# Robust Communicative Multi-Agent Reinforcement Learning with Active Defense

## Lebin Yu, Yunbo Qiu, Quanming Yao, Yuan Shen, Xudong Zhang and Jian Wang [*]

Department of Electronic Engineering, BNRist, Tsinghua University, Beijing 100084, China
{yulb19, qyb18}@mails.tsinghua.edu.cn, quanmingyao@gmail.com, {shenyuan_ee, zhangxd, jianwang}@tsinghua.edu.cn

## Abstract

Communication in multi-agent reinforcement learning (MARL) has been proven to effectively promote cooperation among agents recently. Since communication in real-world scenarios is vulnerable to noises and adversarial attacks, it is crucial to develop robust communicative MARL technique. However, existing research in this domain has predominantly focused on passive defense strategies, where agents receive all messages equally, making it hard to balance performance and robustness. We propose an active defense strategy, where agents automatically reduce the impact of potentially harmful messages on the final decision. There are two challenges to implement this strategy, that are defining unreliable messages and adjusting the unreliable messages' impact on the final decision properly. To address them, we design an Active Defense Multi-Agent Communication framework (ADMAC), which estimates the reliability of received messages and adjusts their impact on the final decision accordingly with the help of a decomposable decision structure. The superiority of ADMAC over existing methods is validated by experiments in three communication-critical tasks under four types of attacks.

## Introduction

In recent years, multi-agent reinforcement learning (MARL) has made remarkable strides in enhancing cooperative robot tasks and distributed control domains, exemplified by traffic lights control (Chu, Chinchali, and Katti 2020) and robots navigation (Han, Chen, and Hao 2020). Given the inherent partial observability in multi-agent tasks, several researchers have explored the integration of communication mechanisms to facilitate information exchange among agents (Ma, Luo, and Pan 2021).

Nevertheless, the utilization of multi-agent communication in real-world applications introduces certain challenges. In such scenarios, agents rely on wireless communication channels to exchange messages, which are susceptible to various sources of noises and interference. These perturbations on messages, especially malicious ones, can severely reduce the performance of multi-agent systems, even though agents get accurate observations of the environment (Sun

et al. 2023). Hence, it becomes imperative to address these issues and devise robust communicative MARL mechanisms.

Adversarial attacks and defenses in communicative MARL receive much less attention when compared to their counterparts in reinforcement learning (Mu et al. 2022). Current frameworks in this area (Ishii, Wang, and Feng 2022; Yuan et al. 2023) commonly follow the principle of passive defense, where agents treat all received messages equally and try to make relative safe decisions. Since perturbed messages are mixed with useful messages, this indiscriminate reception may lead to the result of "garbage in, garbage out".

We notice that robust communicative MARL has an important feature compared with robust RL: The attackers are only allowed to modify a part of the messages, while the modification is unlimited and perturbed messages can be quite different from the original ones. Inspired by noisy learning (Han et al. 2018), we propose an active defense strategy to utilize this feature: agents actively judge the reliability of messages based on their own unperturbed observations and hidden states (which contain history information) and reduce unreliable messages' impact on the final decision.[1] For example, if in a search task an agent receives a message saying that "target is at coordinates (1,1), get to it!", while the agent have searched (1,1) and found nothing, it can realize that this message is fake. We also visualize the difference between active defense and passive defense in Fig. 1 for better demonstration.

Nevertheless, there remain two key challenges to implement active defense. (I) What kind of messages should be defined as unreliable? (II) How to adjust the unreliable messages' impact on the final decision properly? Since the ultimate goal for multi-agent communication is to make agents better finish cooperation tasks, a message beneficial for this goal should be considered reliable. However, common communicative MARL frameworks aggregate received messages with highly nonlinear neural networks, making it hard

---

[*]Corresponding Author

---

[1]An intuitive approach is identifying the perturbed messages and simply remove them from decision-making. However, for an agent, the received messages themselves contain some unknown information, making it infeasible to accurately judge whether a message should be completely trusted or ignored. Besides, it is hard to determine a proper decision threshold for judging whether a message is reliable or unreliable.
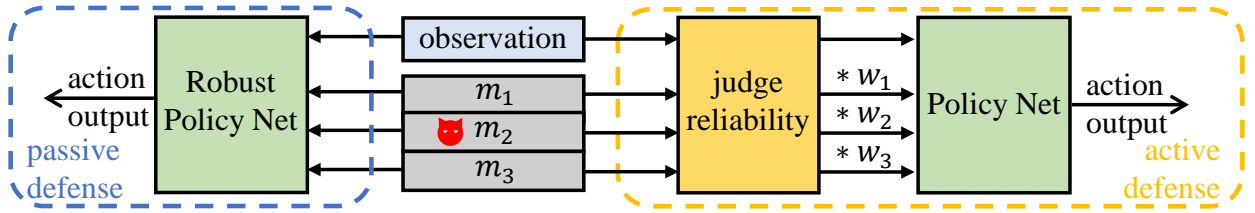
Figure 1: The figure shows the difference between active defense and passive defense. Passive defense takes in harmful messages along with useful ones equally, and try to make robust decisions. This indiscriminate reception makes it difficult to optimize robustness and performance simultaneously. In comparison, active defense assesses the reliability of incoming messages with local information first and reduces the weight of potentially malicious messages. For example, $w_2$ in this figure is expected to be much lower than $w_1$ and $w_3$. By doing so, this framework enables the policy network to efficiently extract information from reliable messages with reduced interference.

to evaluate or adjust the impact of a sole message on the final decision.

To address the aforementioned challenges and achieve active defense, we introduce an Active Defense Multi-Agent Communication (ADMAC) framework comprising two key components. The first is reliability estimator, which outputs the estimated reliability (ranging from 0 to 1) of a certain message according to the agent's observation and hidden state. The second is decomposable message aggregation policy net, which decomposes the impact of each message on the final decision, and allows adjustable weights to control this impact. Below is an introduction of how an agent processes messages and makes decisions using ADMAC: Firstly, it generates a **base action preference vector** based on the observation and hidden state, with each component corresponding to a feasible action. Secondly, for each message received, the agent generates a **message action preference vector** according to the message and the agent's observation. Thirdly, the reliability estimator outputs the reliability values of messages, which then serve as the weights of corresponding message action preference vectors. Lastly, the agent adds all weighted preference vectors to get the total action preference vector and feeds it into Softmax to derive the final action distribution.

In order to substantiate the robustness of ADMAC, we evaluate it alongside three alternative communication frameworks in three communication-critical cooperation tasks. Besides, we implement four distinct types of attacks, including Gaussian attack (He et al. 2023), Monte-Carlo adversarial attack, Fast Gradient Sign Method (Goodfellow, Shlens, and Szegedy 2014) and Projected Gradient Descent (Madry et al. 2018). Experiments confirm the superiority of AD-MAC in terms of its robustness and performance compared to the alternatives. Moreover, we conduct ablation study to further explore the features of components in ADMAC.

Our contributions can be summarized below:

1. We elucidate a significant distinction between robust RL and robust communicative MARL, highlighting the unsuitability of the commonly employed passive defense strategy in robust RL for the latter modeling. Moreover, we propose the idea of active defense which exploits the features of robust communicative MARL.

2. There remain two challenges to implement active defense, that are defining unreliable messages and adjusting the unreliable messages' impact on the final decision properly. To overcome them, we propose ADMAC, which incorporates mechanisms that enable agents to lower the impact of potentially malicious perturbed messages on the final decision, thereby achieving robustness against attacks.

3. We empirically demonstrate the superiority of ADMAC over existing approaches, and conduct ablation study to delve deeper into the unique features and characteristics of it.

## Preliminaries

### Dec-POMDP with Communication

Our modeling is based on a decentralized partially observable Markov decision process (Littman 1994) with $N$ agents in the system. At timestep $t$, the global state is $s^t$, and agent $i$ receives observation $o_i^t$ from the environment and chooses an action $a_i^t$ to perform. Then the environment provides agents with rewards $r_i^t$ and updates the global state from $s^t$ to $s^{t+1}$ according to all agents' actions.

There are numerous communication frameworks in communicative MARL, including scheduler-based (Rangwala and Williams 2020), inquiry-reply-based (Ma, Luo, and Pan 2021) and GNN-based (Pu et al. 2022) ones. They basically follow the communication and decision process specified below, regardless of the difference in communication architectures.

At timestep $t$, agent $i$ with hidden states $h_i^{t-1}$ first receives $o_i^t$, based on which it generates message and communicate with other agents. Then, it updates its hidden states and chooses action $a_i^t$ with a policy network according to its current hidden state, observation and messages from others $m_1^t, m_2^t, ..., m_N^t$. The objective function that needs to be maximized for agent $i$ is the discounted return:

$$J(\theta) = \mathbb{E}[\sum_t \gamma^t r_i^t | \theta], \qquad (1)$$

where $\theta$ denotes the neural network parameters of the agents.

## Attack against Communication

Attack against multi-agent communication has been a hot topic (Xue et al. 2022; Sun et al. 2023) recently since wireless communication is vulnerable to distractions and noises. These works commonly follow the setting that attackers only interfere a part of messages in the multi-agent system to disrupt collaboration. Based on them, we consider the following attack model.

Suppose there are $N$ agents in the system, and at each timestep there are at most $N \times N$ messages. For each message $m_j^t$, the attacker has a probability $p$ to change it to $\hat{m}_j^t$. Then agents receive messages without knowing which are perturbed.

The attack may be adversarial or non-adversarial, depending on how much information the attacker has about the agents. For adversarial attack, we adopt the setting of Zhang et al. where the attacker tries to minimize the chosen probability or preference of the best action, denoted by $\hat{P}(a_{i,best}^t)$:

$$\hat{m}_j^t = f_A(m_j^t) = \arg\min \sum_{i \neq j} \hat{P}(a_{i,best}^t). \qquad (2)$$

Another feasible attack target is to maximize the KL divergence between $a_i^t$ and $\hat{a}_i^t$, with $a_i^t$ representing the action distribution output by agent $i$ when receiving raw message $m_j^t$ and $\hat{a}_i^t$ representing the action distribution output by agent $i$ when receiving perturbed message $\hat{m}_j^t$:

$$\hat{m}_j^t = f_B(m_j^t) = \arg\max \sum_{i \neq j} D_{KL}(\hat{a}_i^t | a_i^t). \qquad (3)$$

Following the setting of (Sun et al. 2023), we consider a strong attack model to better evaluate the robustness of models: The strength of the perturbation, denoted as $||\hat{m}_j^t - m_j^t||$, is not bounded. Consequently, there are multiple ways to implement $f_A(\cdot)$ and $f_B(\cdot)$, which will be detailed in the Experiments Section.

## Active Defense Multi-Agent Communication Framework

We propose an active defense idea for robustness, namely, making agents judge the reliability of received messages based on their local information and reduce unreliable messages' impact on the final decision. However, the implementation brings two challenges: (I) How to define "unreliable" messages? (II) How to adjust the unreliable messages' impact on the final decision? In this section, we propose AD-MAC comprising a decomposable message aggregation policy net and a reliability estimator to address these two challenges, whose visualization is presented in Fig.2.

### Decomposable Message Aggregation Policy Net

The decomposable message aggregation policy net $f_P$ is designed to decompose the impact of each message on the final decision by restricting their influence to action preference vectors. Specifically, $f_P$ consists of three parameterized modules: a GRU module $f_{HP}$ used to update hidden states with observations, a base action generation module

$f_{BP}$ that generates **base action preference vectors** according to the updated hidden states, and a message-observation process module $f_{MP}$ that generates **message action preference vectors** according to the observations and received messages. Suppose there are $K$ actions to choose from, then an action preference vector has $K$ components, each representing the preference for the corresponding action. Use $p_i^t = f_P(o_i^t, h_i^t, m_1^t, ..., m_N^t)$ to denote the final output action distribution for agent $i$ at timestep $t$, where the $k$-th component of $p_i^t$, denoted by $p_i^t[k]$, refers to the probability of choosing the $k$-th action, the decision process is formulated below:

$$
\begin{aligned}
h_i^t &= f_{HP}(h_i^{t-1}, o_i^t), \\
v_i^t &= f_{BP}(h_i^t) + \sum_{j \neq i} w_i(m_j^t) f_{MP}(o_i^t, m_j^t), \\
p_i^t[k] &= e^{v_i^t[k]} / \sum_k e^{v_i^t[k]}.
\end{aligned}
\qquad (4)
$$

where $v_i^t$ is the **total action preference vector**, and $w_i(m_j^t)$ is a weight determined by the reliability estimator detailed in the next subsection. $w_i(m_j^t)$ is set to 1 by default if no robustness is required. Evidently, messages with larger weights have a stronger influence on the final decision. Therefore, for a message considered to be malicious, reducing its weight can effectively attenuate its impact on the final decision. We provide the following proposition to characterize this feature:

**Proposition 1** *For an agent making decisions using (4), if message $m_j^t$ recommends an action most or least recommends an action to the agent, incorporating this message into decision-making must increase or decrease the probability of choosing this action in the final decision, and the magnitude of this effect varies monotonically with the weight $w_i(m_j^t)$.*

The proof is presented in Appendix A. Here "$m_j^t$ recommends the $k_a$-th action most" means $k_a = \arg\max_k f_{MP}(o_i^t, m_j^t)[k]$, and "$m_j^t$ recommends the $k_b$-th action least" means $k_b = \arg\min_k f_{MP}(o_i^t, m_j^t)[k]$.

### Reliability Estimator

The reliability estimator is a classifier $f_R(h_i^t, o_i^t, m_j^t)$ that judges whether a message $m_j^t$ is reliable for agent $i$ with the help of the agent's hidden state $h_i^t$ and observation $o_i^t$. The output of $f_R(h_i^t, o_i^t, m_j^t)$ is a vector of length 2 normalized by Softmax. Let

$$w_i(m_j^t) = f_R(h_i^t, o_i^t, m_j^t)[0]. \qquad (5)$$

$w_i(m_j^t)$ represents the extent to which the reliability estimator thinks $m_j^t$ is reliable for agent $i$, and is used as the weight of $m_j^t$ in (4). Besides, judging the reliability of messages can be treated as a binary classification problem, and the key challenge here is labelling data, i.e., defining what messages are "reliable" and what are not. Since the ultimate goal of communicative MARL is maximizing cooperation performance represented by (1), messages should be labelled
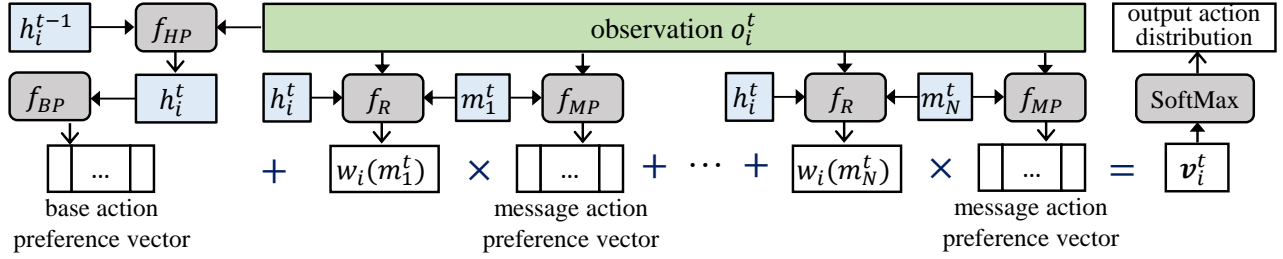
Figure 2: The figure shows how agent $i$ with hidden state $h_i^{t-1}$ generates an output action distribution within ADMAC after receiving an observation $o_i^t$ and messages $m_1^t, m_2^t, ..., m_N^t$ from others. It is noteworthy that the length of actions preference vectors is the same as the number of feasible actions, and each component of the vectors represents the preference for the corresponding action. As depicted in the figure, the impact of each message on the final decision is restricted to the respective action preference vector and can be regulated by the weight $w_i(m_j^t)$.

according to whether they are conducive to achieving this goal.

With the assumption that agents' policies are well-trained, we use the following criteria to label messages: For an agent, if a received message recommends it to choose the best action, then the message is considered to be reliable, otherwise it is bad. Here "best action" refers to the action most likely to be chosen in the absence of perturbations, and "$m_j^t$ recommends the $k_r$-th action" is defined as $f_{MP}(o_i^t, m_j^t)[k_r] > \sum_k f_{MP}(o_i^t, m_j^t)[k]/K$. With the labelled messages, the reliability estimator can be trained in a supervised learning way.

## Instantiation

ADMAC only specifies how agents process messages as well as observations to make robust decisions, and is compatible with numerous communication architectures and training algorithms. In this paper, we adopt a basic broadcast communication mechanism (Singh, Jain, and Sukhbaatar 2018) and use the following paradigm to train our model. More details are provided in Appendix B.

**Stage 1: training policy net**: The training goal of this stage is to optimize the parameters of the decomposable message aggregation policy net and message encoders to maximize the objective function defined in (1). Notably, messages used in our framework are one-dimensional vectors with each component ranging from $-1$ to $1$, and this stage does not involve any attacks or defenses. Consequently, the agents interact with the environment, communicate with each other and update the parameters in a traditional communicative MARL way.

**Stage 2: generating dataset**: After the policies are well-trained, let agents interact with the environment for several episodes to generate training dataset for reliability estimator. To enhance its identification ability, we implement two representational and strong attacks from (Sun et al. 2023) during training:

**Attack I. Random Perturbation**: It suits the scenario where the attacker has no information about the agents. Attackers generate random perturbed messages $\hat{m}_j^t$ to replace the original ones, which means each component of $\hat{m}_j^t$ is a sample from Uniform distribution on $(-1, 1)$.

**Attack II. (L2-normed) Gradient descent adversarial attack**: It suits the scenario where the attacker has full knowledge of the agents. To maximize the adversarial attack objective presented in (2), gradient descent can be utilized to generate a perturbed message:

$$\hat{m}_j^t = m_j^t + \lambda \nabla_{m_j^t} f_A(m_j^t)/||\nabla_{m_j^t} f_A(m_j^t)||_2 \quad (6)$$

When creating dataset for the reliability estimator, we randomly replace 1/3 raw messages with Attack I messages and 1/3 messages with Attack II messages, and let agents make decisions based on them. The decisions will not be truly executed, on the contrary, they are only used to label these messages according to the aforementioned criteria (It is possible that some perturbed messages are labelled reliable, and unperturbed messages are labelled unreliable). All messages and corresponding observations, hidden states, and labels are collected to form a dataset for supervised learning.

**Stage 3: training reliability estimator**: Train the reliability estimator parameterized via MLP based on the dataset using Adam optimizer and cross-entropy loss function. A well-trained reliability estimator is expected to assign low weights to messages that are not conducive to an agent's selection of the optimal action.

## Experiments
### Experimental Environments

We implement three communication-critical multi-agent environments for demonstrative purposes: Food Collector (Sun et al. 2023), Predator Prey (Singh, Jain, and Sukhbaatar 2018), and Treasure Hunt (Freed et al. 2020). They either adhere to a predefined communication setting or a learned communication setting. We test these two kinds of communication settings because the messages within them have different distributions, potentially influencing the performance of attack and defense. Within all environments, one episode ends if all agents finish their tasks or the timestep reaches the upper limit $t_{max}$. Therefore, lower average timesteps indicate better performance. These environments are detailed below and visualized in Fig.3.
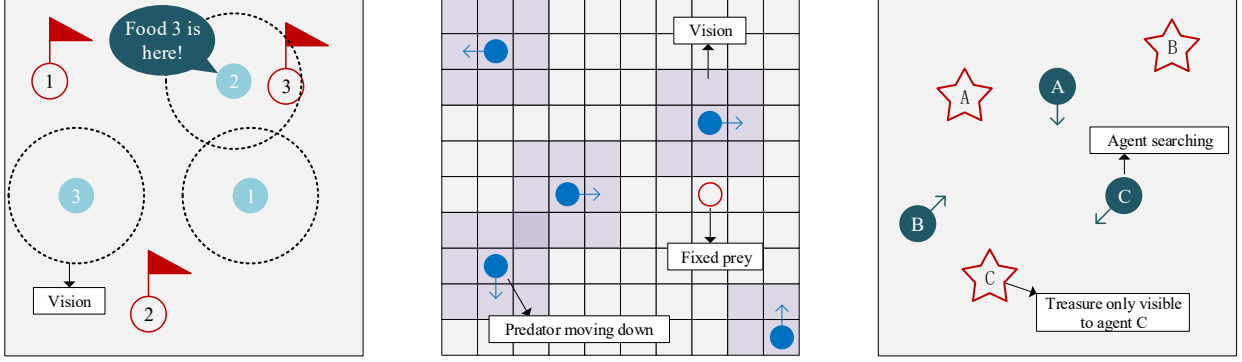
Figure 3: Visualizations of our three experiment environments: Food Collector, Predator Prey, and Treasure Hunt.

**Food Collector (predefined communication)** In this task, $N = 5$ agents with different IDs search for foods with the same IDs in a $1 \times 1$ field. Each agent can observe targets within its vision $d = 0.2$, and moves in eight directions at a speed $v = 0.15$. If an agent finds a target with a different ID, it will kindly broadcasts the coordinates to help the corresponding agent find it.

**Predator Prey (learned communication)** In this task, $N = 5$ agents with vision 1 are required to reach a fixed prey in a grid world of size $10 \times 10$. Due to the severely limited perception, agents must communicate with others to finish tasks earlier. For example, they can delineate their respective search areas through communication, and the first agent to reach the prey can tell others the coordinates.

**Treasure Hunt (learned communication)** In this task, $N = 5$ agents work together to hunt treasures in a field of size $1 \times 1$ with movement speed $v = 0.9$. Each agent obtains the coordinates of its own treasure, which is invisible to others. Note that an agent cannot collect its treasure by itself. Instead, it should help others hunt it through learned communication.

## Tested Algorithms

We evaluate our proposed ADMAC alongside three alternatives introduced below:

**Targeted Multi-Agent Communication (TARMAC)** (Das et al. 2019) It is a baseline communication framework where agents utilize an attention mechanism to determine the weights of receiving messages. It does not include any robust techniques, and its performance demonstrates the impact of attacks without any defense.

**Adversarial Training (AT)** (Pattanaik et al. 2018; Tu et al. 2021) This framework adopts the most frequently used robust learning technique, which is performing adversarial attacks during training to gain robustness.

**Ablated Message Ensemble (AME)** (Sun et al. 2023) It constructs a message-ensemble policy that aggregates multiple randomly ablated message sets. The key idea is that since only a small portion ($< 50\%$) of received messages are harmful, making agents take the consensus of received

messages should provide robustness.

All tested frameworks are trained with an improved version of REINFORCE (Williams 1992). With the assumption that all agents of one task are homogeneous, they share policy network and reliability estimator parameters to accelerate training. More information about training is presented in Appendix B.

## Implemented Attacks

When conducting experimental evaluations of the models, we implement the following four kinds of attacks, including adversarial and non-adversarial ones:

**Attack III Gaussian attack** (He et al. 2023): Attacker adds Gaussian Noises to the messages, i.e.

$$\hat{m}_j^t = m_j^t + \sigma N(0,1). \tag{7}$$

We set $\sigma = 0.5$ in the experiments.

**Attack IV Monte-Carlo adversarial attack**: Randomly generate 10 messages, and find the one that maximizes the attack objective $f(m_j^t)$ defined in (2) or (3).

**Attack V Fast Gradient Sign Method**(Goodfellow, Shlens, and Szegedy 2014):

$$\hat{m}_j^t = m_j^t + \eta sign(\nabla_{m_j^t} f(m_j^t)), \tag{8}$$

Then range of $m$ is $(-1, 1)$, and we set $\eta = 1$ to obtain a strong attack.

**Attack VI Projected Gradient Descent** (Madry et al. 2018): PGD can be treated as a multi-step version of FGSM:

$$\hat{m}_j^{t,(i+1)} = \hat{m}_j^{t,(i)} + \epsilon sign(\nabla_{\hat{m}_j^{t,(i)}} f(\hat{m}_j^{t,(i)})), \tag{9}$$

where $\hat{m}_j^{t,(0)} = m_j^t$. We set $\epsilon = 0.3$ and use 5-step updates to obtain the final perturbed messages.

For adversarial attacks (IV, V and VI), the attack objective can be chosen from $f_A(\cdot)$ and $f_B(\cdot)$, defined accordingly in (2) and (3). It is noteworthy that although they may have the same objectives with attack-II, which is used to train ADMAC, the generated perturbed messages belong to different distributions because they are calculated differently. In other words, our ADMAC framework DOES NOT have prior information of the above attacks used for evaluation.
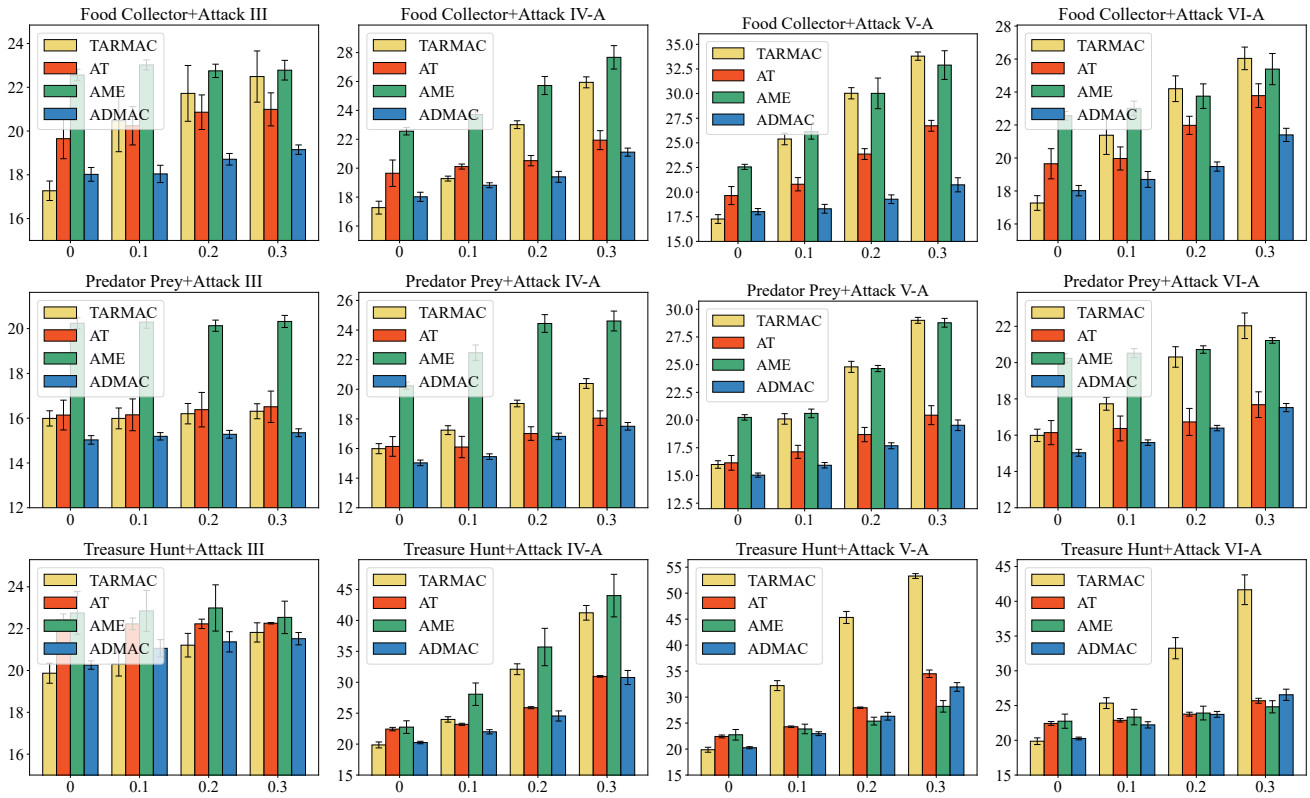
Figure 4: Results for the main experiments. The x-axis represents the attack probability, and the y-axis represents the timesteps required for task completion, where smaller values signify better performance. For each setting (e.g. TARMAC in Treasure Hunt), we train five models with different seeds. During test, we run 500 episodes for each setting and plot the mean as well as standard error of the averaged timesteps. Lower timesteps indicate better performance.

## Performance under Attacks

We evaluate the four multi-agent communication frameworks with the aforementioned tasks and attacks. Additionally, we set the attack probability $p$ to different values to observe how the performance of the agents decays as the attack intensity increases. Besides, since attack objective $f_B(\cdot)$ is not applicable to AME, we only present adversarial attacks with objective function $f_A(\cdot)$ in the main experiments, and provide additional experimental results in Appendix C. From the results presented in Fig.4, the following conclusions can be drawn.

Attack III is the weakest attack, and its attack effect on TARMAC (the baseline non-robust framework) indicates the extent to which messages can influence decision-making, which also reflects the importance of communication for cooperation in specific attacks. It can be seen that communication is important in all tasks, with the order of importance being Treasure Hunt > Food Collector > Predator Prey.

AT is one of the most popular techniques in robust learning, and it is empirically confirmed to be a relatively reliable method. Compared with TARMAC, AT achieves lower timesteps under strong attacks, however, its baseline performance without attacks is slightly worse. This is inevitable as the trade-off between robustness and performance has long been a concern in adversarial training.

Though AME provides robustness against attacks to a certain degree, its baseline performance without attacks is bad, severely dragging down the overall performance. This results from the key feature of AME: adopting the consensus of all incoming messages. It protects the agents from being influenced by a small number of malicious messages, but it also prevents agents from getting exclusive information.

Our proposed framework ADMAC has the best overall performance, and this advantage is provided by the active defense strategy. Compared with AT and AME, ADMAC exploits the features of adversarial multi-agent communication, judging received messages' reliability based on agents' observations and hidden states. Notably, ADMAC has better baseline performance than TARMAC in Predator Prey tasks, which is caused by the feature of decomposable message aggregation policy and will be discussed more in the ablation study.

## Ablation Study

ADMAC consists of two components: the reliability estimator (abbreviated as RE) and the decomposable message aggregation policy net (abbreviated as DPN). To further investigate each component's function, we evaluate the following four methods under attacks: TARMAC, DPN only, DPN+RE

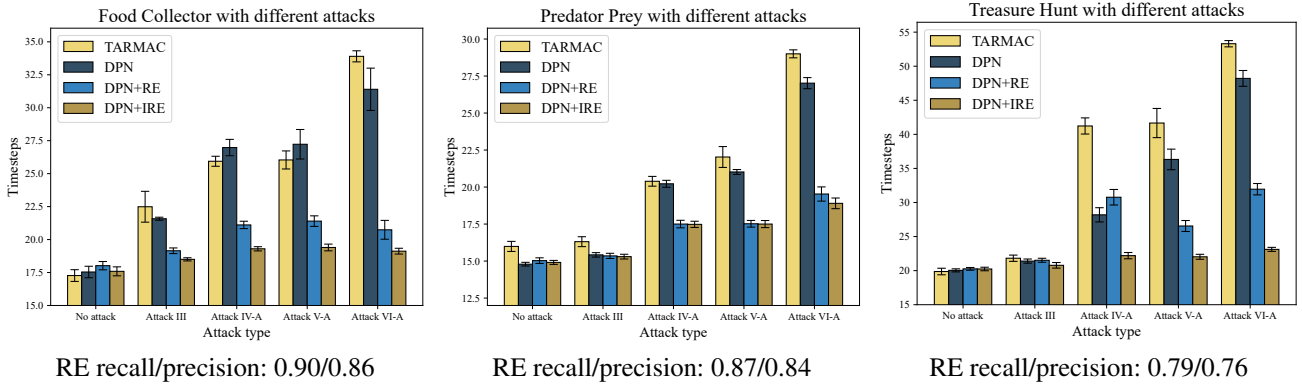RE recall/precision: 0.90/0.86      RE recall/precision: 0.87/0.84      RE recall/precision: 0.79/0.76

Figure 5: Results for the ablation study. The attack probability is set to $p = 0.3$ for all types of attack.

(i.e. ADMAC), and DPN+IRE. Here IRE refers to an ideal reliability estimator, achieving $100\%$ classification accuracy and cannot be implemented in reality. We also list the recall and precision of RE below the figures to analysis its performance from the perspective of supervised learning.

From the results presented in Fig 5, the following conclusions can be drawn. Firstly, the special structure of DPN does not degrades baseline performance, on the contrary, it provides a little robustness sometimes. Secondly, RE provides considerable robustness, especially against adversarial attacks. Thirdly, DPN+IRE has the best robust performance, which again verifies the effectiveness of the way we label messages. Additionally, the recall and precision of RE accounts for the performance gap between DPN+RE and DPN+IRE: The better the classification performance of RE is, the closer the performance of DPN+RE is to DPN+IRE.

## Related Work

TMC (Zhang, Zhang, and Lin 2020) is one of the earliest robust multi-agent communication frameworks. However, it only provides certain robustness against Gaussian noises and random message loss. Some recent researches have considered the existence of malicious attackers (Blumenkamp and Prorok 2021; Tu et al. 2021; Yuan et al. 2023; Sun et al. 2023), who perturb normal messages or send fake messages to disrupt normal operation of multi-agent systems. Agents within these frameworks commonly receive all messages equally and try to make robust decisions, which follows a passive defense idea conventional in robust RL (Havens, Jiang, and Sarkar 2018; Pattanaik et al. 2018). Notably, passive defense strategies fail to utilize one important feature of robust communicative MARL: unperturbed messages and hidden states can help identify fake messages to some extent. Unlike the aforementioned methods, R-MACRL (Xue et al. 2022) takes a rather proactive defense strategy, which is correcting perturbed messages. Since useful messages must contain some information unknown to the agents and attackers may replace the useful messages with manipulated ones, this strategy is not quite practical. Compared with them, while utilizing observations and hidden states, we have made corresponding preparations for the fact that it is impossible to perfectly judge received messages.

Besides, we would like to differentiate our research from robust MARL. Following Sun et al., we make two important assumptions: 1) only a part of the messages might be perturbed; 2) the perturbation power is unlimited and the perturbed messages can be completely different from the original ones. As a comparison, robust MARL commonly assumes the observations of agents (Li et al. 2019) or the environment model (Zhang et al. 2020b) is perturbed. Some robust MARL techniques can be used in robust communicative MARL (e.g. adversarial training), but an active defense strategy can better utilize the features of this problem.

## Conclusion

In this paper, we investigate the issue of robust communicative MARL, where malicious attackers may interrupt the multi-agent communication, and agents are required to make resilient decisions based on unperturbed observations and potentially perturbed messages. Considering that agents have multiple sources of information in this setting, we put forward an active defense strategy, which involves agents actively assessing the reliability of incoming messages and reducing the impact of potentially malicious messages to the final decisions.

Implementing active defense poses two primary challenges: labelling "unreliable" messages and adjusting the unreliable messages' impact on the final decision appropriately. To address them, we introduce ADMAC, a framework comprising an reliability estimator and a decomposable message aggregation policy net. ADMAC is evaluated alongside three alternative multi-agent communication frameworks in three communication-critical environments under attacks of different kinds and intensities, and shows outstanding robustness.

One limitation of ADMAC is that in scenarios where messages are likely to carry unique information, judging whether a message is reliable is hard, leading to a decline in the robustness of ADMAC. As future work, we will try to make agents aggregate information from a wider range of sources in order to better assess received messages, and expand our framework to accommodate scenarios with continuous action spaces.

# References

Blumenkamp, J.; and Prorok, A. 2021. The emergence of adversarial communication in multi-agent reinforcement learning. In *Conference on Robot Learning*, 1394–1414. PMLR.

Chu, T.; Chinchali, S.; and Katti, S. 2020. Multi-agent Reinforcement Learning for Networked System Control. In *International Conference on Learning Representations*. virtual conference.

Das, A.; Gervet, T.; Romoff, J.; Batra, D.; Parikh, D.; Rabbat, M.; and Pineau, J. 2019. Tarmac: Targeted multi-agent communication. In *International Conference on Machine Learning*, 1538–1546. PMLR.

Freed, B.; Sartoretti, G.; Hu, J.; and Choset, H. 2020. Communication learning via backpropagation in discrete channels with unknown noise. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, 7160–7168. New York.

Goodfellow, I. J.; Shlens, J.; and Szegedy, C. 2014. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*.

Han, B.; Yao, Q.; Yu, X.; Niu, G.; Xu, M.; Hu, W.; Tsang, I.; and Sugiyama, M. 2018. Co-teaching: Robust training of deep neural networks with extremely noisy labels. *Advances in neural information processing systems*, 31.

Han, R.; Chen, S.; and Hao, Q. 2020. Cooperative Multi-Robot Navigation in Dynamic Environment with Deep Reinforcement Learning. In *IEEE International Conference on Robotics and Automation*, 448–454.

Havens, A.; Jiang, Z.; and Sarkar, S. 2018. Online robust policy learning in the presence of unknown adversaries. *Advances in neural information processing systems*, 31.

He, S.; Han, S.; Su, S.; Han, S.; Zou, S.; and Miao, F. 2023. Robust Multi-Agent Reinforcement Learning with State Uncertainty. *Transactions on Machine Learning Research*.

Ishii, H.; Wang, Y.; and Feng, S. 2022. An overview on multi-agent consensus under adversarial attacks. *Annual Reviews in Control*.

Li, S.; Wu, Y.; Cui, X.; Dong, H.; Fang, F.; and Russell, S. 2019. Robust multi-agent reinforcement learning via minimax deep deterministic policy gradient. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, 4213–4220.

Littman, M. L. 1994. Markov games as a framework for multi-agent reinforcement learning. In *Proceedings of the Eleventh International Conference on International Conference on Machine Learning*, 157–163. New Brunswick.

Ma, Z.; Luo, Y.; and Pan, J. 2021. Learning selective communication for multi-agent path finding. *IEEE Robotics and Automation Letters*, 7(2): 1455–1462.

Madry, A.; Makelov, A.; Schmidt, L.; Tsipras, D.; and Vladu, A. 2018. Towards Deep Learning Models Resistant to Adversarial Attacks. In *International Conference on Learning Representations*.

Mu, R.; Ruan, W.; Marcolino, L. S.; Jin, G.; and Ni, Q. 2022. Certified Policy Smoothing for Cooperative Multi-Agent Reinforcement Learning. *arXiv preprint arXiv:2212.11746*.

Pattanaik, A.; Tang, Z.; Liu, S.; Bommannan, G.; and Chowdhary, G. 2018. Robust Deep Reinforcement Learning with Adversarial Attacks. In *Proceedings of the 17th International Conference on Autonomous Agents and MultiAgent Systems*, 2040–2042.

Pu, Z.; Wang, H.; Liu, Z.; Yi, J.; and Wu, S. 2022. Attention Enhanced Reinforcement Learning for Multi agent Cooperation. *IEEE Transactions on Neural Networks and Learning Systems*. To be published, doi:10.1109/TNNLS.2022.3146858.

Rangwala, M.; and Williams, R. 2020. Learning multi-agent communication through structured attentive reasoning. *Advances in Neural Information Processing Systems*, 33: 10088–10098.

Singh, A.; Jain, T.; and Sukhbaatar, S. 2018. Learning when to Communicate at Scale in Multiagent Cooperative and Competitive Tasks. In *International Conference on Learning Representations*. Vancouver.

Sun, Y.; Zheng, R.; Hassanzadeh, P.; Liang, Y.; Feizi, S.; Ganesh, S.; and Huang, F. 2023. Certifiably Robust Policy Learning against Adversarial Multi-Agent Communication. In *The Eleventh International Conference on Learning Representations*.

Tu, J.; Wang, T.; Wang, J.; Manivasagam, S.; Ren, M.; and Urtasun, R. 2021. Adversarial attacks on multi-agent communication. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 7768–7777.

Williams, R. J. 1992. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3): 229–256.

Xue, W.; Qiu, W.; An, B.; Rabinovich, Z.; Obraztsova, S.; and Yeo, C. K. 2022. Mis-spoke or mis-lead: Achieving Robustness in Multi-Agent Communicative Reinforcement Learning. In *Proceedings of the 21st International Conference on Autonomous Agents and Multiagent Systems*, 1418–1426.

Yuan, L.; Chen, F.; Zhang, Z.; and Yu, Y. 2023. Communication-Robust Multi-Agent Learning by Adaptable Auxiliary Multi-Agent Adversary Generation. *arXiv preprint arXiv:2305.05116*.

Zhang, H.; Chen, H.; Xiao, C.; Li, B.; Liu, M.; Boning, D.; and Hsieh, C.-J. 2020a. Robust deep reinforcement learning against adversarial perturbations on state observations. *Advances in Neural Information Processing Systems*, 33: 21024–21037.

Zhang, K.; Sun, T.; Tao, Y.; Genc, S.; Mallya, S.; and Basar, T. 2020b. Robust multi-agent reinforcement learning with model uncertainty. *Advances in neural information processing systems*, 33: 10571–10583.

Zhang, S. Q.; Zhang, Q.; and Lin, J. 2020. Succinct and robust multi-agent communication with temporal message control. *Advances in Neural Information Processing Systems*, 33: 17271–17282.

# A. Proof of Proposition 1

Proposition 1 can be formulated as follows:
If

$$v_i = f_{BP}(h_i) + \sum_{j \neq i,u} w_i(m_j) f_{MP}(o_i, m_j), \ \ p_i[k] = e^{v_i[k]}/\sum_k e^{v_i[k]},$$

$$\hat{v}_i = v_i + w f_{MP}(o_i, m_u), \ \ \hat{p}_i[k] = e^{\hat{v}_i[k]}/\sum_k e^{\hat{v}_i[k]}, \tag{10}$$

$$k_{max} = \arg\max_k f_{MP}(o_i, m_u)[k],$$

$$k_{min} = \arg\min_k f_{MP}(o_i, m_u)[k],$$

then

$$\partial \hat{p}_i[k_{max}]/\partial w > 0, \ \ \partial \hat{p}_i[k_{min}]/\partial w < 0 \tag{11}$$

*Proof*:
$f_{MP}(o_i, m_u)$ and $v_i$ both are vectors containing $K$ components, and we denote them as:

$$v_i = [a_1, a_2, ..., a_K]$$
$$f_{MP}(o_i, m_j) = [b_1, b_2, ..., b_K] \tag{12}$$

Then $\hat{p}_i[k_{max}]$ is calculated below:

$$\begin{aligned} \hat{p}_i[k_{max}] &= \frac{e^{a_{k_{max}}+wb_{k_{max}}}}{\sum_{k=1}^K e^{a_k+wb_k}} \\ &= \frac{1}{\sum_{k=1}^K e^{w(b_k-b_{k_{max}})}e^{a_k-a_{k_{max}}}} \\ &= \frac{1}{G_{max}(w)} \end{aligned} \tag{13}$$

Notably,

$$\partial G_{max}(w)/\partial w = \sum_{k=1}^K (b_k - b_{k_{max}})e^{w(b_k-b_{k_{max}})}e^{a_k-a_{k_{max}}} \tag{14}$$

According to (10), $b_k - b_{k_{max}} \leq 0$. Besides, $e^{w(b_k-b_{k_{max}})}e^{a_k-a_{k_{max}}} > 0$. Therefore, $\partial G_{max}(w)/\partial w < 0$, and we get $\partial \hat{p}_i[k_{max}]/\partial w > 0$. Similarly, $b_k - b_{k_{min}} \geq 0$ and $\partial \hat{p}_i[k_{min}]/\partial w < 0$.

# B. Implementation Details

## Training details

Our ADMAC framework consists of two key components, namely, a decomposable message aggregation policy net $f_P$ and a reliability estimator $f_R$. $f_P$ is trained using policy gradient methods specified below.

Use $\theta$ to denote the parameters of $f_P$, and the training goal is to maximize $J(\theta) = \mathbb{E}[\sum_t \gamma^t r_i^t | \theta]$. Use $\pi_\theta(h, o, m, a)$ to represent the probability of choosing action $a$ with hidden state $h$, observation $o$ and message $m$, $\nabla_\theta J(\theta)$ can be calculated using $N_e$ complete episodes consisting of $N_t$ transitions:

$$\nabla_\theta J(\theta) = \sum_{j=1}^{N_t} \pi_\theta(h_j, o_j, m_j, a_j) A(h_j, o_j, m_j, a_j) \nabla_\theta \log \pi_\theta(h_j, o_j, m_j, a_j), \tag{15}$$

where $h_j, o_j, m_j, a_j$ are from the $j$-th transition used for training and $A(h_j, o_j, m_j, a_j)$ is the advantage function calculated as follows:

$$A(h_j, o_j, m_j, a_j) = Q(h_j, o_j, m_j, a_j) - V_\phi(h_j, o_j, m_j), \tag{16}$$

where $Q(h_j, o_j, m_j, a_j)$ can be calculated using complete episodes in the replay buffer and $V_\phi(h_j, o_j, m_j)$ is approximated via neural networks with parameter $\phi$. Suppose the $j$-th transition is the $t_j$-th timestep of the episode it belongs to, then

$$Q(h_j, o_j, m_j, a_j) = \sum_{t=t_j}^{t_{max}} \gamma^{t-t_j} r^t. \tag{17}$$

Algorithm 1: Pseudo Code of Our Framework

---

**Input:** Settings for training policy net: training epochs $T_p$, batch size $t_p$ and optimizer A; Settings for training reliability estimator: dataset size $t_d$, perturb probabilities $p_a$ and $p_b$, training epochs $T_e$, batch size $N_b$ and optimizer B.

Randomly initialize the network parameters $\theta$, $\phi$.

// Stage I: training policy net

**for** $T = 1$ to $T_p$ **do**

    // Interacting with the environment

    **for** $t = 1$ to $t_p$ **do**

        Reset the environment and get observations $o_i^t$ for each agent if necessary;

        Each agent updates hidden states $h_i^t$, generates and broadcasts messages $m_i^t$, receives messages from others, and chooses actions $a_i^t$;

        For each agent, compute value function $v_i^t = V_\phi(h_i^t, o_i^t)$;

        Each agent executes action $a_i^t$ and receives reward $r_i^t$ as well as next observation $o_i^{t+1}$ from the environment;

        Store $(o_t^i, m_i^t, v_i^t, a_i^t, r_i^t, o_i^{t+1})$ into the replay buffer;

    **end for**

    // Updating neural networks

    For each transition in the replay buffer, compute $Q$ value according to (17);

    Calculate $L_V(\phi)$ and corresponding gradients $\nabla_\phi L_V(\phi)$ based on the transitions in the replay buffer according to (18);

    Calculate $\nabla_\theta J(\theta)$ according to (15);

    Update $\theta$ and $\phi$ using $\nabla_\theta J(\theta)$, $\nabla_\phi L_V(\phi)$, and optimizer A;

**end for**

// Stage II: generating dataset

**for** $t = 1$ to $t_d$ **do**

    Reset the environment and get observations $o_i^t$ for each agent if necessary;

    Each agent updates hidden states $h_i^t$, generates and broadcasts messages $m_i^t$, receives messages from others, and chooses actions $a_i^t$;

    Replace messages with probability $p_a$. If a message is replaced, then it has a probability $p_b$ to be replaced by a random message and a probability $1 - p_b$ to be replaced by an adversarial message. Then, all messages, including unperturbed and perturbed ones, are labeled according to their preference for the best action.

    Store $x_{ij} = [h_i, o_i, m_j]$ and $y_{ij}$ into a dataset $D$

    Each agent executes action $a_i^t$ and receives reward $r_i^t$ as well as next observation $o_i^{t+1}$ from the environment;

**end for**

// Stage III: training reliability estimator

**for** $T = 1$ to $T_e$ **do**

    Randomly divide dataset $D$ into batches of size $N_b$, and assume there are $t_e$ batches in total.

    **for** $t = 1$ to $t_e$ **do**

        Use the $t$-th batch to calculate $L_R(\psi)$ with (19)

        Update $\psi$ with $\nabla_\psi L_R(\psi)$ as well as optimizer B.

    **end for**

**end for**

---

$\phi$ is optimized by minimizing the value loss:

$$L_V(\phi) = \sum_{j=1}^{N_t} (Q(h_j, o_j, m_j, a_j) - V(h_j, o_j, m_j))^2. \tag{18}$$

$f_R$ with parameter $\psi$ is trained in a typical supervised learning way. Suppose there are $N_s$ samples in the dataset, and the $i$-th sample comprises $x_i = [h_i, o_i, m_i]$ and $y_i = 0$ or 1, where $y_i = 0$ represents the message is considered reliable and $y_i = 1$ represents the opposite. Then the loss is calculated as follows:

$$L_R(\psi) = -\sum_{i=1}^{N_s} (y_i \log f_R(h_i, o_i, m_i)[1] + (1 - y_i) \log f_R(h_i, o_i, m_i)[0]). \tag{19}$$

We present the pseudo code for our framework in Algorithm 1.

## Training settings

**Neural Network architectures** An agent uses a GRU module to obtain hidden states, and all other neural networks are fully connected with hidden layer size 128 and activation function ReLU. Specifically, observations are linearly transformed into

vectors of length 128 before being used. The base action generation module $f_{BP}$ and the message-observation process module $f_{MP}$ both have one hidden layer. The reliability estimator $f_R$ has two hidden layers.

**Training hyperparameters**   Training epochs $T_p = 2000$, $T_e = 100$. Batch sizes $t_p = 6000$, $N_b = 64$. The training dataset for the estimator contains around 300000 samples. Perturb probabilities $p_a = 0.5$, $p_b = 0.5$. We use Adam optimizer with learning rate 0.0003 during training.

## C. Additional Experimental Results

In our experiments, the implemented adversarial attacks have two attack goals:

$$f_A(m_j^t) = \arg\min \sum_{i \neq j} \hat{P}(a_{i,best}^t). \tag{20}$$

and

$$f_B(m_j^t) = \arg\max \sum_{i \neq j} D_{KL}(\hat{a}_i^t | a_i^t), \tag{21}$$

Since $f_B(m_j^t)$ is not applicable to AME, we only present type-A adversarial attacks in the main manuscript. Fig.6 presents the performance of models under type-B adversarial attacks. Consistent with the primary experimental findings, our approach exhibits excellent defensive performance.
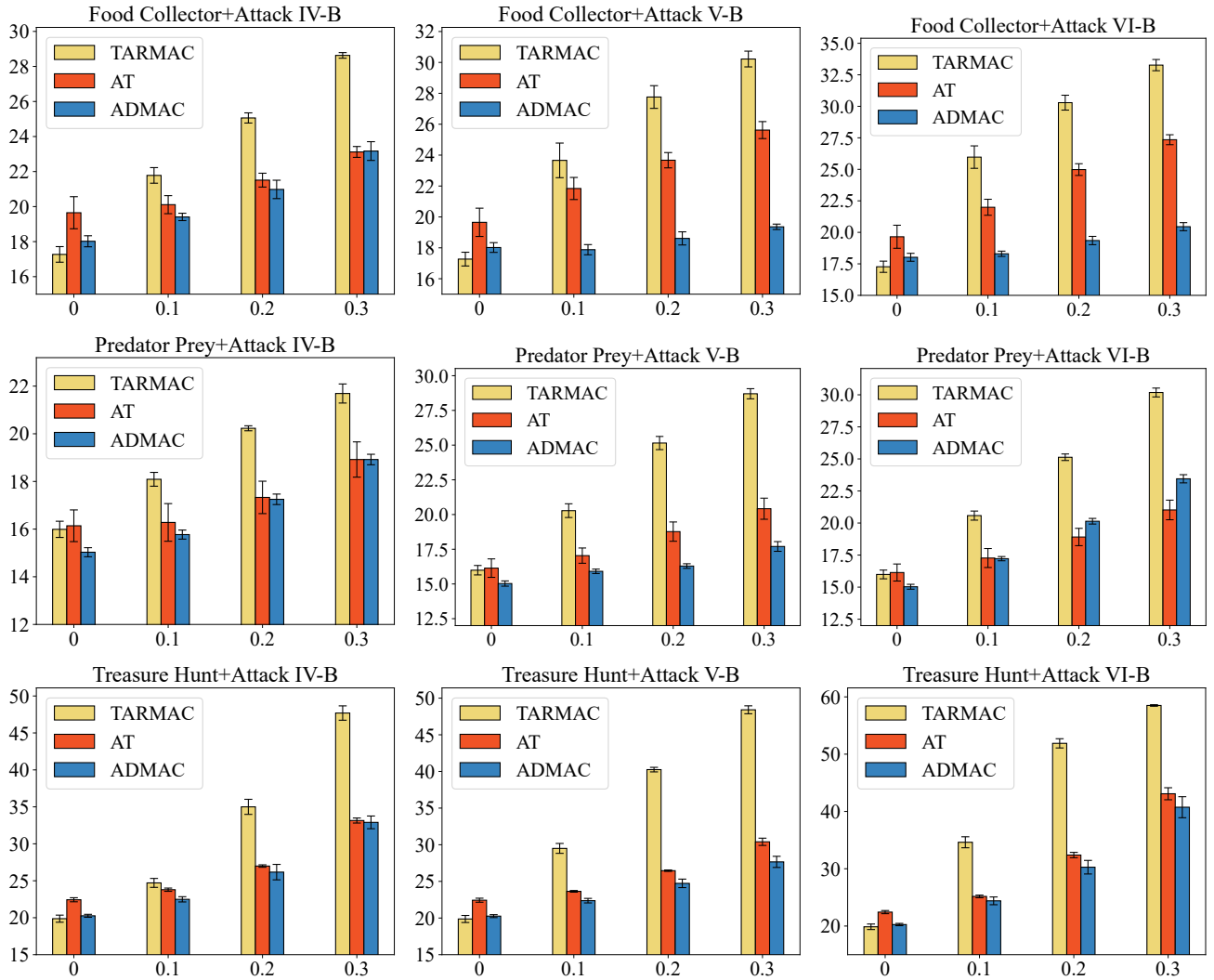


Figure 6: Performance of different models under adversarial attacks.