# Mis-spoke or mis-lead: Achieving Robustness in Multi-Agent Communicative Reinforcement Learning

Wanqi Xue
Nanyang Technological University
Singapore
wanqi001@e.ntu.edu.sg

Wei Qiu*
Nanyang Technological University
Singapore
qiuw0008@e.ntu.edu.sg

Bo An
Nanyang Technological University
Singapore
boan@ntu.edu.sg

Zinovi Rabinovich
Nanyang Technological University
Singapore
zinovi@ntu.edu.sg

Svetlana Obraztsova
Nanyang Technological University
Singapore
lana@ntu.edu.sg

Chai Kiat Yeo
Nanyang Technological University
Singapore
asckyeo@ntu.edu.sg

## ABSTRACT

Recent studies in multi-agent communicative reinforcement learning (MACRL) have demonstrated that multi-agent coordination can be greatly improved by allowing communication between agents. Meanwhile, adversarial machine learning (ML) has shown that ML models are vulnerable to attacks. Despite the increasing concern about the robustness of ML algorithms, how to achieve robust communication in multi-agent reinforcement learning has been largely neglected. In this paper, we systematically explore the problem of adversarial communication in MACRL. Our main contributions are threefold. First, we propose an effective method to perform attacks in MACRL, by learning a model to generate optimal malicious messages. Second, we develop a defence method based on message reconstruction, to maintain multi-agent coordination under message attacks. Third, we formulate the adversarial communication problem as a two-player zero-sum game and propose a game-theoretical method $\mathfrak{R}$-MACRL to improve the worst-case defending performance. Empirical results demonstrate that many state-of-the-art MACRL methods are vulnerable to message attacks, and our method can significantly improve their robustness.

## KEYWORDS

Multi-Agent Reinforcement Learning; Robust Reinforcement Learning; Adversarial Reinforcement Learning

## 1 INTRODUCTION

Cooperative multi-agent reinforcement learning (MARL) has achieved remarkable success in a variety of challenging problems, such as urban systems [31], coordination of robot swarms [11] and real-time strategy video games [38]. To tackle the problems of scalability and non-stationarity in MARL, the framework of centralized training

with decentralized execution (CTDE) is proposed [18, 25], where decentralized policies are learned in a centralized manner so that they can share experiences, parameters, etc., during training. Despite the advantages, CTDE-based methods still perform unsatisfactorily in scenarios where multi-agent coordination is necessary, mainly due to partial observability in decentralized execution. To mitigate partial observability, many multi-agent communicative reinforcement learning (MACRL) methods have been proposed, which allow agents to exchange information such as private observations, intentions during the execution phase. MACRL greatly improves multi-agent coordination in a wide range of tasks [6, 7, 13, 15, 32, 39, 40].

Meanwhile, adversarial machine learning has received extensive attention [2, 10]. Adversarial machine learning demonstrates that machine learning (ML) models are vulnerable to manipulation [8, 9]. As a result, ML models often suffer from performance degradation when under attack. For the same reason, many practical applications of ML models are at high risk. For instance, researchers have shown that, by placing a few small stickers on the ground at an intersection, self-driving cars can be tricked into making abnormal judgements and driving into the opposite lane [34]. Maliciously designed attacks on ML models can have serious consequences.

Unfortunately, despite great importance, adversarial problems, especially adversarial inter-agent communication problems, remain largely uninvestigated in MACRL. Blumenkamp et al. show that, by introducing random noise in communication, agents are able to deceive their opponents in competitive games [3]. However, the attacks are not artificially designed and therefore inefficient. Besides, cooperative cases, where communication is more crucial, are neglected. They also fail to propose an effective defence, but merely retrain the models to adapt to the attacks. Mitchell et al. propose to generate weights of Attention models [37] through Gaussian process for defending against random attacks in attention-based MACRL [22]. However, the applicability of this approach is unsatisfactory, being limited to attention-based MACRL, and its performance on maliciously designed attacks is unclear.

In this paper, we systematically explore the problem of adversarial communication in MACRL, where there are malicious agents that attempt to disrupt multi-agent cooperation by manipulating messages. Our contributions are in three aspects. First, we propose an effective learning approach to model the optimal attacking scheme in MACRL. Second, to defend against the adversary, we propose a two-stage message filter which works by first detecting

the malicious message and then recovering the message. The defending scheme can greatly maintain the coordination of agents under message attacks. Third, to address the problem that the message filter can be exploited by learnable attackers, we formulate the attack-defense problem as a two-player zero-sum game and propose $\mathfrak{R}$-MACRL, based on a game-theoretical framework Policy-Space Response Oracle (PSRO) [19, 23], to approximate a Nash equilibrium policy in the adversarial communication game. $\mathfrak{R}$-MACRL improves the defensive performance under the worst case and thus improves the robustness. Empirical experiments demonstrate that many state-of-the-art MACRL algorithms are vulnerable to attacks and our method can significantly improve their robustness.

## 2 PRELIMINARIES AND RELATED WORK

**Multi-Agent Communicative Reinforcement Learning.** There has been extensive research on encouraging communication between agents to improve performance on cooperative or competitive tasks. Among the recent advances, some design communication mechanisms to address the problem of when to communicate [13, 14, 30]; other lines of works, e.g., TarMac [6], focus on who to communicate. These works determine the two fundamental elements in communication, i.e., the message sender and the receiver. Apart from the two elements, the message itself is another element which is crucial in communication, i.e., what to communicate: Jaques et al. propose to maximize the social influence of messages [12]. Kim et al. encode messages such that they contain the intention of agents [15]. Some other works learn to send succinct messages to meet the limitations of communication bandwidth [39, 40, 44]. Despite significant progress in MACRL, if some agents are adversarial and send maliciously designed messages, multi-agent coordination will rapidly disintegrate as these messages propagate.

**Adversarial Training.** Adversarial training is a prevalent paradigm for training robust models to defend against potential attacks [9, 33]. Recent literature has considered two types of attacks [5, 26, 35]: black-box attack and white-box attack. In black-box attack, the attacker does not have access to information about the attacked deep neural network (DNN) model; whereas in white-box attack, the attacker has complete knowledge, e.g., the architecture, the parameters and potential defense mechanisms, about the DNN model. We consider the black-box attack in our problem formulation, because the setting of the white-box attack is too idealistic and may not be applicable to many realistic adversarial scenarios. In adversarial training, the attacker tries to attack a DNN by corrupting the input via $\ell_p$-norm ($p \in \{1, 2, \infty\}$) attack [9]. The attacker carefully generates artificial perturbations to manipulate the input of the model. In doing so, the DNN will be fooled into making incorrect predictions or decisions. The attacker finds the optimal perturbation $\delta$ by optimizing:

$$\max_{\delta} \mathcal{L}_{\text{predict}}\left(f(x), f(x + \delta)\right) \quad \text{s.t.} \quad \min \mathcal{L}_{\text{norm}}\left(x, x + \delta\right)$$

where $x$ is the input, $f$ is the DNN model, $\mathcal{L}_{\text{predict}}$ is a metric to measure the distance between the outputs of the DNN model w/ and w/o being attacked, $\mathcal{L}_{\text{norm}}$ is used to measure that for the inputs.

**Adversarial Reinforcement Learning (RL).** Recent advances in adversarial machine learning motivate researchers to investigate the adversarial problem in RL [8, 20, 41]. SA-MDP [43] characterizes the problem of decision making under adversarial attacks on state observations. Lin et al. propose two tactics of attacks, i.e., the strategically-timed attack and the enchanting attack, which attack by injecting noise to states and luring the agent to a designated target [20]. Gleave et al. consider the problem of taking adversarial actions that change the environment and consequentially change the observation of agents [8]. ATLA [42] propose to train the optimal adversary to perturb state observations and improve the worst-case agent reward. The settings of these works are different from ours: we consider the multi-agent scenario and restrict the attacking approach to adversarial messages, which makes the detection of anomalies difficult. Tu et al. propose to attack on multi-agent communication [36]. However, their focus is on the representation-level, whereas we focus on the policy-level. Recently, there are some works considering a similar setting as ours [3, 22]. However, they either focus on random attacks in specific competitive games or the defence of specific communication methods.

## 3 ACHIEVING ROBUSTNESS IN MACRL

In this section, we propose our method for achieving robustness in MACRL. We begin by proposing the problem formulation for adversarial communication in MACRL. Then, we introduce the method to build the optimal attacking scheme in MACRL. Next, we propose a learnable two-stage message filter to defend against possible attacks. Finally, we propose to formulate the attack-defense problem as a two-player zero-sum game, and design a game-theoretic method $\mathfrak{R}$-MACRL to approximate a Nash equilibrium policy for the defender. In this way, we can improve the worst-case performance of the defender and thus enhance robustness.

### 3.1 Problem Formulation: Adversarial Communication in MACRL

An MACRL problem can be modeled by *Decentralised Partially Observable Markov Decision Process with Communication* (*Dec-POMDP-Com*) [24], which is defined by a tuple $\langle \mathcal{S}, \mathcal{M}, \mathcal{A}, \mathcal{P}, R, \Upsilon, O, C, \mathcal{N}, \gamma \rangle$. $\mathcal{S}$ denotes the state of the environment and $\mathcal{M}$ is the message space. Each agent $i \in \mathcal{N} := \{1, ..., N\}$ chooses an action $a_i \in \mathcal{A}$ at a state $s \in \mathcal{S}$, giving rise to a joint action vector, $\boldsymbol{a} := [a_i]_{i=1}^N \in \mathcal{A}^N$. $\mathcal{P}(s'|\boldsymbol{s}, \boldsymbol{a}) : \mathcal{S} \times \mathcal{A}^N \times \mathcal{S} \mapsto \mathcal{P}(\mathcal{S})$ is a Markovian transition function. Every agent shares the same joint reward function $R(\boldsymbol{s}, \boldsymbol{a}) : \mathcal{S} \times \mathcal{A}^N \mapsto \mathbb{R}$, and $\gamma \in [0, 1)$ is the discount factor. Due to partial observability, each agent has individual partial observation $\upsilon \in \Upsilon$, according to the observation function $O(\boldsymbol{s}, i) : \mathcal{S} \times \mathcal{N} \mapsto \Upsilon$. Each agent holds two policies, i.e., action policy $p_i(a_i|\tau_i, m_i^{in}) : \mathcal{T} \times \mathcal{M} \times \mathcal{A} \mapsto [0, 1]$ and message policy $v_i(m_i^{out}|\tau_i, m_i^{in}) : \mathcal{T} \times \mathcal{M} \times \mathcal{M} \mapsto [0, 1]$, both of which are conditioned on the action-observation history $\tau_i \in \mathcal{T} := (\Upsilon \times \mathcal{A})$ and incoming messages $m_i^{in}$ aggregated by the communication protocol $C(m_i^{in}|\boldsymbol{m}^{out}, i) : \mathcal{M}^{|\mathcal{N}|} \times \mathcal{N} \times \mathcal{M} \mapsto [0, 1]$.

In adversarial communication where there are $N_{adv}$ malicious agents, we assume that each malicious agent holds the third private adversarial policy, $\xi(\delta_i^{adv}|\tau_i, m_i^{in}, m_i^{out}) : \mathcal{T} \times \mathcal{M} \times \mathcal{M} \times \mathcal{M} \mapsto [0, 1]$, which generates adversarial message $\delta_i^{adv}$ based on its action-observation history $\tau_i$, received messages $m_i^{in}$ and the message
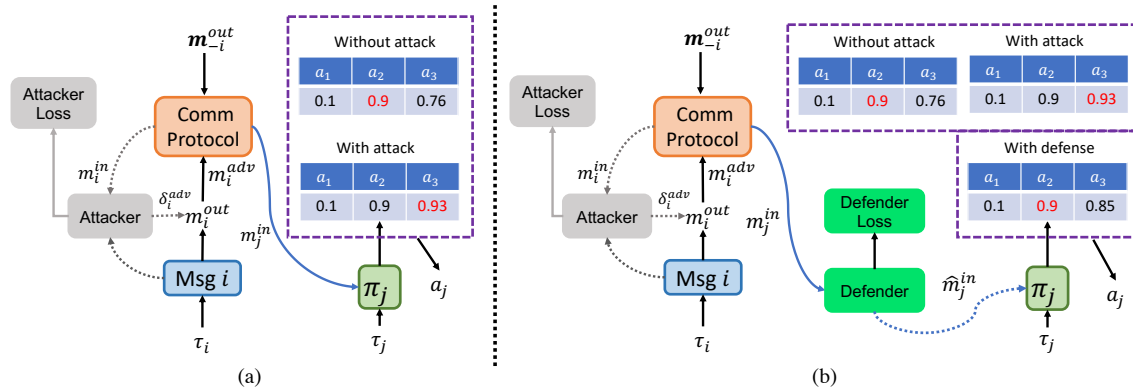
Figure 1: General framework of attack (left) and defence (right) in MACRL. *Attack:* The agent $i$, taken over by the attacker, generates malicious message $m_i^{adv}$ to disrupt multi-agent cooperation. The incoming message $m_j^{in}$ and estimated Q-values of the benign agent $j$ will be affected due to $m_i^{adv}$, which may lead to incorrect decisions. *Defence:* A learnable defender is cascaded to the communication protocol module, which is used to reconstruct the contaminated message ($m_j^{in}$ to $\hat{m}_j^{in}$). The benign agent $j$ estimates Q-values based on $\hat{m}_j^{in}$, which can reduce the probability of selecting sub-optimal actions.

$m_i^{out}$ intended to be sent. Malicious agents could send messages by convexly combining their original messages with adversarial messages, i.e., $m_i^{adv} = (1 - \omega) \times m_i^{out} + \omega \times \delta_i^{adv}$, or simply summing up the messages, i.e., $m_i^{adv} = m_i^{out} + \delta_i^{adv}$. To reduce the likelihood of being detected, apart from the adversarial policy $\xi$, malicious agents strictly follow their former action policy and message policy, trying to behave like benign agents. Fig. 1(a) presents the overall attacking procedure. The agent $i$ (attacker) is malicious and tries to generate adversarial message $m_i^{adv}$ to disrupt cooperation. The adversarial message sent by agent $i$ together with normal messages sent by other agents (denoted by $\boldsymbol{m}_{-i}^{out}$) are processed by the communication protocol (algorithm-related), generating a contaminated incoming message $m_j^{in}$ for a benign agent $j$. From agent $j$'s perspective, under such attack, the estimated Q-values will change. If the action with the highest Q-value shifts, agent $j$ will make incorrect decisions, leading to suboptimality. To perform an effective attack in MACRL, we propose to optimize the adversarial policy $\xi$ by minimizing the joint accumulated rewards, i.e., $\min_\xi \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^t r_t\right]$. We make two assumptions to make the adversarial communication problem both practical and tractable.

ASSUMPTION 1. *(Byzantine Failure [17]) Agents have imperfect information on who are malicious.*

ASSUMPTION 2. *(Concealment) Malicious agents do not communicate or cooperate with each other when performing attacks in order to be covert.*

To defend against the attacker, as in Fig. 1(b), we propose to cascade a learnable defender to the communication protocol module. The defender performs a transformation from $m_j^{in}$ to $\hat{m}_j^{in}$ to reconstruct the contaminated messages, to avoid distributing the contaminated message $m_j^{in}$ directly to agent $j$. With such transformation, the benign agent $j$ can estimate the Q-value for each action more properly and reduce the probability of selecting sub-optimal actions.

## 3.2 Learning the Attacking Scheme

To model the attack scheme in adversarial communication, we use a deep neural network (DNN) $f_\mu$, parameterized by $\theta_\mu$, to generate adversarial messages for a malicious agent. The adversarial policy $\xi$ is a multivariate Gaussian distribution whose parameters are determined by the DNN $f_\mu$, i.e., $\xi = \mathcal{N}(f_\mu(\cdot | \tau, m^{in}, m^{out}; \theta_d), \Lambda)$ where $\Lambda$ is a fixed covariance matrix. The reason of using Gaussian distribution as the prior is that it is the maximum entropy distribution under constraints of mean and variance. Each malicious agent generates adversarial messages by sampling from its adversarial policy. The optimization objective of the adversarial policy is to minimize the accumulated team rewards subject to a constraint on the distance between $m^{out}$ and $m^{adv}$. We utilize Proximal Policy Optimization (PPO) [29] to optimize the adversarial policy by maximizing the following objective:

$$\mathcal{J}_\xi(\theta_\mu) = \mathbb{E}_{(\tau, m^{in}, m^{out}, \delta^{adv})}\left[\min(\rho \cdot A^\xi, \text{clip}(\rho, 1 \pm \epsilon) \cdot A^\xi)\right]$$
$$- \alpha \cdot \mathbb{E}_{(m^{out}, \delta^{adv})}\left[(m^{out} - m^{adv})^2\right] \quad (1)$$

where $\rho = \xi(\delta^{adv} | \tau, m^{in}, m^{out}) / \xi^{old}(\delta^{adv} | \tau, m^{in}, m^{out})$, $\xi^{old}$ is the policy in the previous learning step, $\epsilon$ is the clipping ratio, and $A^\xi(\delta^{adv}, \tau, m^{in}, m^{out})$ is the advantage function. Let $e = \langle \tau, m^{in}, m^{out} \rangle$ denotes the input of the value function and we define the reward of the attacker to be negative team reward, then $A^\xi(\delta^{adv}, e) = -r + V(e') - V(e)$ where $V$ is the value function, $e'$ denotes the next step state, and $r$ is the immediate reward. We learn the value function $V(e)$ by minimizing $\mathbb{E}[(V(e) + \sum_t \gamma^t r_t)^2]$.

## 3.3 Defending against Adversarial Communication

We can train a DNN to model the attack scheme in adversarial communication. However, the defence of that is non-trivial because i) benign agents have no prior knowledge on which agents are malicious; ii) the malicious agents can inconsistently pretend to be cooperative or non-cooperative to avoid being detected; and iii)
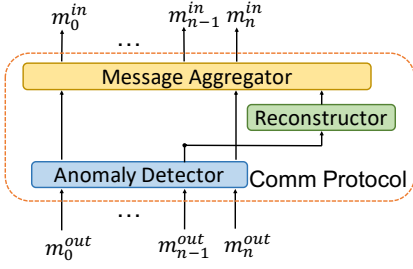
**Figure 2: Structure of the communication protocol with the message filter (defender). The anomaly detector and the reconstructor are designed to determine and recover the potential malicious messages, respectively.**

recovering useful information from the contaminated messages is difficult. To address these challenges, we design a two-stage message filter for the communication protocol to perform defence. The message filter works by first determining the messages that are likely to be malicious and then recovering the potential malicious messages before distributing them to the corresponding agents. As in Fig. 2, the message filter $\zeta(h_d, g_r)$ consists of an anomaly detector $h_d$ and a message reconstructor $g_r$. The anomaly detector, parameterized by $\theta_d$, outputs the probability that each incoming message needs to be recovered, i.e., $h_d(\cdot|m_i^{out}; \theta_d) : \mathcal{M} \mapsto \Psi$, where $\Psi$ denotes a binomial distribution, $m_i^{out}$ denotes the message sent by agent $i$. We perform sampling from the generated distributions to determine the anomaly messages $x \sim h_d(\cdot|\boldsymbol{m}^{out}; \theta_d)$ [1]. Here, $x$ is an indicator that records whether a message is predicted to be malicious or not. The predicted malicious messages $\dot{m}$ are recovered by the reconstructor $g_r(\cdot|\dot{m}_i; \theta_r) : \mathcal{M} \mapsto \mathcal{M}$ separately. The recovered message and other normal messages are aggregated and determine the messages that each agent will receive.

The optimization objective of the message filter is to maximize the joint accumulated rewards under attack, i.e., $\max_\zeta \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^t \tilde{r}_t | \xi\right]$, where $\tilde{r}_t$ is the team reward after performing the defending strategy. We utilize PPO to optimize the strategy. To mitigate the local optima induced by unguided exploration in large message space, we introduce a regularizer which is optimized under the guidance of the ground-truth labels of messages (whether they are malicious). For the reconstructor, we train it by minimizing the distance between the messages before and after being attacked. To improve the tolerance of the reconstructor to errors made by the detector, apart from malicious messages, we also use benign messages as training data, in which case the reconstructor is an identical mapping. Formally, the two-stage message filter is optimized by maximizing the following function:

$$\mathcal{J}_\zeta(\theta_d, \theta_r) = \mathbb{E}_{(\boldsymbol{m}, \boldsymbol{x})}\left[\sum_i^{|\boldsymbol{m}|} \min\left(\rho_i \cdot A_i^{h_d}, \text{clip}\left(\rho_i, 1 \pm \epsilon\right) \cdot A_i^{h_d}\right)\right]$$
$$+ \beta_1 \cdot \mathbb{E}_{\boldsymbol{m}}\left[\hat{\boldsymbol{y}} \cdot \log(h_d(\cdot|\boldsymbol{m}; \theta_d))\right]$$
$$- \beta_2 \cdot \mathbb{E}_{\boldsymbol{m}}\left[(\hat{\boldsymbol{m}} - g_r(\cdot|\boldsymbol{m}; \theta_r))^2\right] \quad (2)$$

where $\rho_i = h_d(x_i|m_i; \theta_d)/h_d^{old}(x_i|m_i; \theta_d)$, $A_i^{h_d}(x_i, m_i)$ is the advantages which are estimated similarly as in the attack, $\hat{\boldsymbol{y}}$ is the one-hot

---
[1] We abuse $h_d(\cdot|\boldsymbol{m}^{out}; \theta_d)$ to represent that messages $m_i^{out} \in \boldsymbol{m}^{out}$ are fed into $h_d$ in batch.



**Figure 3: Workflow of $\mathfrak{R}$-MACRL.**

ground-truth labels of messages, $\hat{m}$ is the messages that have not been attacked. $\beta_1$, $\beta_2$ and $\epsilon$ are hyperparameters.

### 3.4 Achieving Robust MACRL

Despite the defensive message filter, the effectiveness of the defence system can rapidly decrease if malicious agents are aware of the defender and adjust their attack schemes. To mitigate this, we formulate the attack-defence problem as a two-player zero-sum game (malicious agents are controlled by the attacker) and improve the performance of the defender in the worst case. We define the adversarial communication game by a tuple $\langle \Pi, U \rangle$, where $\Pi = \langle \Pi_\xi, \Pi_\zeta \rangle$ is the joint policy space of the players ($\Pi_\xi$ and $\Pi_\zeta$ denote the policy space of the defender and the attacker respectively), $U : \Pi \mapsto \mathbb{R}^2$ is the utility function which is used to calculate utilities for the attacker and the defender given their joint policy $\pi = \langle \pi_\xi, \pi_\zeta \rangle \in \Pi$. The utility of the defender is defined as the expected team return. The adversarial communication game is zero-sum, i.e., the utility of the defender $U_\zeta(\pi)$ must be the negative of the utility of the attacker $U_\xi(\pi)$ for $\forall \pi \in \Pi$. The solution concept of the adversarial communication game is Nash equilibrium (NE) and we can approach an NE by optimizing the following MaxMin objective:

$$\mathcal{J}_{\xi,\zeta} = \max_{\pi_\zeta \in \Pi_\zeta} U_\zeta(\text{Br}(\pi_\zeta), \pi_\zeta) = \max_{\pi_\zeta \in \Pi_\zeta} \min_{\pi_\xi \in \Pi_\xi} \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^t \tilde{r}_t | \pi_\xi, \pi_\zeta\right]$$
$$(3)$$

where $\text{Br}(\pi_\zeta)$ is the best response policy of the defender, i.e., $\text{Br}(\pi_\zeta) = \arg\min_{\pi_\xi} U_\zeta(\pi_\xi, \pi_\zeta)$.

We propose $\mathfrak{R}$-MACRL to optimize the objective in the adversarial communication game. $\mathfrak{R}$-MACRL is developed based on Policy-Space Response Oracle (PSRO) [19], a deep learning-based extension of the classic Double Oracle algorithm [21]. The workflow of $\mathfrak{R}$-MACRL is depicted in Fig. 3. At each iteration, $\mathfrak{R}$-MACRL keeps a population (set) of policies $\Pi^t \subset \Pi$, e.g., $\Pi_\zeta^t = (\pi_{11}, \pi_{12})$ and $\Pi_\xi^t = (\pi_{21}, \pi_{22})$. We can evaluate the utility table $U(\Pi^t)$ for the current iteration and calculate a Nash equilibrium $\pi_*^t = \langle \Delta(\Pi_\xi^t), \Delta(\Pi_\zeta^t) \rangle$ for the game restricted to policies in $\Pi^t$ by, for example, linear programming, replicator dynamics, etc. $\Delta$ denotes the meta-strategy which is an arbitrary categorical distribution. Next, we calculate the best response policy $\text{Br}(\pi_{*,-i}^t)$ to the NE in this restricted game for each player $i$ (player $-i$ denotes the opponent of player $i$, $i \in \{\xi, \zeta\}$), e.g., $\pi_{13} = \text{Br}(\pi_{*,\xi}^t)$ and $\pi_{23} = \text{Br}(\pi_{*,\zeta}^t)$. We extend the population by adding the best response policies to the policy set:

**Algorithm 1:** $\mathfrak{R}$-MACRL

---

1 **Input:** initial policy sets for both players $\Pi^0$, maximum number of iterations $T$;

2 Empirically estimate utilities $U(\Pi^0)$ for each joint policy $\pi \in \Pi^0$;

3 Initialize mixed strategies for both players $\pi^0_{*,i} = \text{Uniform}(\Pi^0_i)$ for $i \in \{\xi, \zeta\}$;

4 Initialize number of iterations $t = 0$;

5 **while** *not converge and* $t \le T$ **do**

6      **for** *player* $i \in \{\xi, \zeta\}$ **do**

7          Train $\pi^{t+1}_i$ by letting it exploit $\pi^t_{*,-i}$;

8          Extend policy set $\Pi^{t+1}_i = \Pi^t_i \cup \{\pi^{t+1}_i\}$;

9      Check convergence and $t = t + 1$;

10      Estimate missing entries in $U(\Pi^{t+1})$ ;

11      Compute mixed strategies $\pi^{t+1}_*$ by solving the matrix game defined by $U(\Pi^{t+1})$ ;

12 **Output:** mixed strategies $\pi_*$ for both players ;

---

$\Pi^{t+1}_i = \Pi^t_i \cup \text{Br}(\pi^t_{*,-i})$ for $i \in \{\xi, \zeta\}$. After extending the population, we complete the utility table $U(\Pi^{t+1})$ and perform the next iteration. The algorithm converges if the best response policies are already in the population. Practically, $\mathfrak{R}$-MACRL approximates the utility function $U(\cdot)$ by having each policy in the population $\Pi^t_i$ playing with each other policy in $\Pi^t_{-i}$ and tracking the average utilities. The approximation of the best response policies is performed by maximizing $\mathcal{J}_\xi(\theta_\mu | \pi^t_{*,\zeta})$ and $\mathcal{J}_\zeta(\theta_d, \theta_r | \pi^t_{*,\xi})$ for the attacker and the defender, where the full expressions of the two objectives are described in Eq. 1 and Eq. 2, respectively.

The overall algorithm of $\mathfrak{R}$-MACRL is presented in Algorithm 1. We first initialize the meta-game and the mixed strategies (line 2 and line 3). Then, at each step, we train the attacker and the defender alternately by letting them best respond to their corresponding opponent (line 7). Next, the learned new policy $\pi^{t+1}_i$ is added to the policy set $\Pi^{t+1}_i$ for the two players respectively (line 8). After extending the policy set, we can check the convergence (line 9) and calculate the missing entries in $U(\Pi^{t+1})$ (line 10). Finally, we solve the new meta-game (line 11) and repeat the iteration.

## 4 EXPERIMENTS

We conduct extensive experiments on a variety of state-of-the-art MACRL algorithms to answer: **Q1:** *Are MACRL methods vulnerable to message attacks and whether the two-stage message filter is able to consistently recover multi-agent coordination?* **Q2:** *Whether $\mathfrak{R}$-MACRL is able to improve the robustness of MACRL algorithms under message attacks?* **Q3:** *Whether our method is able to scale to scenarios where there are multiple attackers?* **Q4:** *Which components contribute to the performance of the method and how does the proposed method work?* We first categorize existing MACRL algorithms and then select representative algorithms to perform the evaluation. All experiments are conducted on a server with 8 NVIDIA Tesla V100 GPUs and a 44-core 2.20GHz Intel(R) Xeon(R) E5-2699 v4 CPU.

| Categories | Methods | Environments |
|---|---|---|
| **CD** | CommNet [32] | Predator Prey (PP) [4] |
| **LC** | TarMAC [6] | Traffic Junction (TJ) [32] |
| **CC** | NDQ [40] | StarCraft II (SCII) [28] |

**Table 1: The chosen algorithms and environments.**

### 4.1 Experimental Setting

MACRL methods are commonly categorized by whether they are implicit or explicit [1, 24]. In implicit communication, the actions of agents influence the observations of the other agents. Whereas in explicit communication, agents have a separate set of communication actions and exchange information via communication actions. In this paper, we focus on explicit communication because in implicit communication, to carry out an attack, the attacker's behaviour is often bizarre, making the attack trivial and easily detectable. We consider the following three realistic types of explicit communication:

- **Communication with delay (CD):** Communication in real world is usually not real-time. We can model this by assuming that it takes one or a few time steps for the messages being received by the targeted agents [32].
- **Local communication (LC):** Messages sometimes cannot be broadcast to all agents due to communication distance limitations or privacy concerns. Therefore, agents need to learn to communicate locally, affecting only those agents within the same communication group [4, 6].
- **Communication with cost or limited bandwidth (CC):** Agents should avoid sending redundant messages because communication in real world is costly and communication channels often have a limited bandwidth [40, 44].

Following the above taxonomy, we select some representative state-of-the-art algorithms to perform evaluation[2]. As in Table 1, we select CommNet [32], TarMAC [6] and NDQ [40] for CD, LC and CC respectively. A brief introduction to each of the chosen algorithms is provided in Appendix A. We evaluate in the following environments, which are similar to those in the paper that first introduced the algorithms:

**Predator Prey (PP).** There are 3 predators, trying to catch 6 prey in a $7 \times 7$ grid. Each predator has local observation of a $3 \times 3$ sub-grid centered around it and can move in one of the 4 compass directions, remain still, or perform catch action. The prey moves randomly and is caught if at least two nearby predators try to catch it simultaneously. The predator will obtain a team reward at 10 for each successful catch and will be punished $p$ for each losing catch action. There are two tasks with $p = 0$ and $p = -0.5$, respectively. A prey will be removed from the grid after it being caught. An episode ends after 60 steps or all preys have been caught.

**Traffic Junction (TJ).** In TJ, there are $N_{max}$ cars and the aim of each car is to complete its pre-assigned route. Each car can observe of a $3 \times 3$ region around it, but is free to communicate with other cars. The action space for each car at every time step is {gas, brake}. The reward function is $-0.01\tau + r_{\text{collision}}$, where $\tau$ is the number of time steps since the activation of the car, and $r_{\text{collision}} = -10$ is a

---

[2]Sweeping the whole list of algorithms for each category is impossible and unnecessary since there have been a lot of algorithms proposed and algorithms in each category usually share common characteristics.
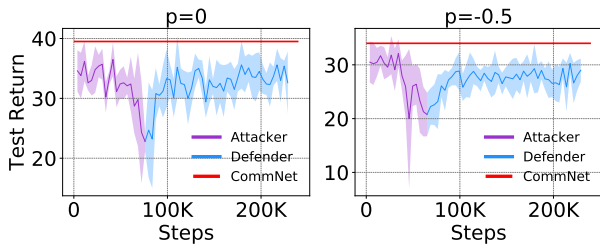
Figure 4: Attack and defence on CommNet.

| | Easy | Hard |
|---|---|---|
| TarMAC | $99.9 \pm 0.1\%$ | $94.9 \pm 0.2\%$ |
| TarMAC w/ $\pi_\xi$ | $87.2 \pm 4.68\%$ | $88.75 \pm 7.29\%$ |
| TarMAC w/ $\pi_\zeta$ | $96.41 \pm 1.38\%$ | $93.23 \pm 8.11\%$ |

Table 2: Attack and defence on TarMAC.



Figure 5: Attack and defence on NDQ.



Figure 6: Exploiting the message filter in CommNet.

collision penalty. Cars become available to be sampled and put back to the environment with new assigned routes once they complete their routes.

**StarCraft II (SCII).** We consider SMAC [28] combat scenarios where the enemy units are controlled by StarCraft II built-in AI (difficulty level is set to hard), and each of the ally units is controlled by a learning agent. The units of the two groups can be asymmetric. The action space contains no-op, move[direction], attack[enemy id], and stop. Agents receive a globally shared reward which calculates the total damage made to the enemy units at each time step. We conduct experiments on four SMAC tasks: 3bane_vs_hM, 4bane_vs_hM, 1o_2r_vs_4r and 1o_3r_vs_4r. More details about the tasks are in Appendix B.

For each of the tasks, we first train the chosen MACRL methods to obtain the action policy and the message policy for each agent. After that, we randomly select an agent to be malicious and train its adversarial policy to examine whether MACRL methods are vulnerable to message attacks. Then we perform defence by training the message filter, during which the adversarial policy of the malicious agent is fixed but keeps working. To show that a single message filter is brittle and can be easily exploited if the attacker adapts to it, we freeze the learned message filter and retrain the adversarial policy. Finally, we integrate the message filter into the framework of $\Re$-MACRL and justify if $\Re$-MACRL is helpful to improve the robustness. All experiments are carried out with five random seeds and results are shown with a 95% confidence interval.

### 4.2 Recovering Multi-Agent Coordination

We first evaluate the performance of our attacking method on the three selected MACRL algorithms. Then we try to recover multi-agent coordination for the attacked algorithms by applying the message filter.

**CommNet.** The experiments for CommNet are conducted on predator prey (PP) [4]. We set the punishment as $p = 0$ and $p = -0.5$ to create two PP tasks with different difficulties. As in Fig. 4, at the beginning of the attack, the performance of CommNet does not have obvious decrease, indicating that injecting random noise into the message is hard to disrupt agent coordination. As we gradually train the adversarial policy, there is a significant drop in the test return, with 40% and 33% decreases in the task of $p = 0$ and $p = -0.5$, respectively. Multi-agent cooperation has been severely affected due to the malicious messages. When the test return decreases to preset thresholds, i.e., 23 for $p = 0$ and 20 for $p = -0.5$, we freeze the adversarial policy network and start to train the message filter. As shown in the blue curves in Fig 4, with the message filter,
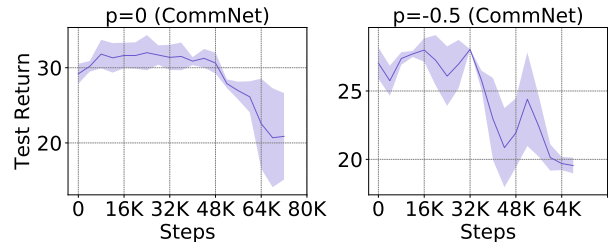
test return steadily approaches the converged value of CommNet (red line), indicating that the message filter can effectively recover multi-agent coordination under attack.

**TarMAC.** We conduct message attack and defence on the Traffic Junction (TJ) environment [32]. There are two modes in TJ, i.e., easy and hard. The easy task has one junction of two one-way roads on a $7 \times 7$ grid with $N_{max} = 5$. In the hard task, its map has four junctions of two-way roads on a $18 \times 18$ grid and $N_{max} = 20$. As shown in Table 2, under attack, the success rate of TarMAC decreases in both the easy and the hard scenarios, demonstrating the vulnerability of TarMAC under malicious messages. After that, we equip TarMAC with the defender, then its performance improves considerably, demonstrating the merit of the message filter.

**NDQ.** We further examine the adversarial communication problem in NQD. We perform evaluation on two StarCraft II Multi-Agent Challenge (SMAC) [28] scenarios, i.e., 3bane_vs_hM and 1o_2r_vs_4r, in which the communication between cooperative agents is necessary [40]. As presented in Fig. 5, under attack, the test win rate of NDQ decreases dramatically, demonstrating that NDQ is also vulnerable to message attacks. After applying the message filter as the defender, we can find the test win rate quickly reaches to around 60% in 3bane_vs_hM and 55% in 1o_2r_vs_4r, demonstrating that, once again, our defence methods succeed in restoring multi-agent coordination under message attacks.

### 4.3 Improving Robustness with $\Re$-MACRL

We have demonstrated that many state-of-the-art MACRL algorithms are vulnerable to message attacks, and after applying the

| Methods | Scenarios | $u_\zeta^\Re$ | $u_\zeta^{vn}$ |
|---------|-----------|---------------|----------------|
| CommNet | $p = 0$ | **41.75 ± 0.00** | 40.74 ± 0.23 |
|         | $p = -0.5$ | **35.38 ± 1.28** | 31.91 ± 0.94 |
| TarMAC | Easy | **−4.08 ± 0.04** | −4.24 ± 0.08 |
|        | Hard | **−9.33 ± 0.29** | −9.80 ± 0.08 |
| NDQ | 3bane_vs_hM | **12.36 ± 0.26** | 11.32 ± 0.15 |
|     | 1o_2r_vs_4r | **17.70 ± 0.16** | 16.33 ± 0.54 |

**Table 3: Expected utilities for the defender trained with $\Re$-MACRL and the vanilla approach.**

message filter, multi-agent coordination can be recovered. In this part, we integrate the message filter into the framework of $\Re$-MACRL and show that the robustness of MACRL algorithms can be significantly improved with $\Re$-MACRL.

**Exploiting the message filter.** We first show that a single message filter is brittle and can be easily exploited if the attacker adapts to it. We perform experiments on CommNet by freezing the message filter and retraining the adversarial policy. As depicted in Fig. 6, the test return of the team gradually decreases as the training proceeds, with around 30% and 20% decreases in the task of $p = 0$ and $p = -0.5$ respectively. We also conduct experiments on NDQ and similar phenomenon is observed (see Appendix D). We conclude that even though the designed message filter is able to recover multi-agent cooperation under message attacks, its performance can degrade if faced with an adaptive attacker.

**Improving robustness.** To illustrate the improvement in robustness, we make comparisons between the defender trained by $\Re$-MACRL and the vanilla defending method. For the defender trained with $\Re$-MACRL, a population of attack policies are learned together with the defender, whose policy is also a mixture of sub-policies. In the vanilla training method, only a single defending policy is trained for the defender. We use the expected utility value (the accumulated team return) as the metric to compare the performance of the defender. Specifically, the larger the expected utility, the better the robustness. We denote the expected utility of the defender trained by $\Re$-MACRL as $u_\zeta^\Re$ and the result for the vanilla one as $u_\zeta^{vn}$. As shown in Table 3, $\Re$-MACRL consistently outperforms the vanilla method over all the algorithms and environments. The improvement of expected utilities indicates that the defender trained by $\Re$-MACRL is more robust. Intuitively, the defender benefits from exploring a wider range of the policy space with $\Re$-MACRL, which enables the defender to maintain multi-agent coordination when faced with attacks.

### 4.4 Scaling to Multiple Attackers

To examine the performance of the method in scenarios where there is more than one attacker, we conduct experiments on NDQ in two new challenge tasks: 4bane_vs_hM and 1o_3r_vs_4r. In the former maps, i.e., 3bane_vs_hM and 1o_2r_vs_4r, there are only three agents in the team, making multiple attackers unreasonable. To mitigate this, we create 4bane_vs_hM and 1o_3r_vs_4r based on 3bane_vs_hM and 1o_2r_vs_4r by increasing one more agent to the team. We randomly sample two agents to be malicious. According to our assumptions (agents have no knowledge about who are
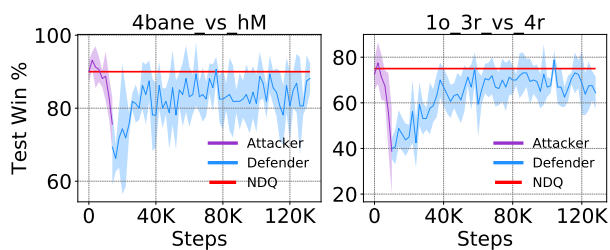


**Figure 7:** *Multiple attackers:* **Attack and defence on NDQ.**

| Methods | Scenarios | $u_\zeta^\Re$ | $u_\zeta^{vn}$ |
|---------|-----------|---------------|----------------|
| NDQ | 4bane_vs_hM | **15.86 ± 0.08** | 15.50 ± 0.29 |
|     | 1o_3r_vs_4r | **18.39 ± 0.80** | 17.03 ± 0.53 |

**Table 4:** *Multiple attackers:* **Expected utilities for the defender trained with $\Re$-MACRL and the vanilla approach.**

malicious and malicious agents cannot communicate with each other), we can directly apply the attacking method to each of the malicious agents. As shown in Fig. 7, the attackers can quickly learn effective adversarial policies with less learning steps. On the other hand, the message filter can still successfully recover multi-agent cooperation under multiple attackers, though it takes much more learning steps to converge. We further evaluate the effectiveness of $\Re$-MACRL on the two tasks and consistent results are observed: trained with $\Re$-MACRL, the agents can obtain larger expected utilities. More results on other algorithms and environments are provided in Appendix D.

### 4.5 Ablation Study and Analysis

We have shown that by integrating the message filter into the framework of $\Re$-MACRL, the robustness of MACRL algorithms can be improved. However, the performance of $\Re$-MACRL directly is affected by the message filter. In this part, we will take a deeper look at how the message filter works.

**Which components contribute to the performance?** The message filter consists of two components: the anomaly detector and the message reconstructor. Here, we alternatively disable these two components in the message filter to study how they contribute to the performance. We run experiments on three scenarios, with each scenario corresponding to an MACRL algorithm. Following the former training procedure, we first train the original MACRL algorithm to obtain the action policy and the message policy; then we sample an attacker and train its adversarial policy by PPO; next we freeze the policy of the attacker and train the message filter (with the anomaly detector or the reconstructor disabled). As in Fig. 8, if we disable the message reconstructor in the message filter and replace it with a random message generator, after being attacked, multi-agent coordination is hard to be recovered, demonstrating the importance of the reconstructor. We further disable the anomaly detector and randomly choose an agent to reconstruct its message, as in Fig. 9, the defending performance is also poor. The ablation illustrates that both the anomaly detector and the reconstructor are critical in the message filter.
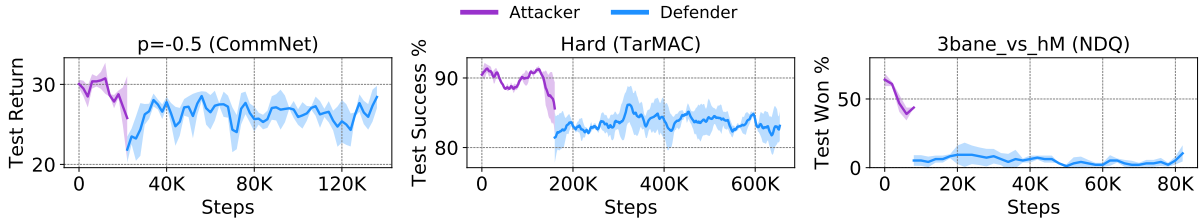
Figure 8: Defending by the message filter with the disabled message reconstructor.
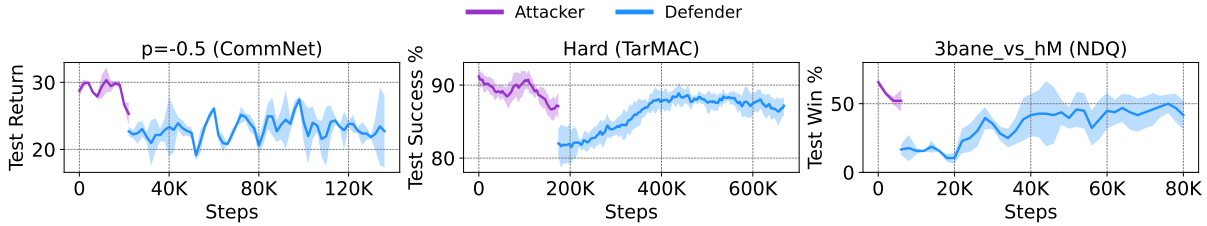


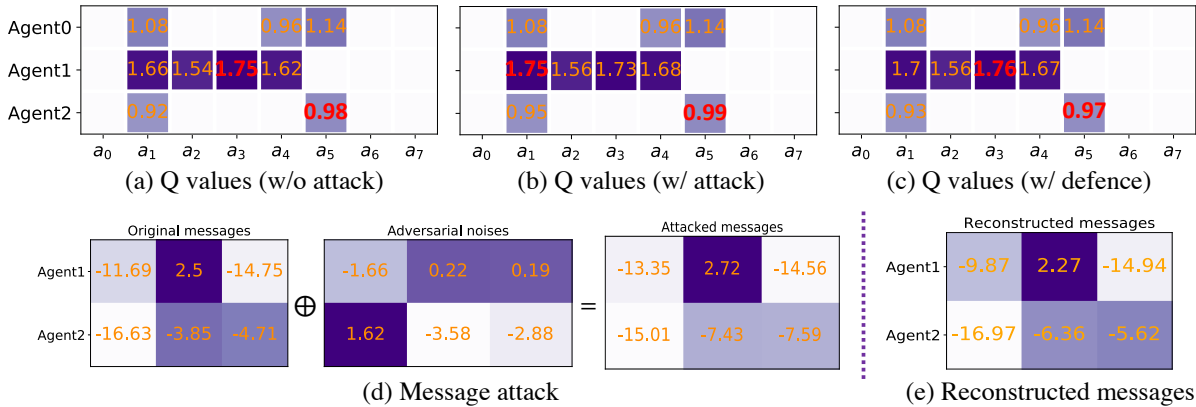Figure 9: Defending by the message filter with the disabled anomaly detector.



Figure 10: Q values and received messages under attack and defence. (a)-(c): Under attack, the optimal action of the benign agent 1 shifts from $a_3$ to $a_1$. After applying the message filter, its optimal action is recovered. (d)-(e): The attacked messages become quite different from the original messages, while the reconstructed messages are to similar the original.

**How does the proposed method work?** Now we take a deeper look at how the message filter works. We conduct experiments on NDQ in the 3bane_vs_hM task. There are three agents in the team, of which agent 0 is the attacker and the others, namely agent 1 and agent 2, are benign agents. As shown in Fig. 10 (a)-(c), the action space of each agent contains eight elements. White cells in Q values correspond to illegal actions at the current state, and the action with red Q value is optimal, e.g., $a_5$ of agent 2. When attacked by the agent 0, the optimal action of agent 1 shifts from $a_3$ to $a_1$, leading to sub-optimality. After applying the message filter, the decision of agent 1 is corrected to $a_3$. We further examine the messages received by agent 1 and agent 2. As in Fig. 10 (d)-(e), the attacked messages are quite different from the original messages, while after applying the message filter, the distance between the original messages and the reconstructed messages becomes closer.

## 5 CONCLUSION

In this paper, we systematically examine the problem of adversarial communication in MACRL. The problem is of importance but has been largely neglected before. We first provide the formulation of adversarial communication. Then we propose an effective method to model message attacks in MACRL. Following that, we design a two-stage message filter to defend against message attacks. Finally, to improve robustness, we formulate the adversarial communication problem as a two-player zero-sum game and propose $\Re$-MACRL to improve the robustness. Experiments on a wide range of algorithms and tasks show that many state-of-the-art MACRL methods are vulnerable to message attacks, while our algorithm can consistently recover multi-agent cooperation and improve the robustness of MACRL algorithms under message attacks.

# REFERENCES

[1] Sanjeevan Ahilan and Peter Dayan. 2021. Correcting Experience Replay for Multi-Agent Communication. *ICLR* (2021).

[2] Marco Barreno, Blaine Nelson, Anthony D Joseph, and J Doug Tygar. 2010. The security of machine learning. *Machine Learning* 81, 2 (2010), 121–148.

[3] Jan Blumenkamp and Amanda Prorok. 2020. The emergence of adversarial communication in multi-agent reinforcement learning. *CoRL* (2020).

[4] Wendelin Böhmer, Vitaly Kurin, and Shimon Whiteson. 2020. Deep coordination graphs. In *ICML*. 980–991.

[5] Nicholas Carlini and David Wagner. 2017. Towards evaluating the robustness of neural networks. In *2017 IEEE Symposium on Security and Privacy (SP)*. 39–57.

[6] Abhishek Das, Théophile Gervet, Joshua Romoff, Dhruv Batra, Devi Parikh, Mike Rabbat, and Joelle Pineau. 2019. Tarmac: Targeted multi-agent communication. In *ICML*. 1538–1546.

[7] Jakob N Foerster, Yannis M Assael, Nando de Freitas, and Shimon Whiteson. 2016. Learning to communicate with Deep multi-agent reinforcement learning. In *NeurIPS*. 2145–2153.

[8] Adam Gleave, Michael Dennis, Cody Wild, Neel Kant, Sergey Levine, and Stuart Russell. 2020. Adversarial Policies: Attacking Deep Reinforcement Learning. In *ICLR*.

[9] Ian Goodfellow, Nicolas Papernot, Sandy Huang, Yan Duan, Pieter Abbeel, and Jack Clark. 2017. Attacking machine learning with adversarial examples. *OpenAI. https://blog. openai. com/adversarial-example-research* (2017).

[10] Ling Huang, Anthony D Joseph, Blaine Nelson, Benjamin IP Rubinstein, and J Doug Tygar. 2011. Adversarial machine learning. In *Proceedings of the 4th ACM Workshop on Security and Artificial Intelligence*. 43–58.

[11] Maximilian Hüttenrauch, Adrian Šošić, and Gerhard Neumann. 2017. Guided deep reinforcement learning for swarm systems. *arXiv preprint arXiv:1709.06011* (2017).

[12] Natasha Jaques, Angeliki Lazaridou, Edward Hughes, Caglar Gulcehre, Pedro Ortega, DJ Strouse, Joel Z Leibo, and Nando De Freitas. 2019. Social influence as intrinsic motivation for multi-agent deep reinforcement learning. In *ICML*. 3040–3049.

[13] Jiechuan Jiang and Zongqing Lu. 2018. Learning attentional communication for multi-agent cooperation. In *NeurIPS*. 7254–7264.

[14] Daewoo Kim, Sangwoo Moon, David Hostallero, Wan Ju Kang, Taeyoung Lee, Kyunghwan Son, and Yung Yi. 2019. Learning to schedule communication in multi-agent reinforcement learning. *ICLR* (2019).

[15] Woojun Kim, Jongeui Park, and Youngchul Sung. 2021. Communication in multi-Agent reinforcement learning: Intention Sharing. *ICLR* (2021).

[16] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).

[17] Hubert Kirrmann. 2015. *Fault Tolerant Computing in Industrial Automation.* Switzerland: ABB Research Center.

[18] Landon Kraemer and Bikramjit Banerjee. 2016. Multi-agent reinforcement learning as a rehearsal for decentralized planning. *Neurocomputing* 190 (2016), 82–94.

[19] Marc Lanctot, Vinicius Zambaldi, Audrūnas Gruslys, Angeliki Lazaridou, Karl Tuyls, Julien Pérolat, David Silver, and Thore Graepel. 2017. A unified game-theoretic approach to multiagent reinforcement learning. In *NeurIPS*. 4193–4206.

[20] Yen-Chen Lin, Zhang-Wei Hong, Yuan-Hong Liao, Meng-Li Shih, Ming-Yu Liu, and Min Sun. 2017. Tactics of adversarial attack on deep reinforcement learning agents. In *IJCAI*. 3756–3762.

[21] H Brendan McMahan, Geoffrey J Gordon, and Avrim Blum. 2003. Planning in the presence of cost functions controlled by an adversary. In *ICML*. 536–543.

[22] Rupert Mitchell, Jan Blumenkamp, and Amanda Prorok. 2020. Gaussian Process Based Message Filtering for Robust Multi-Agent Cooperation in the Presence of Adversarial Communication. *arXiv preprint arXiv:2012.00508* (2020).

[23] Paul Muller, Shayegan Omidshafiei, Mark Rowland, Karl Tuyls, Julien Perolat, Siqi Liu, Daniel Hennes, Luke Marris, Marc Lanctot, Edward Hughes, et al. 2019. A generalized training approach for multiagent learning. In *ICLR*.

[24] Frans A Oliehoek, Christopher Amato, et al. 2016. *A Concise Introduction to Decentralized POMDPs.* Vol. 1. Springer.

[25] Frans A Oliehoek, Matthijs TJ Spaan, and Nikos Vlassis. 2008. Optimal and approximate Q-value functions for decentralized POMDPs. *JAIR* 32 (2008), 289–353.

[26] Nicolas Papernot, Patrick McDaniel, Ian Goodfellow, Somesh Jha, Z Berkay Celik, and Ananthram Swami. 2017. Practical black-box attacks against machine learning. In *Proceedings of the 2017 ACM on Asia Conference on Computer and Communications Security*. 506–519.

[27] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In *NeurIPS*. 8026–8037.

[28] Mikayel Samvelyan, Tabish Rashid, Christian Schroeder de Witt, Gregory Farquhar, Nantas Nardelli, Tim G. J. Rudner, Chia-Man Hung, Philiph H. S. Torr, Jakob Foerster, and Shimon Whiteson. 2019. The StarCraft multi-agent challenge. *CoRR* abs/1902.04043 (2019).

[29] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347* (2017).

[30] Amanpreet Singh, Tushar Jain, and Sainbayar Sukhbaatar. 2018. Learning when to communicate at scale in multiagent cooperative and competitive tasks. *ICLR* (2018).

[31] Arambam James Singh, Akshat Kumar, and Hoong Chuin Lau. 2020. Hierarchical Multiagent Reinforcement Learning for Maritime Traffic Management. In *AAMAS*. 1278–1286.

[32] Sainbayar Sukhbaatar, Arthur Szlam, and Rob Fergus. 2016. Learning multiagent communication with backpropagation. In *NeurIPS*. 2252–2260.

[33] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. 2014. Intriguing properties of neural networks. In *ICLR*.

[34] Keen Security Lab Tencent. 2019. Experimental Security Research of Tesla Autopilot. (2019). https://keenlab.tencent.com/en/whitepapers/Experimental_Security_Research_of_Tesla_Autopilot.pdf

[35] Florian Tramèr, Alexey Kurakin, Nicolas Papernot, Ian Goodfellow, Dan Boneh, and Patrick McDaniel. 2018. Ensemble adversarial training: attacks and defenses. In *ICLR*.

[36] James Tu, Tsunhsuan Wang, Jingkang Wang, Sivabalan Manivasagam, Mengye Ren, and Raquel Urtasun. 2021. Adversarial attacks on multi-agent communication. *arXiv preprint arXiv:2101.06560* (2021).

[37] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *NeurIPS*. 5998–6008.

[38] Oriol Vinyals, Igor Babuschkin, Wojciech M Czarnecki, Michaël Mathieu, Andrew Dudzik, Junyoung Chung, David H Choi, Richard Powell, Timo Ewalds, Petko Georgiev, et al. 2019. Grandmaster level in StarCraft II using multi-agent reinforcement learning. *Nature* 575, 7782 (2019), 350–354.

[39] Rundong Wang, Xu He, Runsheng Yu, Wei Qiu, Bo An, and Zinovi Rabinovich. 2020. Learning efficient multi-agent communication: An information bottleneck approach. In *ICML*. 9908–9918.

[40] Tonghan Wang, Jianhao Wang, Chongyi Zheng, and Chongjie Zhang. 2019. Learning Nearly Decomposable Value Functions Via Communication Minimization. In *ICLR*.

[41] Hang Xu, Rundong Wang, Lev Raizman, and Zinovi Rabinovich. 2021. Transferable environment poisoning: Training-time attack on reinforcement learning. In *AAMAS*. 1398–1406.

[42] Huan Zhang, Hongge Chen, Duane S Boning, and Cho-Jui Hsieh. 2021. Robust reinforcement learning on state observations with learned optimal adversary. In *ICLR*.

[43] Huan Zhang, Hongge Chen, Chaowei Xiao, Bo Li, Mingyan Liu, Duane Boning, and Cho-Jui Hsieh. 2020. Robust deep reinforcement learning against adversarial perturbations on state observations. In *NeurIPS*. 21024–21037.

[44] Sai Qian Zhang, Jieyu Lin, and Qi Zhang. 2020. Succinct and robust multi-agent communication with temporal message control. *arXiv preprint arXiv:2010.14391* (2020), 17271–17282.

## A INTRODUCTION TO THE SELECTED MACRL ALGORITHMS

In this section, we introduce the selected MACRL algorithms for evaluation:

**CommNet [32]** is a simple yet effective multi-agent communicative algorithm where incoming messages of each agent are generated by averaging the messages sent by all agents in the last time step. We present the architecture of CommNet in Fig. 13.

**TarMAC [6]** extends CommNet by allowing agents to pay attention to important parts of the incoming messages via Attention [37] network. Concretely, each agent generates signature and query vectors when sending message. In the message receiving phase, attention weights of incoming messages are calculated by comparing the similarity between the query vector and the signature vector of each incoming message. Then a weighted sum over all incoming messages is performed to determine the message an agent will receive for decentralized execution. The architecture of TarMAC is in Fig. 14.

**NDQ [40]** achieves nearly decomposable Q-functions via communication minimization. Specifically, it trains the communication model by maximizing mutual information between agents' action selection and communication messages while minimizing the entropy of messages between agents. Each agent broadcasts messages to all other agents. The sent messages are sampled from learned distributions which are optimized by mutual information. The overview of NDQ is in Fig. 15.
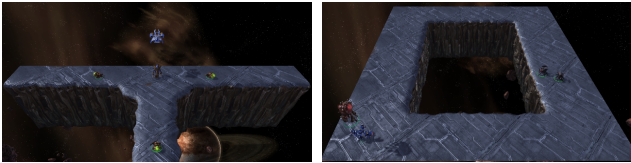
## B THE SMAC MAPS



**Figure 11: Snapshots of the StarCraft II scenarios. Left: 3bane_vs_hM. Right: 1o_2r_vs_4r.**

For the SCII environments, we conduct experiments on the following four maps to examine the attack and defence on NDQ:

**3bane_vs_hM:** 3 Banelings try to kill a Hydralisk assisted by a Medivac. 3 Banelings together can just blow up the Hydralisk. Therefore, to win the game, 3 Banelings ally units should not give the Hydralisk changes to rest time when the Medivac can restore its health.

**4bane_vs_hM:** Similar to 3bane_vs_hM, the main difference in 3bane_vs_hM is there are 4 Banelings.

**1o_2r_vs_4r:** Ally units consist of 1 Overseer and 2 Roaches. The Overseer can find 4 Reapers and then notify its teammates to kill the invading enemy units, the Reapers, in order to win the game. At the start of an episode, ally units spawn at a random point on the map while enemy units are initialized at another random point. Given that only the Overseer knows the position of the enemy unit, the ability to learn to deliver this message to its teammates is vital for effectively winning the combat.

**1o_3r_vs_4r:** Similar to 1o_2r_vs_4r, the main difference in 1o_3r_vs_4r is there are 1 Overseer and 3 Roaches.

## C HYPERPARAMETERS

We train all the selected MACRL methods with the same hyperparameters as in the corresponding papers, to reproduce the performances. We use Adam [16] optimizer with a learning rate of 0.001 to train the attacker model as well as the anomaly detector and the message reconstructor. We use the default hyperparameters of Adam optimizer. $\alpha$, $\beta_1$ and $\beta_2$ are all 0.001. The clip ratio in PPO is 0.5. The attacker model contains 3 linear layers and each layer has 64 neurons with ReLU activation. Same deep neural networks are used in the anomaly detector and the message reconstructor. We use PyTorch [27] to implement all the deep neural network models.

## D ADDITIONAL RESULTS

We train an adaptive attacker in NDQ to exploit the message filter. As in Fig. 12, the test won rate decrease significantly, illustrating that a single message filter is brittle and can easily be exploited.
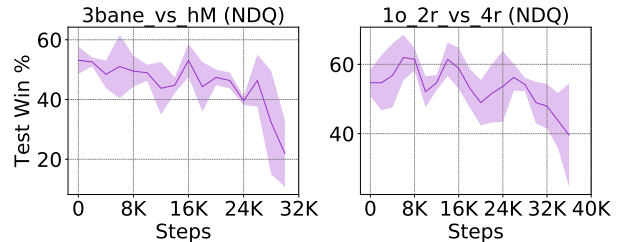


**Figure 12: Exploiting the message filter in NDQ.**

We also conduct experiments of multi attackers on the original environments. Two agents are selected randomly to be malicious. As in Table. 5, $\Re$ outperforms the vanilla method in most of the tasks, demonstrating the effectiveness of our method.

| Methods | Scenarios | $u_\zeta^\Re$ | $u_\zeta^{vn}$ |
|---------|-----------|---------------|----------------|
| CommNet | $p = 0$ | $41.64 \pm 1.90$ | $41.25 \pm 1.21$ |
|         | $p = -0.5$ | $\mathbf{36.83 \pm 1.71}$ | $35.27 \pm 0.66$ |
| TarMAC  | Easy | $\mathbf{-0.70 \pm 0.00}$ | $-0.80 \pm 0.00$ |
|         | Hard | $\mathbf{-0.90 \pm 0.00}$ | $-1.00 \pm 0.00$ |
| NDQ     | 3bane_vs_hM | $\mathbf{12.14 \pm 0.35}$ | $11.46 \pm 0.23$ |
|         | 1o_2r_vs_4r | $\mathbf{17.72 \pm 0.84}$ | $16.77 \pm 0.62$ |

**Table 5:** *Multiple attackers:* **Expected utilities for the defender trained with $\Re$-MACRL and the vanilla approach.**
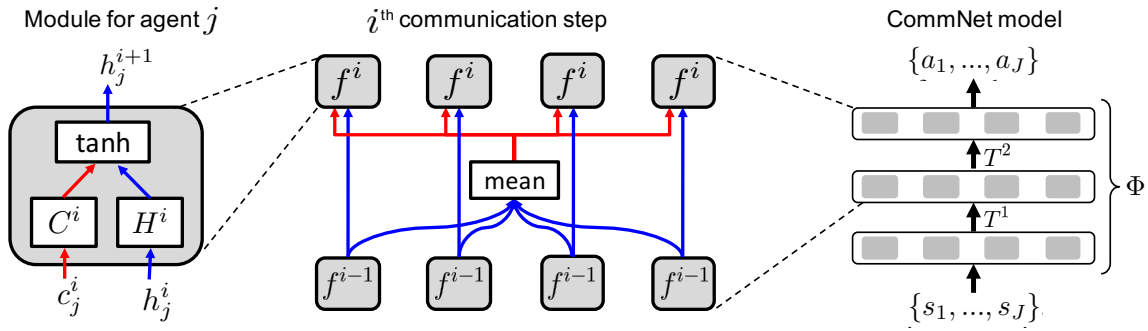
**Figure 13: The Architecture of CommNet from [32].** An overview of our CommNet model. **Left:** view of module $f^i$ for a single agent $j$. Note that the parameters are shared across all agents. **Middle:** a single communication step, where each agents modules propagate their internal state $h$, as well as broadcasting a communication vector $c$ on a common channel (shown in red). **Right:** full model $\Phi$, showing input states $s$ for each agent, two communication steps and the output actions for each agent.
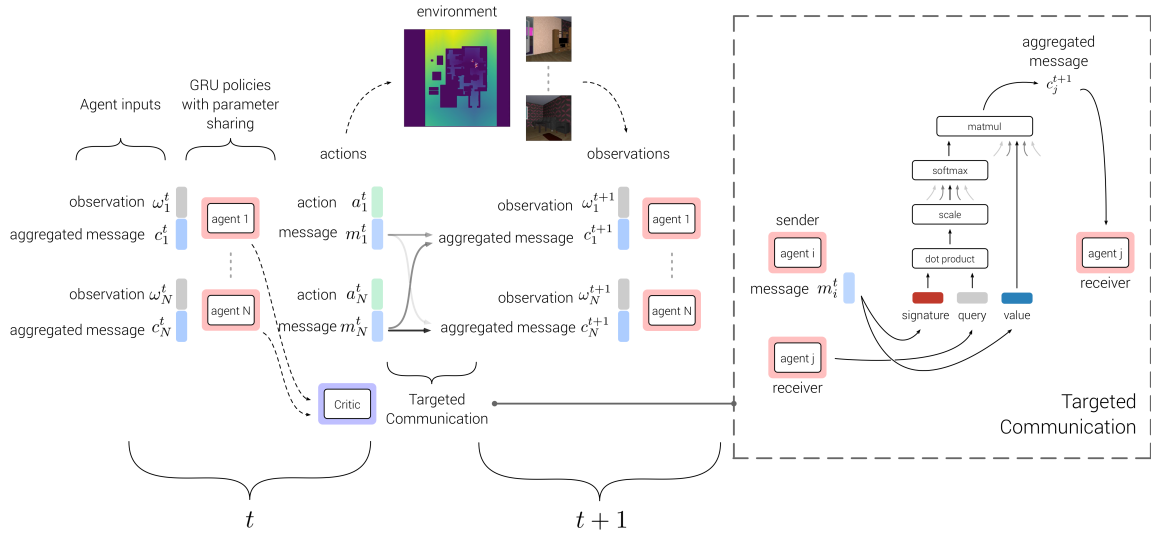


**Figure 14: Overview of TarMAC from [6]. Left:** At every timestep, each agent policy gets a local observation $\omega_i^t$ and aggregated message $c_i^t$ as input, and predicts an environment action $a_i^t$ and a targeted communication message $m_i^t$. **Right:** Targeted communication between agents is implemented as a signature-based soft attention mechanism. Each agent broadcasts a message $m_i^t$ consisting of a signature $k_i^t$, which can be used to encode agent-specific information and a value $v_i^t$, which contains the actual message. At the next timestep, each receiving agent gets as input a convex combination of message values, where the attention weights are obtained by a dot product between sender's signature $k_i^t$ and a query vector $q_j^{t+1}$ predicted from the receiver's hidden state.
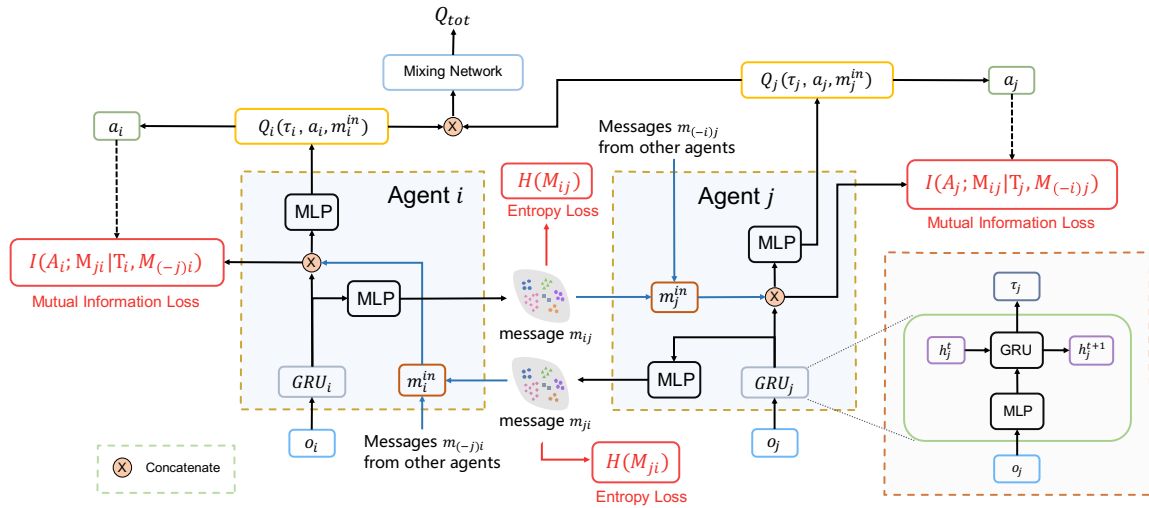
Figure 15: Overview of NDQ from [40]. The message encoder generates an embedding distribution that is sampled and concatenated with the current local history to serve as an input to the local action-value function. Local action values are fed into a mixing network to to get an estimation of the global action value..