. Multi-Scenario Experiments

Multi-agent reinforcement learning is sensitive to distortions and does not work well if the trained policies are deployed on scenarios other than the scenario on which the policies are trained on. This problem is generally referred as single-

TABLE IV: Average cumulative return of the local reward of multi-scenario learning over 1000 episodes over 8 seeds. Maximum value for each task is bolded. ± corresponds to 95% confidence interval over seeds.

| Environment | DE-MAUPG (TD3+PER) | DE-MAUPG (TD3) | DE-MAUPG (PER) | DE-MAUPG | MAUPG | UVFA |
|---|---|---|---|---|---|---|
| Shopping Mall | -0.56 ± 0.02 | **-0.27 ± 0.03** | -0.22 ± 0.03 | -0.25 ± 0.07 | -1.04± 0.06 | -0.96± 0.05 |
| Rand. Landm. | -0.18 ± 0.02 | -0.12 ± 0.02 | -0.22 ± 0.02 | **-0.11 ± 0.02** | -0.89± 0.08 | -0.71± 0.08 |
| Street | -0.16 ± 0.02 | -0.15 ± 0.02 | -0.17 ± 0.02 | **-0.13 ± 0.02** | -0.64± 0.07 | -0.60± 0.07 |
| Pie-in-the-Face | **-0.21 ± 0.01** | -0.43 ± 0.01 | -0.24 ± 0.01 | -0.27 ± 0.01 | -0.27± 0.007 | -0.22± 0.009 |

task multi-scenario learning. The goal here is to learn a joint policy $\pi$ that performs equally well as scenario-dependant policy. In [16] have introduced *Multi-Agent Universal Policy Gradients* (MAUPG) to solve the multi-scenario learning and evaluated it on the VIP protection environments similar to our experiments. The main idea behind MAUPG is to replace the standard centralized Q-functions in MADDPG with Universal Value Function Approximators (UVFA). Given the similarity between MADDPG and MAUPG, we replaced all the critics in DE-MADDPG with UVFAs. For brevity, we will refer our multi-scenario solution as *Decomposed-Multi-Agent Universal Policy Gradients* (DE-MAUPG) For our experiments, we kept our settings identical to DE-MADDPG experiments and trained it for 20,000 episodes. The point to note here is that unlike scenario-dependant training, where every scenario was trained for 20,000 episodes, all scenarios are trained in parallel using one joint policy $\pi$.

Figure 6 shows the learning curves of the DE-MAUPG and its variants. Similar to our results from Section V-B, decomposition based learning solutions outperform the baseline MAUPG. The point to be noted here is that not our solutions learn to achieve higher reward but they also learn faster. This can be easily seen in Figure 6a and Figure 6b where DE-MAUPG (TD3 + PER) reaches the maximum performance in less than 2500 episodes where baselines MAUPG reaches its peak performance at 12,500 episodes in *Random Landmarks* environment and does not even reach its peak performance before 19,000 episodes for the *Shopping Mall* environment.

Table III quantifies the improvement of DE-MAUPG in maximzing global reward when compared to MAUPG and UVFA. We find that the decomposition based approaches achieve 81% on the *Shopping Mall* environment and 41% on the *Random Landmarks* environment. Similarly Table IV quantifies the improvement of DE-MAUPG in maximzing local reward when compared to MAUPG and UVFA.

## E. Computational Evaluations

In this section we evaluate the growth of parametric space as the number of agents increases. Moreover, we empirically evaluate the computational time for DE-MADDPG and compare it with MADDPG and DDPG. The main component of the MADDPG that fuels its learning are the distributed centralized Q-functions. As the number of the agents grow, the input space of those Q-functions increase quadratically. Concretely, assuming all the agents have identical observation and action space, the number of trainable parameters can be represented by $\mathcal{O}(n^2(odim + adim))$. Where $n$ is the number of agents,

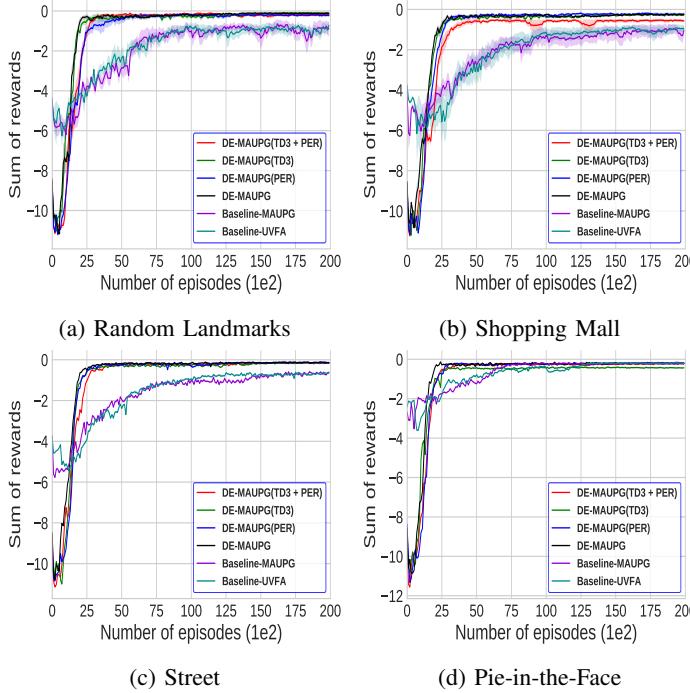(a) Random Landmarks      (b) Shopping Mall

(c) Street      (d) Pie-in-the-Face

Fig. 7: Learning curves representing the sum of local rewards of multi-scenario learning experiments.
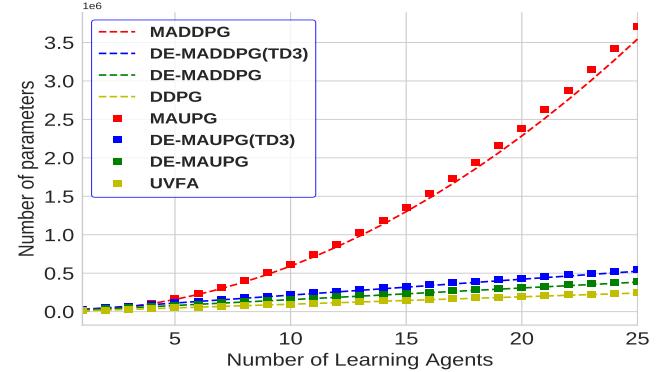


Fig. 8: Number of trainable parameters in the main Q-networks. Notice that 25 agents can be trained by DE-MADDPG approaches using the same number of parameters as compared to 10 agents if MADDPG is used.

TABLE V: Average training time for experiments in minutes over 8 different seeds. Minimum value for each task is bolded. $\pm$ corresponds to standard deviation across seeds

| Environment | DE-MADDPG (TD3+PER) | DE-MADDPG (TD3) | DE-MADDPG (PER) | DE-MADDPG | MADDPG |
|---|---|---|---|---|---|
| Shopping Mall | 168.25 ± 7.22 | 152.85 ± 6.31 | 167.25 ± 4.40 | **149.5 ± 1.65** | 268.5± 21.21 |
| Rand. Landm. | 168.87 ± 0.78 | 151.25 ± 1.19 | 164.14 ± 1.35 | **146.875 ± 1.16** | 262.25± 2.86 |
| Street | 397.62 ± 4.71 | 352.87 ± 5.63 | 356.87 ± 13.97 | **340.75 ± 16.74** | 462.75± 43.14 |
| Pie-in-the-Face | 64.25 ± 1.23 | 59.25 ± 0.59 | 63.5± 1.09 | **55.12 ± 1.05** | 72.79± 1.16 |

TABLE VI: The parameters used for various variations of DE-MADDPG and the baseline algorithm MADDPG in the experiments.

| Parameter | DE-MADDPG (TD3+PER) | DE-MADDPG (TD3) | DE-MADDPG (PER) | DE-MADDPG | MADDPG | DDPG |
|---|---|---|---|---|---|---|
| Episodes | 20k | 20k | 20k | 20k | 20k | 20k |
| Replay buffer | $10^6$ | $10^6$ | $10^6$ | $10^6$ | $10^6$ | $10^6$ |
| Minibatch size | 2048 | 2048 | 2048 | 2048 | 2048 | 2048 |
| Steps per train $Q_g$ | 4 | 4 | 4 | 4 | N/A | N/A |
| Steps per train $Q_l$ | 2 | 2 | 2 | 2 | 2 | 2 |
| Max env steps | 25 | 25 | 25 | 25 | 25 | 25 |
| PER $\alpha$ | 0.6 | N/A | 0.6 | N/A | N/A | N/A |
| PER $\beta$ | 0.4 | N/A | 0.4 | N/A | N/A | N/A |
| PER $\epsilon$ | 1e-6 | N/A | 1e-6 | N/A | N/A | N/A |
| PER $\beta$ decay | 10000 | N/A | 10000 | N/A | N/A | N/A |

*odim* and *adim* represents the dimensionality of observation and action space respectively. Alternatively, DE-MADDPG solves this scalability problem by having a shared global centralized Q-function whose parametric space increases linearly and can be represented as $\mathcal{O}(n(odim+adim))$. In Figure 8, we show the growth of number of parameters of all the main Q-networks of MADDPG, DE-MADDPG and its TD3 variant. It can be seen that our solution takes significantly small number of parameters as MADDPG and MAUPG. Note that the figure only represents the number of trainable parameters in the main Q-networks and does not include target or policy networks as the growth of the policy network parameters are identical and no gradient based learning happens in the target networks. All these statements hold true for their UVFA variants.

We empirically verified the benefits of parametric reduction by measuring the time taken to train the experiments. It can be seen in Table V that DE-MADDPG based approaches always train faster than baseline MADDPG. The complete details about the network architecture can be seen in Section VI. We do not provide time comparisons for DDPG and UVFA variants as the experiments were run on multiple machines with a variety of hardware thus making it difficult for a fair comparison.

## VI. EXPERIMENTAL DETAILS

### A. Network Architecture

Both actor and critic networks for all agents consists of 2 hidden layers containing 64 units in each layer. The hidden layers uses ReLU activation function while the output layers of both networks use linear activation function. Both networks are initialized using Xavier normal initializers however, the output layer of the target critics were initialized with uniform random values between (-0.01 and 0.01) to enable one-step look ahead learning of the critics after each training cycle. To keep experiments as fair as possible, we initialized the network parameters with same seed across all the experiments. Complete hyperparameter details can be seen in Table VI. Note that same configurations were used for multi-scenario learning experiments.

## VII. CONCLUSIONS

In this paper, we focused on the dual-reward MARL: a collaborative setting where a group of learning agents have to simultaneously learn to maximize the collective global reward and individual local reward. To solve the problem, we proposed the *Decomposed Multi-Agent Deep Deterministic Policy Gradient (DE-MADDPG)* algorithm and applied it to the problem of defensive escort team: how can agent learn a policy

to maintain an optimal formation around the VIP to protect him/her from possible physical attack. We first demonstrated that decomposing a multi-objective reward function leads to higher and more stable performance. We compared our results with the MADDPG algorithm and achieved at least 50% better performance in terms of the collected reward. Additionally, we showed that our solution is computationally efficient and requires a significantly lower number of parameters while achieving better performance than the baseline. Finally, we showed that by replacing the standard critics with UVFAs, our solution also outperforms MAUPG which is a baseline algorithm for single-task multi-scenario learning in multi-agent reinforcement learning.

REFERENCES

[1] J. Z. Leibo, V. Zambaldi, M. Lanctot, J. Marecki, and T. Graepel, "Multi-agent reinforcement learning in sequential social dilemmas," in *Proc. of the 16th Int'l Conf. on Autonomous Agents and Multiagent Systems (AAMAS-2017)*, pp. 464–473, 2017.

[2] Y. Cao, W. Yu, W. Ren, and G. Chen, "An overview of recent progress in the study of distributed multi-agent coordination," *IEEE Transactions on Industrial informatics*, vol. 9, no. 1, pp. 427–438, 2012.

[3] H. U. Sheikh, M. Razghandi, and L. Bölöni, "Learning distributed cooperative policies for security games via deep reinforcement learning," in *2019 IEEE 43rd Annual Computer Software and Applications Conference (COMPSAC-2019)*, vol. 1, pp. 489–494, Jul 2019.

[4] OpenAI, *OpenAI Five*, 2018. https://blog.openai.com/openai-five/.

[5] M. Jaderberg, W. M. Czarnecki, I. Dunning, L. Marris, G. Lever, A. G. Castaneda, C. Beattie, N. C. Rabinowitz, A. S. Morcos, A. Ruderman, N. Sonnerat, T. Green, L. Deason, J. Z. Leibo, D. Silver, D. Hassabis, K. Kavukcuoglu, and T. Graepel, "Human-level performance in first-person multiplayer games with population-based deep reinforcement learning," *arXiv preprint arXiv: 1807.01281*, 2018.

[6] O. Vinyals, I. Babuschkin, W. M. Czarnecki, M. Mathieu, A. Dudzik, J. Chung, D. H. Choi, R. Powell, T. Ewalds, P. Georgiev, J. Oh, D. Horgan, M. Kroiss, I. Danihelka, A. Huang, L. Sifre, T. Cai, J. P. Agapiou, M. Jaderberg, A. S. Vezhnevets, R. Leblond, T. Pohlen, V. Dalibard, D. Budden, Y. Sulsky, J. Molloy, T. L. Paine, C. Gulcehre, Z. Wang, T. Pfaff, Y. Wu, R. Ring, D. Yogatama, D. Wünsch, K. McKinney, O. Smith, T. Schaul, T. Lillicrap, K. Kavukcuoglu, D. Hassabis, C. Apps, and D. Silver, "Grandmaster level in starcraft ii using multi-agent reinforcement learning," *Nature*, 2019.

[7] B. Baker, I. Kanitscheider, T. Markov, Y. Wu, G. Powell, B. McGrew, and I. Mordatch, "Emergent tool use from multi-agent autocurricula," *arXiv preprint arXiv:1909.07528*, 2018.

[8] J. Foerster, G. Farquhar, T. Afouras, N. Nardelli, and S. Whiteson, "Counterfactual multi-agent policy gradients," in *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence (AAAI-2018)*, 2018.

[9] T. Rashid, M. Samvelyan, C. Schroeder, G. Farquhar, J. Foerster, and S. Whiteson, "QMIX: Monotonic value function factorisation for deep multi-agent reinforcement learning," in *Proceedings of the 35th International Conference on Machine Learning*, pp. 4295–4304, 2018.

[10] R. Lowe, Y. Wu, A. Tamar, J. Harb, P. Abbeel, and I. Mordatch, "Multi-agent actor-critic for mixed cooperative-competitive environments," in *Advances in Neural Information Processing Systems 30*, pp. 6379–6390, 2017.

[11] S. Li, Y. Wu, X. Cui, H. Dong, F. Fang, and S. Russell, "Robust multi-agent reinforcement learning via minimax deep deterministic policy gradient," in *Proc. of the Thirty-Third AAAI Conference on Artificial Intelligence (AAAI-2019)*, 2019.

[12] J. Yang, A. Nakhaei, D. Isele, H. Zha, and K. Fujimura, "CM3: cooperative multi-goal multi-stage multi-agent reinforcement learning," *arXiv preprint arXiv:1809.05188*, 2018.

[13] H. U. Sheikh and L. Bölöni, "Designing a multi-objective reward function for creating teams of robotic bodyguards using deep reinforcement learning," in *Proc. of 1st Workshop on Goal Specifications for Reinforcement Learning (GoalsRL-2018) at ICML 2018*, July 2018.

[14] T. Schaul, J. Quan, I. Antonoglou, and D. Silver, "Prioritized experience replay," in *International Conference on Learning Representations (ICLR-2016)*, 2016.

[15] S. Fujimoto, H. van Hoof, and D. Meger, "Addressing function approximation error in actor-critic methods," in *Proc. of the 35th International Conference on Machine Learning (ICML-2018)*, pp. 1587–1596, 2018.

[16] H. U. Sheikh and L. Bölöni, "Emergence of scenario-appropriate collaborative behaviors for teams of robotic bodyguards," in *Proc. of the 18th International Conference on Autonomous Agents and Multi-Agent Systems (AAMAS-2019)*, pp. 2189–2191, 2019.

[17] M. Tan, "Multi-agent reinforcement learning: Independent vs. cooperative agents," in *Proceedings of the tenth international conference on machine learning*, pp. 330–337, 1993.

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

[19] J. Hu and M. P. Wellman, "Nash Q-learning for general-sum stochastic games," *Journal of Machine Learning Research*, vol. 4, pp. 1039–1069, Nov 2003.

[20] I. Mordatch and P. Abbeel, "Emergence of grounded compositional language in multi-agent populations," in *Proc. of AAAI International Conference on Artificial Intelligence (AAAI-2017)*, 2017.

[21] L. Panait and S. Luke, "Cooperative multi-agent learning: The state of the art," *Autonomous Agents and Multi-Agent Systems*, vol. 11, no. 3, pp. 387–434, 2005.

[22] J. L. Austerweil, S. Brawner, A. Greenwald, E. Hilliard, M. Ho, M. L. Littman, J. MacGlashan, and C. Trimbach, "How other-regarding preferences can promote cooperation in non-zero-sum grid games," in *Proceedings of the AAAI Symposium on Challenges and Opportunities in Multiagent Learning for the Real World*, 2016.

[23] K. Zhang, Z. Yang, H. Liu, T. Zhang, and T. Basar, "Fully decentralized multi-agent reinforcement learning with networked agents," in *Proc. of the 35th International Conference on Machine Learning (ICML-2018)*, pp. 5872–5881, 2018.

[24] S. Omidshafiei, J. Pazis, C. Amato, J. P. How, and J. Vian, "Deep decentralized multi-task multi-agent reinforcement learning under partial observability," in *International Conference on Machine Learning*, pp. 2681–2690, 2017.

[25] A. S. Vezhnevets, S. Osindero, T. Schaul, N. Heess, M. Jaderberg, D. Silver, and K. Kavukcuoglu, "Feudal networks for hierarchical reinforcement learning," in *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pp. 3540–3549, JMLR. org, 2017.

[26] C. Wu, A. Rajeswaran, Y. Duan, V. Kumar, A. M. Bayen, S. Kakade, I. Mordatch, and P. Abbeel, "Variance reduction for policy gradient with action-dependent factorized baselines," in *Proc. of the 6th Int'l Conf. on Learning Representations (ICLR)*, 2018.

[27] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *arXiv preprint arXiv:1707.06347*, 2017.

[28] D. Silver, G. Lever, N. Heess, T. Degris, D. Wierstra, and M. Riedmiller, "Deterministic policy gradient algorithms," in *Proc. of the 31st Int'l Conf. on Machine Learning (ICML-2014)*, pp. 387–395, 2014.

[29] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra, "Continuous control with deep reinforcement learning," in *Proc. of the 3rd Int'l Conf. on Learning Representations (ICLR-2015)*, 2015.

[30] H. Van Hasselt, A. Guez, and D. Silver, "Deep reinforcement learning with double Q-learning," in *Proc of the Thirtieth AAAI Conference on Artificial Intelligence (AAAI-2016)*, 2016.

[31] J. Ackermann, V. Gabler, T. Osa, Alec, and M. Sugiyama, "Reducing overestimation bias in multi-agent domains using double centralized critics," *arXiv preprint arXiv:1910.01465*, 2019.

[32] T. Vicsek, A. Czirók, E. Ben-Jacob, I. Cohen, and O. Shochet, "Novel type of phase transition in a system of self-driven particles," *Physical Review Letters*, vol. 75, no. 6, p. 1226, 1995.