

# DPMAC: DIFFERENTIALLY PRIVATE COMMUNICATION FOR COOPERATIVE MULTI-AGENT REINFORCEMENT LEARNING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Communication lays the foundation for cooperation in human society and in multi-agent reinforcement learning (MARL). Humans also desire to maintain their privacy when communicating with others, yet such privacy concern has not been considered in existing works in MARL. We propose the *differentially private multi-agent communication* (DPMAC) algorithm, which protects the sensitive information of individual agents by equipping each agent with a local message sender with rigorous  $(\epsilon, \delta)$ -differential privacy (DP) guarantee. In contrast to directly perturbing the messages with predefined DP noise as commonly done in privacy-preserving scenarios, we adopt a stochastic message sender for each agent respectively and incorporate the DP requirement into the sender, which automatically adjusts the learned message distribution to alleviate the instability caused by DP noise. Further, we prove the existence of a Nash equilibrium in cooperative MARL with privacy-preserving communication, which suggests that this problem is game-theoretically learnable. Extensive experiments demonstrate a clear advantage of DPMAC over baseline methods in privacy-preserving scenarios.

## 1 INTRODUCTION

Multi-agent reinforcement learning (MARL) has shown remarkable achievements in many real-world applications such as sensor networks (Zhang & Lesser, 2011), autonomous driving (Shalev-Shwartz et al., 2016b), and traffic control (Wei et al., 2019). To mitigate non-stationarity when training the multi-agent system, centralized training and decentralized execution (CTDE) paradigm is proposed. The CTDE paradigm yet faces the hardness to enable complex cooperation and coordination for agents during execution due to the inherent partial observability in multi-agent scenarios (Wang et al., 2020b). To make agents cooperate more efficiently in complex partial observable environments, communication between agents has been considered. Numerous works proposed differentiable communication methods between agents, which can be trained in an end-to-end manner, for more efficient cooperation among agents (Foerster et al., 2016; Jiang & Lu, 2018; Das et al., 2019; Ding et al., 2020; Kim et al., 2021; Wang et al., 2020b). The communication can be either broadcast (Das et al., 2019; Jiang & Lu, 2018; Wang et al., 2020b), where the connection between agents can be modeled as a complete graph, or one-to-one as a general graph (Ding et al., 2020).

However, the advantages of communication, resulting from full information sharing, come with the possible privacy leakage of individual agents for both broadcasted and one-to-one messages. Therefore, in practice, one agent may be unwilling to fully share its private information with other agents even though in *cooperative* scenarios. For instance, if we train and deploy an MARL-based autonomous driving system, each autonomous vehicle involved in this system could be regarded as an agent and all vehicles work together to improve the safety and efficiency of the system. Hence, this can be regarded as a cooperative MARL scenario (Shalev-Shwartz et al., 2016a; Yang et al., 2020). However, owners of autonomous vehicles may not allow their vehicles to send private information to other vehicles without any desensitization since this may divulge their private information such as their personal life routines (Hassan et al., 2020). Hence, a natural question arises:

*Can the MARL algorithm with communication under the CTDE framework be endowed with both the rigorous privacy guarantee and the empirical efficiency?*

To answer this question, we start with a simple motivating example called *single round binary sums*, where several players attempt to guess the bits possessed by others and they can share their own information by communication. In Section 4, we show that a local message sender using the randomized response mechanism allows an analytical receiver to correctly calculate the binary sum in a privacy-preserving way. From the example we gain two insights: 1) The information is not supposed to be aggregated likewise in previous communication methods in MARL (Das et al., 2019; Ding et al., 2020), as a trusted data curator is not available in general. On the contrary, privacy is supposed to be achieved locally for every agent; 2) Once the agents know a priori, that certain privacy constraint exists, they could adjust their inference on the noised message. These two insights indicate the principles of our privacy-preserving communication structure that we desire a *privacy-preserving local sender* and a *privacy-aware analytical receiver*.

Our algorithm, *differentially private multi-agent communication* (DPMAC), instantiates the described principles. More specifically, for the sender part, each agent is equipped with a *local sender* which ensures differential privacy (DP) (Dwork, 2006) by performing an additive Gaussian noise. The message sender in DPMAC is local in the sense that each agent is equipped with its own message sender, which is only used to send its own messages. Equipped with this local sender, DPMAC is able to not only protect the privacy of communications between agents but also satisfy different privacy levels required from different agents. In addition, the sender adopts the Gaussian distribution to represent the message space and sample the stochastic message from the learned distribution. However, it is known that the DP noise may impede the original learning process (Dwork et al., 2014; Alvim et al., 2011), resulting in unstable or even divergent algorithms, especially for deep-learning-based methods (Abadi et al., 2016; Chen et al., 2020). To cope with this issue, we incorporate the noise variance into the representation of the message distribution, so that the agents could learn to adjust the message distribution automatically according to varying noise scales. For the receiver part, because of the gradient chain between the sender and the receiver, our receiver naturally utilizes the privacy-relevant information hidden in the gradients. This implements the privacy-aware analytical receiver described in the motivating example.

When protecting the privacy in communication is required in a cooperative game, the game is *not* purely cooperative anymore since each player involved will face a trade-off between the team utility and its personal privacy. To analyze the convergence of cooperative games with privacy-preserving communication, we first define a single-step game, namely the *collaborative game with privacy* (CGP). We prove that under some mild assumptions of the players’ value functions, CGP could be transformed into a potential game (Monderer & Shapley, 1996), subsequently leading to the existence of a Nash equilibrium (NE). With this property, NE could also be proved to exist in the single round binary sums game. Furthermore, we extend the single round binary sums into a multi-step game called *multiple round sums* using the notion of Markov potential game (MPG) (Leonardos et al., 2021). Inspired by Macua et al. (2018) and modeling the privacy-preserving communication as part of the agent action, we prove the existence of NE, which indicates that the multi-step game with privacy-preserving communication could be learnable.

To validate the effectiveness of DPMAC, extensive experiments are conducted in multi-agent particle environment (MPE) (Lowe et al., 2017), including cooperative navigation, cooperative communication and navigation, and predator-prey. Specifically, in privacy-preserving scenarios, DPMAC significantly outperforms baselines. Moreover, even without any privacy constraints, DPMAC could gain competitive performance against baselines.

To sum up, the contributions of this work are threefold:

- To the best of our knowledge, we make the first attempt to develop a framework for private communication in MARL, named *DPMAC*, with the theoretical guarantee of  $(\epsilon, \delta)$ -DP.
- We prove the existence of the Nash equilibrium for the cooperative games with privacy-preserving communication, which shows that these games are learnable.
- Experiments on the MPE show that DPMAC clearly outperforms other algorithms in privacy-preserving scenarios and gains competitive performance in non-private scenarios.

## 2 RELATED WORK

**Learning to communicate in MARL** Learning communication protocols in MARL by backpropagation and end-to-end training has achieved great advances in recent years (Sukhbaatar et al., 2016; Foerster et al., 2016; Jiang & Lu, 2018; Das et al., 2019; Wang et al., 2020b; Ding et al., 2020; Kim et al., 2021; Rangwala & Williams, 2020; Zhang et al., 2019; Singh et al., 2019; Zhang et al., 2020; 2021; Lin et al., 2021; Peng et al., 2017). Amongst these works, Sukhbaatar et al. (2016) propose CommNet as the first differentiable communication framework for MARL. Further, TarMAC (Das et al., 2019) and ATOC (Jiang & Lu, 2018) utilize the attention mechanism to extract useful information as messages. I2C (Ding et al., 2020) makes the first attempt to enable agents to learn one-to-one communication via causal inference. Wang et al. (2020b) propose NDQ, which learns nearly decomposable value functions to reduce the communication overhead. Kim et al. (2021) consider sharing an imagined trajectory as an intention for effectiveness. Besides, to communicate in the scenarios with limited bandwidth, some works consider learning to send compact and informative messages in MARL via minimizing the entropy of messages between agents using information bottleneck methods (Wang et al., 2020a; Tucker et al., 2022; Tian et al., 2021; Li et al., 2021). While learning effective communication in MARL has been extensively investigated, existing communication algorithms potentially leave the privacy of each agent vulnerable to information attacks.

**Privacy preserving in RL** With wide attention on reinforcement learning (RL) algorithms and applications in recent years, so have concerns about their privacy. Sakuma et al. (2008) consider privacy in the distributed RL problem and utilize cryptographic tools to protect the private state-action-state triples. Algorithmically, Balle et al. (2016) make the first attempt to establish a policy evaluation algorithm with differential privacy (DP) guarantee, where the Monte-Carlo estimates are perturbed with Gaussian noise. Wang & Hegde (2019) generalize the results to Q-learning, where functional noise is added to protect the reward function. Theoretically, Garcelon et al. (2021) analyze the regret bound of finite-horizon MDPs in the tabular case. In a large or continuous state space where function approximation is required, Zhou (2022) subsequently takes the first step to establish the sublinear regret in linear mixture Markov decision processes (MDPs). Zhao et al. (2022) propose the differentially private version of the temporal difference learning with nonlinear function approximation. Meanwhile, a large number of works focus on preserving privacy in multi-armed bandits (Tao et al., 2022; Tenenbaum et al., 2021; Dubey, 2021; Zheng et al., 2020; Dubey & Pentland, 2020; Tossou & Dimitrakakis, 2017).

Privacy is also studied in recent literature on MARL and multi-agent system. Ye et al. (2020) study differential advising for value-based agents, which share action values as the advice, largely differing in both the communication framework and the CTDE framework. Dong et al. (2020) propose an average consensus algorithm with a DP guarantee in the multi-agent system.

## 3 PRELIMINARIES

We consider a fully cooperative MARL problem where  $N$  agents work collaboratively to maximize the joint rewards. The underlying environment can be captured by a decentralized partially observable Markov decision process (Dec-POMDP), denoted by the tuple  $\langle \mathcal{S}, \mathcal{A}, \mathcal{O}, \mathcal{P}, \mathcal{R}, \gamma \rangle$ . Specifically,  $\mathcal{S}$  is the global state space,  $\mathcal{A} = \prod_{i=1}^N \mathcal{A}_i$  is the joint action space,  $\mathcal{O} = \prod_{i=1}^N \mathcal{O}_i$  is the joint observation space,  $\mathcal{P}(s' | s, \mathbf{a}) := \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$  determines the state transition dynamics,  $\mathcal{R}(s, \mathbf{a}) : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  is the reward function, and  $\gamma \in [0, 1]$  is the discount factor. Given a joint policy  $\pi = \{\pi_i\}_{i=1}^N$ , the joint action-value function at time  $t$  is  $Q^\pi(s^t, \mathbf{a}^t) = \mathbb{E}[G^t | s^t, \mathbf{a}^t, \pi]$ , where  $G^t = \sum_{i=0}^{\infty} \gamma^i \mathcal{R}^{t+i}$  is the cumulative reward, and  $\mathbf{a}^t = \{a_i^t\}_{i=1}^N$  is the joint action. The ultimate goal of the agents is to find an optimal policy  $\pi^*$  which maximizes  $Q^\pi(s^t, \mathbf{a}^t)$ .

Under the aforementioned cooperative setting, we study the case where agents are allowed to communicate with a joint message space  $\mathcal{M} = \prod_{i=1}^N \mathcal{M}_i$ . When the communication is unrestricted, the problem is reduced to a single-agent RL problem, which effectively solves the challenge posed by partially observable states, but puts the individual agent’s privacy at risk. To overcome the challenges of privacy and partial observable states simultaneously, we investigate algorithms that maximize the cumulative rewards while satisfying differential privacy (DP), given in Definition 3.1.

**Definition 3.1** ( $(\epsilon, \delta)$ -DP, [Dwork \(2006\)](#)). A randomized mechanism  $f : \mathcal{D} \rightarrow \mathcal{R}$  satisfies  $(\epsilon, \delta)$ -differential privacy if for any neighbouring datasets  $D, D' \in \mathcal{D}$  and  $S \subset \mathcal{R}$ , it holds that  $\Pr[f(D) \in S] \leq e^\epsilon \Pr[f(D') \in S] + \delta$ .

DP offers a mathematically rigorous way to quantify the privacy of an algorithm ([Dwork, 2006](#)). An algorithm is said to be “privatized” under the DP model if it is statistically hard to infer the presence of an individual data point in the dataset by observing the output of the algorithm. More intuitively, an algorithm satisfies DP if it provides nearly the same outputs given the neighbouring input datasets (i.e.,  $\Pr[f(D) \in S] \approx \Pr[f(D') \in S]$ ), which hence protects the sensitive information from the curious attacker.

With DP, each agent  $i$  is assigned with a privacy budget  $\epsilon_i$ , which is negatively correlated to the level of privacy protection. Then we have  $\epsilon = \{\epsilon_i\}_{i=1}^N$  as the set of all privacy budgets. In addition to maximizing the joint rewards as usually required in cooperative MARL, the messages sent from agent  $i$  are also required to satisfy the privacy budget  $\epsilon_i$  with probability at least  $1 - \delta$ .

## 4 MOTIVATING EXAMPLE

Before introducing our communication framework, we first investigate a motivating example, which is a *cooperative* game and inspires the design principles of private communication mechanisms in MARL. The motivating example is a simple yet interesting game, called *single round binary sums*. The game is extended from the example provided in [Cheu \(2021\)](#) for analyzing the shuffle model, while we illustrate the game from the perspective of multi-agent systems. We note that though this game is one-step, which is different from the sequential decision process like MDP, it is illustrative enough to show how the communication protocol works as a tool to achieve a better trade-off between privacy and utility.

Assume that there are  $N$  agents involved in this game. Each agent  $i \in [N]$  has a bit  $b_i \in [0, 1]$  and can tell other agents the information about its bit by communication. The objective of the game is for every agent to guess  $\sum_i b_i$ , the sum of the bits of all agents. Namely, each agent  $i$  makes a guess  $g_i$  and the utility of the agent is to maximize  $r_i = -|\sum_j b_j - \mathbb{E}[g_i]|$ . The (global) reward of this game is the sum of the utility over all agents, i.e.,  $\sum_i r_i$ .

Without loss of generality, we write the guess  $g_i$  into  $g_i = \sum_{j \neq i} y_{ij} + b_i$ , where  $y_{ij}$  is the guessed bit of agent  $j$  by agent  $i$ . If all agents share their bits without covering up, the guessed bit  $y_{ij}$  will obviously be equal to  $b_j$  and all agents attain an optimal return. Hence this game is fully cooperative under no privacy constraints. However, the optimal strategy is under the assumption that *everyone is altruistic to share their own bits*.

To preserve the privacy in communication, the message (i.e., the sent bit) could be randomized using *randomized response*, which perturbs the bit  $b_i$  with probability  $p$ , as shown below:

$$x_i = \mathcal{R}_{RR}(b_i) := \begin{cases} \text{Ber}(1/2) & \text{with probability } p \\ b_i & \text{otherwise} \end{cases},$$

where  $x_i$  is the random message and  $\text{Ber}$  indicates the Bernoulli distribution. Under our context,  $\mathcal{R}_{RR}$  is a *privacy-preserving message sender*, whose privacy guarantee is shown in Proposition 4.1.

**Proposition 4.1** ([Beimel et al. \(2008\)](#)). Setting  $p = \frac{2}{e^\epsilon + 1}$  in  $\mathcal{R}_{RR}$  suffices for  $(\epsilon, 0)$ -differential privacy.

When each agent is equipped with such a privacy-preserving sender  $\mathcal{R}_{RR}$  while adhering to the originally optimal strategy (i.e., believing what others tell and doing the guess), all agents would make an inaccurate guess. The bias of the guess denoted as  $\text{err}_i$  caused by  $\mathcal{R}_{RR}$  is then

$$\text{err}_i = \mathbb{E}[g_i] - \sum_i b_i = \sum_{j \neq i} \mathbb{E}[x_j - b_j] = p \sum_{j \neq i} \left(\frac{1}{2} - b_j\right) = \frac{p(N-1)}{2} - p \sum_{j \neq i} b_j.$$

Without any priori knowledge, the bias could not be reduced for  $(\epsilon, 0)$ -DP algorithms. However, if the probability  $p$  of perturbation is set as a prior common knowledge for all agents before the game

starts, the story will be different. One could transform the biased guess into

$$g_i^A = \mathcal{A}_{RR}(\vec{x}_{-i}) := \frac{1}{1-p} \left( \sum_{j \neq i} x_j - (N-1)p/2 \right),$$

where  $\vec{x}_{-i} = [x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_N]^\top$  denote the messages received by agent  $i$ . Then the estimate will be unbiased as

$$\mathbb{E}[g_i^A] = \frac{1}{1-p} \left( \mathbb{E} \left[ \sum_{j \neq i} x_j \right] - \frac{p(N-1)}{2} \right) + b_i = \sum_i b_i.$$

This example inspires that a communication algorithm could be both privacy-preserving and efficient. From the perspective of privacy, by the post-processing lemma of DP, any post-processing does not affect the original privacy level. From the perspective of utility, we could eliminate the bias  $\text{err}_i$  if the agent is equipped with the receiver  $\mathcal{A}_{RR}$  and the prior knowledge  $p$  is given.

In general, our motivating example gives two principles for designing privacy-preserving communication frameworks. First, to prevent the sensitive information from being inferred by other curious agents, we equip each agent with a local message sender with certain privacy constraints. Second, given a priori knowledge about the privacy requirement of other agents, the receiver could strategically analyze the received noisy messages to statistically reduce error due to the noisy communication. These two design principles correspond to two parts of our DPMAC framework respectively, *i.e.*, a *privacy-preserving local sender* and a *privacy-aware receiver*.

## 5 METHODOLOGY

Based on our design principles, we now introduce our DPMAC framework, as shown in Figure 1. Our framework is general and flexible, which makes it compatible with any CTDE method.

### 5.1 PRIVACY-PRESERVING LOCAL SENDER WITH STOCHASTIC GAUSSIAN MESSAGES

In this section, we present the sender’s perspective on the privacy guarantee. At time  $t$ , for agent  $i$ , a message function  $f_i^s$  is used to generate a message for communication.  $f_i^s$  takes a subset of transitions in local trajectory  $\tau_i^t$  as input, where the subset is sampled uniformly without replacement from  $\tau_i^t$  (denote the sampling rate as  $\gamma_1$ ). This message is perturbed by the Gaussian mechanism with variance  $\sigma_i^2$  (Dwork, 2006). Agent  $i$  then samples a subset of other agents to share this message (denote the sampling rate as  $\gamma_2$ ). The following theorem guarantees differential privacy.

**Theorem 5.1** (Privacy guarantee for DPMAC). *Let  $\gamma_1, \gamma_2 \in (0, 1)$ , and  $C$  be the  $\ell_2$  norm of the message functions. For any  $\delta > 0$  and privacy budget  $\epsilon_i$ , the communication of agent  $i$  satisfies  $(\epsilon_i, \delta)$ -DP when  $\sigma_i^2 = \frac{14\gamma_2\gamma_1^2NC^2\alpha}{\beta\epsilon_i}$ , if we have  $\alpha = \frac{\log \delta^{-1}}{\epsilon_i(1-\beta)} + 1 \leq 2\sigma'^2 \log(1/\gamma_1\alpha(1+\sigma'^2)) / 3 + 1$  with  $\beta \in (0, 1)$  and  $\sigma'^2 = \sigma_i^2 / (4C^2) \geq 0.7$ .*

With Theorem 5.1, one can directly translate a non-private MARL with a communication algorithm into a private one. However, as we shall see in our experiment section, directly injecting the privacy noise into existing MARL with communication algorithms may lead to serious performance degradation. In fact, the injected noise might jeopardize the useful information incorporated in the messages, or even leads to meaningless messages. To alleviate the negative impacts of the injected privacy noise on the cooperation between agents, we adopt a stochastic message sender in the sense that the messages sent by our sender are sampled from a learned message distribution. This makes DPMAC different from existing works in MARL that communicate through deterministic messages (Sukhbaatar et al., 2016; Foerster et al., 2016; Jiang & Lu, 2018; Das et al., 2019; Ding et al., 2020; Kim et al., 2021).

In the following, we drop the dependency of parameters on  $t$  when it is clear from the context. Without loss of generality, let the message distribution be multivariate Gaussian and let  $p_i$  be the message sampled from the message distribution  $\mathcal{N}(\mu_i, \Sigma_i)$ , where  $\mu_i = f_i^\mu(o_i, a_i; \theta_i^\mu)$  and  $\Sigma_i = f_i^\sigma(o_i, a_i; \theta_i^\sigma)$  are the mean vector and covariance matrix learned by the sender, and  $\theta_i^\mu$  and  $\theta_i^\sigma$  are the



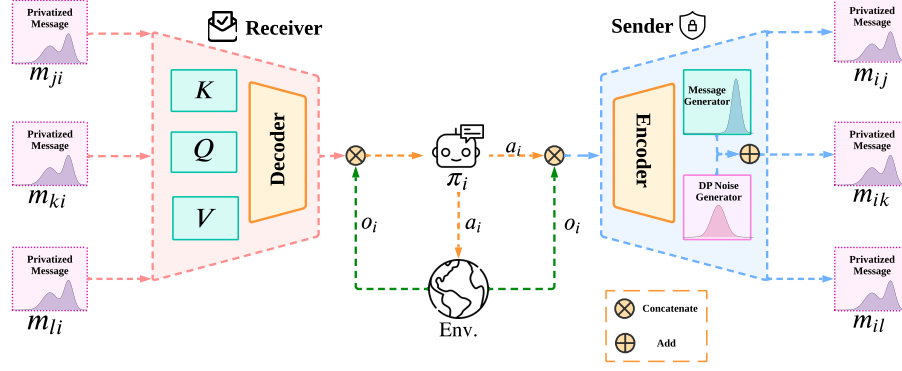


Figure 1: The overall structure of DPMAC. The message receiver of agent  $i$  integrates other agents' messages  $\{m_{ji}, m_{ki}, m_{li}\}$  with the self-attention mechanism and the integrated message is fed into the policy  $\pi_i$  together with the observation  $o_i$ . Agent  $i$  interacts with the environment by taking action  $a_i$ . Then  $o_i$  and  $a_i$  are concatenated and encoded by a privacy-preserving message sender and sent to other agents.

parameters of the sender's neural networks. Then  $\theta_i^{\mu\top}$  and  $\theta_i^{\sigma\top}$  will be optimized towards making all the agents to send more effective messages to encourage better team cooperation and gain higher team rewards. For notational convenience, let  $\theta_i^s = [\theta_i^{\mu\top}, \theta_i^{\sigma\top}]^\top$ . Then the sent privatized message  $m_i = p_i + u_i$  where  $u_i \sim \mathcal{N}(0, \sigma_i^2 \mathbf{I}_d)$  is an additional noise. It is clear that  $m_i \sim \mathcal{N}(\mu_i, \Sigma_i + \sigma_i^2 \mathbf{I}_d)$  since  $p_i$  is independent from  $u_i$ . Counterfactually, let  $m'_i \sim \mathcal{N}(\mu'_i, \Sigma'_i)$ , where  $\mu'_i = f_i^\mu(o_i, a_i; \theta_i^{\mu'})$  and  $\Sigma'_i = f_i^\sigma(o_i, a_i; \theta_i^{\sigma'})$  is the sent message when it was not under any privacy constraint.

Let the optimal message distribution be  $\mathcal{N}(\mu_i^*, \Sigma_i^*)$ . We are interested to characterize  $\theta_i^{s'}$  and  $\theta_i^s$ . By the optimality of  $\mu_i^*, \Sigma_i^*$ ,

$$\begin{aligned} \theta_i^{s'} &= \arg \min_{\theta} D_{\text{KL}}(\mathcal{N}(\mu'_i, \Sigma'_i) \| \mathcal{N}(\mu_i^*, \Sigma_i^*)) \\ &= \arg \min_{\theta} \log \frac{|\Sigma_i^*|}{|\Sigma'_i|} + \text{tr}\{\Sigma_i^{*-1} \Sigma'_i\} + \|\mu'_i - \mu_i^*\|_{\Sigma_i^{*-1}}^2. \end{aligned} \quad (1)$$

Then under the privacy constraints, the stochastic sender will learn  $\theta_i^s$  such that

$$\begin{aligned} \theta_i^s &= \arg \min_{\theta} D_{\text{KL}}(\mathcal{N}(\mu_i, \Sigma_i + \sigma_i^2 \mathbf{I}_d) \| \mathcal{N}(\mu_i^*, \Sigma_i^*)) \\ &= \arg \min_{\theta} \log \frac{|\Sigma_i^*|}{|\Sigma_i + \sigma_i^2 \mathbf{I}_d|} + \text{tr}\{\Sigma_i^{*-1} (\Sigma_i + \sigma_i^2 \mathbf{I}_d)\} + \|\mu_i - \mu_i^*\|_{\Sigma_i^{*-1}}^2. \end{aligned} \quad (2)$$

Through Equation (2), it is possible to directly incorporate the distribution of privacy noise into the optimization process of the sender to help to learn  $\theta_i^s$  such that  $D_{\text{KL}}(\mathcal{N}(\mu_i, \Sigma_i + \sigma_i^2 \mathbf{I}_d) \| \mathcal{N}(\mu_i^*, \Sigma_i^*)) \leq D_{\text{KL}}(\mathcal{N}(\mu'_i, \Sigma'_i) \| \mathcal{N}(\mu_i^*, \Sigma_i^*))$ , which means that the sender could learn to send private message  $m_i = p_i + u_i$  that is at least as effective as the non-private message  $m'_i$ . In this manner, the performance degradation is expected to be well alleviated.

## 5.2 PRIVACY-AWARE MESSAGE RECEIVER

As shown in our motivating example, the message receiver with knowledge a priori could statistically reduce the communication error in privacy-preserving scenarios. In the practical design, this motivation could be naturally instantiated with the gradient flow between the message sender and the message receiver.

Specifically, agent  $i$  first concatenates all the received privatized messages as  $\mathbf{m}_{(-i)i} := \{m_{ji}\}_{j=1, j \neq i}^N$  and then encodes  $\mathbf{m}_{(-i)i}$  into an aggregated message  $q_i = f_i^r(\mathbf{m}_{(-i)i} | \theta_i^r)$  with the decoding function  $f_i^r$  parameterized by  $\theta_i^r$ . Then a similar argument to the policy gradient theorem (Sutton et al., 1999) states that the gradient of the receiver is

$$\nabla_{\theta_i^r} \mathcal{J}(\theta_i^r) = \mathbb{E}_{\tau, \mathbf{o}, \mathbf{a}} [\mathbb{E}_{\pi_i} [\nabla_{\theta_i^r} f_i^r(q_i | \mathbf{m}_{(-i)i}) \nabla_{q_i} \log \pi_i(a_i | o_i, q_i) Q^\pi(\mathbf{a}, \mathbf{o})]],$$

where  $\mathcal{J}(\theta_i^r) = \mathbb{E}[G^1 \mid \pi]$  is the cumulative discounted reward from the starting state. In this way, the receiver could utilize the prior knowledge  $\sigma_i$  of the privacy-preserving sender encoded in the gradient during the optimization process. Please refer to Appendix D for the detailed optimization process of the message senders and receivers.

## 6 PRIVACY-PRESERVING EQUILIBRIUM ANALYSIS

Many cooperative multi-agent games enjoy the existence of a unique NE, which ensures the convergence of iterative algorithms. Under the privacy constraints, however, the existence of a unique Nash equilibrium can no longer be guaranteed even if the original game admits a unique equilibrium. As the convergence of MARL algorithms could depend on the existence of an equilibrium, we investigate such existence in single-step games and extend the result to multi-step games.

### 6.1 SINGLE-STEP GAMES

We study a class of two-player collaborative games, denoted as *collaborative game with privacy (CGP)*. The game involves two agents, each equipped with a privacy parameter  $p_n, n \in \{1, 2\}$ . The value of  $p_n$  represents the importance of privacy to agent  $n$ , with the larger value referring to greater importance. Let  $\mathcal{M}$  be some message mechanism. We denote the privacy loss by  $c^{\mathcal{M}}(p_n)$ , which measures the quantity of the potential privacy leakage and is formally defined in Definition B.2. Besides, let  $b(V_n, V_n^{\mathcal{M}}(p_1, p_2))$  be the utility gained by measuring the gap between private value function  $V_n^{\mathcal{M}}(p_1, p_2)$  and non-private value function  $V_n$ . Then the trade-off between the utility and the privacy is depicted by the total utility function  $u_n(p_1, p_2)$  in Equation (3). The formal definition of CGP is given in Definition 6.1. See more details in Appendix B.1.

**Definition 6.1** (Collaborative game with privacy (CGP)). *The collaborative game with privacy is denoted by a tuple  $\langle \mathcal{N}, \Sigma, \mathcal{U} \rangle$ , where  $\mathcal{N} = \{1, 2\}$  is the set of players,  $\Sigma = \{p_1, p_2\}$  is the action set with  $p_1, p_2 \in [0, 1]$  representing the privacy level, and  $\mathcal{U} = \{u_1, u_2\}$  is the set of utility functions satisfying  $\forall n \in \mathcal{N}$ ,*

$$u_n(p_1, p_2) = B_n \cdot b(V_n, V_n^{\mathcal{M}}(p_1, p_2)) - C_n^{\mathcal{M}} \cdot c^{\mathcal{M}}(p_n). \quad (3)$$

Then the following theorem shows that if changes in the value function of each player can be expressed as a change in their own privacy parameter, then CGP is a potential game and a pure NE thereafter exists. The proof is deferred to Appendix B.1.

**Theorem 6.1** (CGP’s NE guarantee). *The collaborative game with privacy has at least one non-trivial pure-strategy Nash equilibrium if  $\partial_{p_1}^i V_1 = \partial_{p_2}^i V_2, \forall i \in \{1, 2\}$ .*

**Equilibrium in single round binary sums** Let us revisit our motivating example. Armed with the CGP framework, it is immediate that *the single round binary sums game guarantees the existence of a NE*. This result is formalized in Theorem B.2 in Appendix B.1.

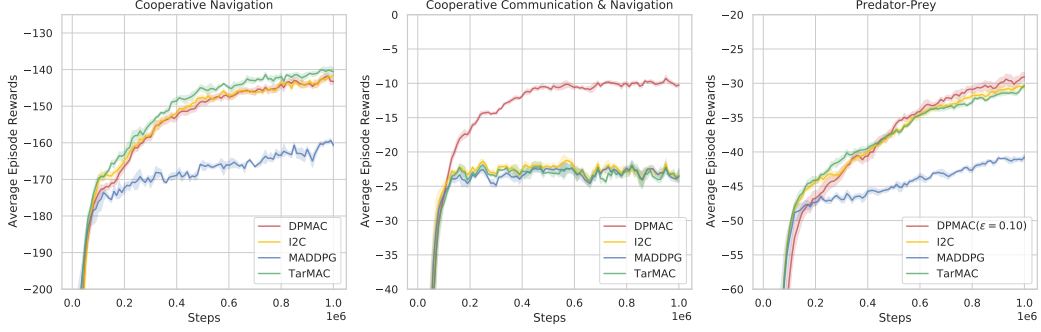
### 6.2 MULTI-STEP GAMES

We now consider an extended version of single round binary sums named *multiple round sums*. Consider an  $N$ -player game where player  $i$  owns a saving  $x_{i,t}$ . Rather than sending a binary bit, the agent can choose to give out  $b_{i,t}$  at round  $t$ . Meanwhile, each player  $i$  selects privacy level  $p_{i,t}$  and sends messages to each other with a sender  $f_i^s$  encoding the information of  $b_{i,t}$  with the privacy level  $p_{i,t}$ . The reward of the agent is designed to find a good trade-off between privacy and utility. The setting of the game is thus similar to the empirical implementation of DPMAC.

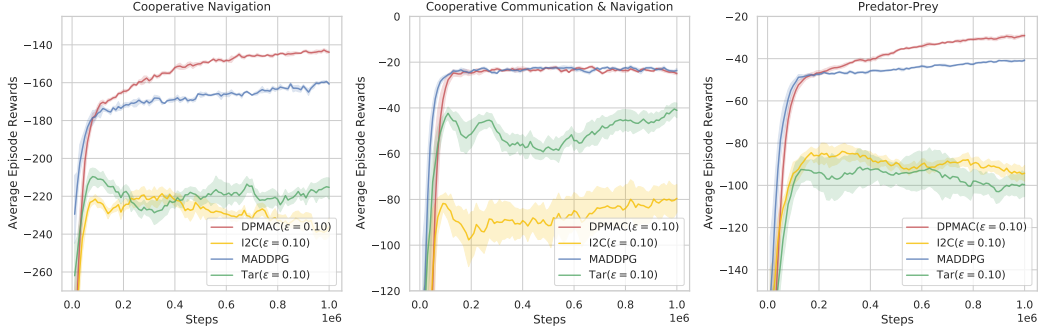
We first transform this game into a Markov potential game (MPG), with the reward of each agent transformed into a combination of the team reward and the individual reward. Then with existing theoretical results from Macua et al. (2018), we present the following result while deferring its proof to Appendix B.2.

**Theorem 6.2** (NE guarantee in multiple round sums). *If Assumptions 1, 2, 3, 4 (see Appendix B.2) are satisfied, our MPG has a NE with potential function  $J$  defined as,*

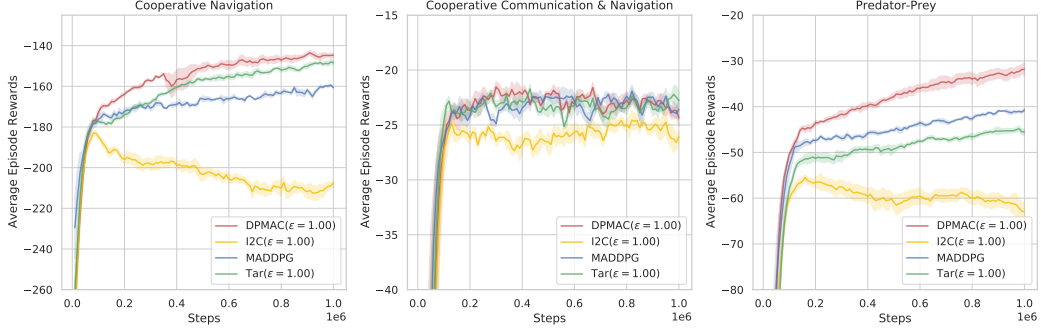
$$J(x_t, \pi(x_t)) = \sum_{j \in [N]} ((1 - p_{j,t})b_{j,t} + \alpha x_{j,t} + \beta p_{i,t}). \quad (4)$$



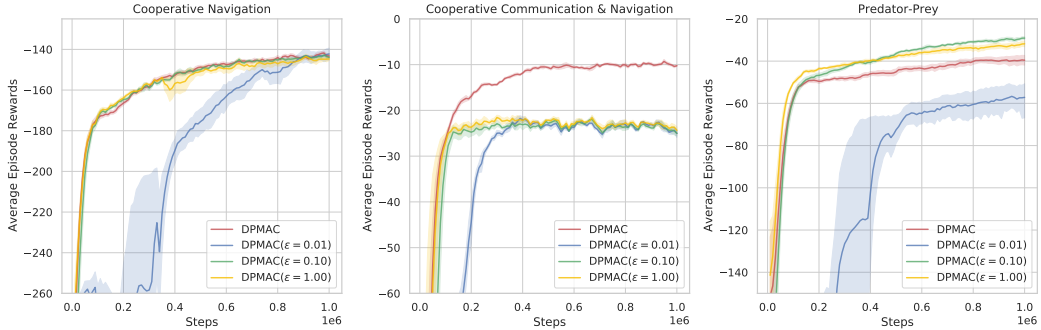
(a) Learning curves of DPMAC, TarMAC, I2C, and MADDPG on three MPE tasks. Note that on the PP task DPMAC ( $\epsilon = 0.10$ ) is shown.



(b) Learning curves of different algorithms under the privacy budget  $\epsilon = 0.10$ . MADDPG (non-private) is also displayed for comparison.



(c) Learning curves of different algorithms under the privacy budget  $\epsilon = 1.0$ . MADDPG (non-private) is also displayed for comparison.



(d) Learning curves of different privacy budgets ( $\epsilon = 0.01, 0.10, 1.00$ ) for DPMAC.

Figure 2: Learning curves of DPMAC and baseline algorithms. The curves are averaged over 5 seeds. Shaded areas denote 1 standard deviation.



## 7 EXPERIMENTS

In this section, we present the experiment results and corresponding experiment analyses. Please see Appendix G for more detailed analyses of experiment results.

**Baselines** We implement our DPMAC and evaluate it against TarMAC (Das et al., 2019), I2C (Ding et al., 2020), and MADDPG (Lowe et al., 2017). All Algorithms are tested with and without the privacy requirement except for MADDPG, which involves no communication among agents. Since TarMAC and I2C do not have a local sender and have no DP guarantee, we add Gaussian noise to their receiver according to the noise variance specified in Theorem 5.1 for a fair comparison. Please see Appendix D for more training details. *We remark that the code will be made publicly available once this manuscript is accepted.*

**Environments** We evaluate the algorithms on the multi-agent particle environment (MPE) (Mordatch & Abbeel, 2017), which is with continuous observation and discrete action space. This environment is commonly used among existing literature (Lowe et al., 2017; Jiang & Lu, 2018; Ding et al., 2020; Kim et al., 2021). We evaluate a wide range of tasks in MPE, including cooperative navigation (CN), cooperative communication and navigation (CCN), and predator prey (PP). More details on the environmental settings are given in Appendix E.

**Experiment results without privacy** DPMAC is first compared with TarMAC, I2C, and MADDPG on three MPE tasks without the privacy requirement, as shown in Figure 2a. DPMAC outperforms baselines on CCN & PP and has competitive performance on CN. Note that for the PP task we pick DPMAC with  $\epsilon = 0.10$  due to even better performance over its non-private variant. The comparison between DPMAC (non-private) and baselines is provided in Appendix F.

**Experiment results with privacy** We further add the privacy constraint on the communication algorithms. We set  $\delta = 10^{-4}$  on all tasks. Figure 2b and Figure 2c show the performance under the privacy budget  $\epsilon = 0.10$ ,  $\epsilon = 1.0$  and both with  $\delta = 10^{-4}$ . We include MADDPG as a non-communication baseline method. We observe that DPMAC with the privacy requirement could still maintain a good result compared to MADDPG, while the performance of TarMAC and I2C drops greatly. Figure 2d further gives the comparison between the performance of DPMAC under different privacy budgets. When  $\epsilon = 0.01$ , DPMAC still gains remarkable performance, while other baselines’ performance degraded greatly, as shown in Figure 2b.

**Variance adjustment of DPMAC** Experiments with privacy also support our claim that DPMAC could automatically adjust the variance of our stochastic message sender so that it learns a noise-robust representation. As shown in Figure 2d, DPMAC gains very close performance when  $\epsilon = 0.1$  and  $\epsilon = 1.0$ , though the privacy requirements of  $\epsilon = 0.1$  and  $\epsilon = 1.0$  differ by one order of magnitude. However, one can see large gaps for the same baseline algorithms under different  $\epsilon$  from Figure 2b and Figure 2c. Please see Figure 4 and Figure 5 for direct presentations of these gaps.

## 8 CONCLUSION

In this paper, we study the privacy-preserving communication in MARL. Motivated by a simple yet effective example of the binary sums game, we propose DPMAC, a new efficient communicating MARL algorithm that preserves agents’ privacy through differential privacy. Our algorithm is justified both theoretically and empirically. Besides, to show that the privacy-preserving communication problem is learnable, we analyze the single-step game and the multi-step game via the notion of Markov potential games (MPG) and show the existence of the Nash equilibrium. This existence further implies the learnability of several instances of MPG under privacy constraints. Extensive experiments are conducted on MPE and show the effectiveness of DPMAC when compared to baseline methods on multiple tasks both with and without the privacy constraints.

Though we make the first step to establish an efficient MARL algorithm with differential private communication, some interesting questions remain open. The first question is that it is still unclear for us whether there exists the Nash equilibrium in private competitive games. Besides, on the empirical side, investigating the performance of DPMAC in competitive games with privacy-preserving communication might also be interesting and valuable.

## REFERENCES

- Martin Abadi, Andy Chu, Ian Goodfellow, H Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC conference on computer and communications security*, pp. 308–318, 2016. 2
- Mário S Alvim, Miguel E Andrés, Konstantinos Chatzikokolakis, Pierpaolo Degano, and Catuscia Palamidessi. Differential privacy: on the trade-off between utility and information leakage. In *International Workshop on Formal Aspects in Security and Trust*, pp. 39–54. Springer, 2011. 2
- Borja Balle, Maziar Gomrokchi, and Doina Precup. Differentially private policy evaluation. In *International Conference on Machine Learning*, pp. 2130–2138. PMLR, 2016. 3
- Amos Beimel, Kobbi Nissim, and Eran Omri. Distributed private data analysis: Simultaneously solving how and what. In David A. Wagner (ed.), *Advances in Cryptology - CRYPTO 2008, 28th Annual International Cryptology Conference, Santa Barbara, CA, USA, August 17-21, 2008. Proceedings*, volume 5157 of *Lecture Notes in Computer Science*, pp. 451–468. Springer, 2008. 4
- Dingfan Chen, Tribhuvanesh Orekondy, and Mario Fritz. GS-WGAN: A gradient-sanitized approach for learning differentially private generators. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020. 2
- Albert Cheu. Differential privacy in the shuffle model: A survey of separations. *arXiv preprint arXiv:2107.11839*, 2021. 4
- Abhishek Das, Théophile Gervet, Joshua Romoff, Dhruv Batra, Devi Parikh, Mike Rabbat, and Joelle Pineau. Tarmac: Targeted multi-agent communication. In *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pp. 1538–1546. PMLR, 2019. 1, 2, 3, 5, 9
- Ziluo Ding, Tiejun Huang, and Zongqing Lu. Learning individually inferred communication for multi-agent cooperation. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020. 1, 2, 3, 5, 9, 19
- Tao Dong, Xiangyu Bu, and Wenjie Hu. Distributed differentially private average consensus for multi-agent networks by additive functional laplace noise. *Journal of the Franklin Institute*, 357(6):3565–3584, 2020. 3
- Abhimanyu Dubey. No-regret algorithms for private gaussian process bandit optimization. In *The 24th International Conference on Artificial Intelligence and Statistics, AISTATS 2021, April 13-15, 2021, Virtual Event*, volume 130 of *Proceedings of Machine Learning Research*, pp. 2062–2070. PMLR, 2021. 3
- Abhimanyu Dubey and Alex ‘Sandy’ Pentland. Differentially-private federated linear bandits. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020. 3
- Cynthia Dwork. Differential privacy. In *International Colloquium on Automata, Languages, and Programming*, pp. 1–12. Springer, 2006. 2, 4, 5
- Cynthia Dwork, Aaron Roth, et al. The algorithmic foundations of differential privacy. *Found. Trends Theor. Comput. Sci.*, 9(3-4):211–407, 2014. 2, 14
- Jakob N. Foerster, Yannis M. Assael, Nando de Freitas, and Shimon Whiteson. Learning to communicate with deep multi-agent reinforcement learning. In *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, pp. 2137–2145, 2016. 1, 3, 5

- Evrard Garcelon, Vianney Perchet, Ciara Pike-Burke, and Matteo Pirodda. Local differential privacy for regret minimization in reinforcement learning. In *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pp. 10561–10573, 2021. [3](#)
- Muneeb Ul Hassan, Mubashir Husain Rehmani, and Jinjun Chen. Differential privacy techniques for cyber physical systems: A survey. *IEEE Commun. Surv. Tutorials*, 22(1):746–789, 2020. [1](#)
- Matthew Hausknecht and Peter Stone. Deep recurrent q-learning for partially observable mdps. In *2015 aaai fall symposium series*, 2015. [18](#)
- Jiechuan Jiang and Zongqing Lu. Learning attentional communication for multi-agent cooperation. In *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pp. 7265–7275, 2018. [1](#), [3](#), [5](#), [9](#)
- Woojun Kim, Jongeui Park, and Youngchul Sung. Communication in multi-agent reinforcement learning: Intention sharing. In *International Conference on Learning Representations*, 2021. [1](#), [3](#), [5](#), [9](#), [19](#)
- Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. In Yoshua Bengio and Yann LeCun (eds.), *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*, 2014. [18](#)
- Valli Kumari and Srinivasa Chakravarthy. Cooperative privacy game: a novel strategy for preserving privacy in data publishing. *Human-centric Computing and Information Sciences*, 6(1):1–20, 2016. [17](#)
- Stefanos Leonardos, Will Overman, Ioannis Panageas, and Georgios Piliouras. Global convergence of multi-agent policy gradient in markov potential games. *arXiv preprint arXiv:2106.01969*, 2021. [2](#), [15](#), [17](#)
- Ziyan Li, Quan Yuan, Guiyang Luo, and Jinglin Li. Learning effective multi-vehicle cooperation at unsignalized intersection via bandwidth-constrained communication. In *94th IEEE Vehicular Technology Conference, VTC Fall 2021, Norman, OK, USA, September 27-30, 2021*, pp. 1–7. IEEE, 2021. doi: 10.1109/VTC2021-Fall52928.2021.9625057. [3](#)
- Toru Lin, Jacob Huh, Christopher Stauffer, Ser-Nam Lim, and Phillip Isola. Learning to ground multi-agent communication with autoencoders. In *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pp. 15230–15242, 2021. [3](#)
- Ryan Lowe, Yi Wu, Aviv Tamar, Jean Harb, Pieter Abbeel, and Igor Mordatch. Multi-agent actor-critic for mixed cooperative-competitive environments. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pp. 6379–6390, 2017. [2](#), [9](#), [18](#), [19](#)
- Sergio Valcarcel Macua, Javier Zazo, and Santiago Zazo. Learning parametric closed-loop policies for markov potential games. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018. [2](#), [7](#), [15](#), [17](#)
- Ilya Mironov. Rényi differential privacy. In *2017 IEEE 30th Computer Security Foundations Symposium (CSF)*, pp. 263–275. IEEE, 2017. [14](#)
- Dov Monderer and Lloyd S Shapley. Potential games. *Games and economic behavior*, 14(1):124–143, 1996. [2](#), [15](#), [16](#)
- Igor Mordatch and Pieter Abbeel. Emergence of grounded compositional language in multi-agent populations. *arXiv preprint arXiv:1703.04908*, 2017. [9](#), [19](#)
- John F Nash et al. Equilibrium points in n-person games. *Proceedings of the national academy of sciences*, 36(1):48–49, 1950. [15](#), [17](#)

- Balazs Pejo, Qiang Tang, and Gergely Biczók. Together or alone: The price of privacy in collaborative learning. *Proc. Priv. Enhancing Technol.*, 2019(2):47–65, 2019. 15, 16, 17
- Peng Peng, Quan Yuan, Ying Wen, Yaodong Yang, Zhenkun Tang, Haitao Long, and Jun Wang. Multiagent bidirectionally-coordinated nets for learning to play starcraft combat games. *CoRR*, abs/1703.10069, 2017. 3
- Murtaza Rangwala and Ryan Williams. Learning multi-agent communication through structured attentive reasoning. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (eds.), *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020. 3
- Jun Sakuma, Shigenobu Kobayashi, and Rebecca N Wright. Privacy-preserving reinforcement learning. In *Proceedings of the 25th international conference on Machine learning*, pp. 864–871, 2008. 3
- Shai Shalev-Shwartz, Shaked Shammah, and Amnon Shashua. Safe, multi-agent, reinforcement learning for autonomous driving. *CoRR*, abs/1610.03295, 2016a. URL <http://arxiv.org/abs/1610.03295>. 1
- Shai Shalev-Shwartz, Shaked Shammah, and Amnon Shashua. Safe, multi-agent, reinforcement learning for autonomous driving. *CoRR*, abs/1610.03295, 2016b. 1
- Amanpreet Singh, Tushar Jain, and Sainbayar Sukhbaatar. Learning when to communicate at scale in multiagent cooperative and competitive tasks. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019. 3
- Sainbayar Sukhbaatar, Arthur Szlam, and Rob Fergus. Learning multiagent communication with backpropagation. In Daniel D. Lee, Masashi Sugiyama, Ulrike von Luxburg, Isabelle Guyon, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, pp. 2244–2252, 2016. 3, 5
- Richard S. Sutton, David A. McAllester, Satinder P. Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. In *Advances in Neural Information Processing Systems 12, [NIPS Conference, Denver, Colorado, USA, November 29 - December 4, 1999]*, pp. 1057–1063. The MIT Press, 1999. 6
- Youming Tao, Yulian Wu, Peng Zhao, and Di Wang. Optimal rates of (locally) differentially private heavy-tailed multi-armed bandits. In *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics*, volume 151 of *Proceedings of Machine Learning Research*, pp. 1546–1574. PMLR, 28–30 Mar 2022. 3
- Jay Tenenbaum, Haim Kaplan, Yishay Mansour, and Uri Stemmer. Differentially private multi-armed bandits in the shuffle model. In *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pp. 24956–24967, 2021. 3
- Qi Tian, Kun Kuang, Baoxiang Wang, Furui Liu, and Fei Wu. Multi-agent communication with graph information bottleneck under limited bandwidth. *CoRR*, abs/2112.10374, 2021. 3
- Aristide Charles Yedia Tossou and Christos Dimitrakakis. Achieving privacy in the adversarial multi-armed bandit. In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, February 4-9, 2017, San Francisco, California, USA*, pp. 2653–2659. AAAI Press, 2017. 3
- Mycal Tucker, Julie Shah, Roger Levy, and Noga Zaslavsky. Towards human-agent communication via the information bottleneck principle. *CoRR*, abs/2207.00088, 2022. doi: 10.48550/arXiv.2207.00088. 3
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pp. 5998–6008, 2017. 18

- Baoxiang Wang and Nidhi Hegde. Privacy-preserving q-learning with functional noise in continuous spaces. *Advances in Neural Information Processing Systems*, 32, 2019. 3
- Lingxiao Wang, Bargav Jayaraman, David Evans, and Quanquan Gu. Efficient privacy-preserving nonconvex optimization. *arXiv e-prints*, pp. arXiv–1910, 2019. 14
- Rundong Wang, Xu He, Runsheng Yu, Wei Qiu, Bo An, and Zinovi Rabinovich. Learning efficient multi-agent communication: An information bottleneck approach. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pp. 9908–9918. PMLR, 2020a. 3
- Tonghan Wang, Jianhao Wang, Chongyi Zheng, and Chongjie Zhang. Learning nearly decomposable value functions via communication minimization. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020b. 1, 3
- Hua Wei, Nan Xu, Huichu Zhang, Guanjie Zheng, Xinshi Zang, Chacha Chen, Weinan Zhang, Yanmin Zhu, Kai Xu, and Zhenhui Li. Colight: Learning network-level cooperation for traffic signal control. In *Proceedings of the 28th ACM International Conference on Information and Knowledge Management, CIKM 2019, Beijing, China, November 3-7, 2019*, pp. 1913–1922. ACM, 2019. 1
- Jiachen Yang, Alireza Nakhaei, David Isele, Kikuo Fujimura, and Hongyuan Zha. CM3: cooperative multi-goal multi-stage multi-agent reinforcement learning. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020. 1
- Dayong Ye, Tianqing Zhu, Zishuo Cheng, Wanlei Zhou, and S Yu Philip. Differential advising in multiagent reinforcement learning. *IEEE Transactions on Cybernetics*, 2020. 3
- Chongjie Zhang and Victor R. Lesser. Coordinated multi-agent reinforcement learning in networked distributed pomdps. In *Proceedings of the Twenty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2011, San Francisco, California, USA, August 7-11, 2011*. AAAI Press, 2011. 1
- Sai Qian Zhang, Qi Zhang, and Jieyu Lin. Efficient communication in multi-agent reinforcement learning via variance based control. In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pp. 3230–3239, 2019. 3
- Sai Qian Zhang, Qi Zhang, and Jieyu Lin. Succinct and robust multi-agent communication with temporal message control. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020. 3
- Xin Zhang, Zhuqing Liu, Jia Liu, Zhengyuan Zhu, and Songtao Lu. Taming communication and sample complexities in decentralized policy evaluation for cooperative multi-agent reinforcement learning. In *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pp. 18825–18838, 2021. 3
- Canzhe Zhao, Yanjie Ze, Jing Dong, Baoxiang Wang, and Shuai Li. Differentially private temporal difference learning with stochastic nonconvex-strongly-concave optimization. *arXiv preprint arXiv:2201.10447*, 2022. 3
- Kai Zheng, Tianle Cai, Weiran Huang, Zhenguo Li, and Liwei Wang. Locally differentially private (contextual) bandits learning. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020. 3
- Xingyu Zhou. Differentially private reinforcement learning with linear function approximation. *Proc. ACM Meas. Anal. Comput. Syst.*, 6(1):8:1–8:27, 2022. 3



## A PRIVACY ANALYSIS

In this section, we first present the proof of Theorem 5.1, which guarantees the  $(\epsilon_i, \delta)$ -DP for communication at each step. Then we give Corollary A.1, which provides episode-level  $(\epsilon_i, \delta)$ -DP guarantee for communication, together with its proof.

### A.1 PROOF OF THEOREM 5.1

We first present some necessary definitions and lemmas. We start by introducing Rényi differential privacy (RDP) and  $\ell_2$ -sensitivity.

**Definition A.1** (Rényi differential privacy, [Mironov \(2017\)](#)). For  $\alpha > 1$  and  $\rho > 0$ , a randomized mechanism  $f : \mathcal{D} \rightarrow \mathcal{R}$  is said to have  $\rho$ -Rényi differential privacy of order  $\alpha$ , or  $(\alpha, \rho)$ -RDP for short, if for any neighbouring datasets  $D, D' \in \mathcal{D}$  differing by one element, it holds that  $D_\alpha(f(D) \| f(D')) := \log \mathbb{E}(f(D)/f(D'))^\alpha / (\alpha - 1) \leq \rho$ .

**Definition A.2** ( $\ell_2$ -sensitivity, [Dwork et al. \(2014\)](#)). The  $\ell_2$ -sensitivity  $\Delta(q)$  of a function  $q$  is defined as  $\Delta(q) = \sup_{D, D'} \|q(D) - q(D')\|_2$ , for any two neighbouring datasets  $D, D' \in \mathcal{D}$  differing by one element.

With the properly added Gaussian noise, Lemma A.1 shows that the Gaussian mechanism could satisfy RDP.

**Lemma A.1** (Lemma 3.7, [Wang et al. \(2019\)](#)). For function  $q : \mathcal{S}^n \rightarrow \mathcal{R}$ , the Gaussian mechanism  $\mathcal{M} = q(S) + \mathbf{u}$  with  $\mathbf{u} \sim N(0, \sigma^2 \mathbf{I})$  satisfies  $(\alpha, \alpha \Delta^2(q) / (2\sigma^2))$ -RDP. Additionally, if  $\mathcal{M}$  is applied to a subset of all the samples which are uniformly sampled from the whole datasets without replacement using sampling rate  $\gamma$ , then  $\mathcal{M}$  satisfies  $(\alpha, 3.5\gamma^2 \Delta^2(q) \alpha / \sigma^2)$ -RDP with  $\sigma'^2 = \sigma^2 / \Delta^2(q) \geq 0.7$  and  $\alpha \leq 2\sigma'^2 \log(1/\gamma \alpha (1 + \sigma'^2)) / 3 + 1$ .

We now present the following two propositions regarding RDP. The first proposition shows that a composition of  $k$  mechanisms satisfying RDP is also a mechanism that satisfies RDP.

**Proposition A.1** (Proposition 1, [Mironov \(2017\)](#)). If  $k$  randomized mechanisms  $f_i : \mathcal{D} \rightarrow \mathcal{R}$  for all  $i \in [k]$ , satisfy  $(\alpha, \rho_i)$ -RDP, then their composition  $(f_1(D), \dots, f_k(D))$  satisfies  $(\alpha, \sum_{i=1}^k \rho_i)$ -RDP.

The second proposition provides the transformation from a RDP guarantee to a corresponding DP guarantee.

**Proposition A.2** (Proposition 3, [Mironov \(2017\)](#)). If a randomized mechanism  $f : \mathcal{D} \rightarrow \mathcal{R}$  satisfies  $(\alpha, \rho)$ -RDP, then  $f$  satisfies  $((\rho + \log(1/\delta)) / (\alpha - 1), \delta)$ -DP for all  $\delta \in (0, 1)$ .

We are now ready to give the formal proof of Theorem 5.1. Recall  $p_i = f_i^s(o_i, a_i; \theta_i^s) \sim \mathcal{N}(\mu_i, \Sigma_i)$  is the generated stochastic message before injecting privacy noise  $u_i \sim \mathcal{N}(0, \sigma_i^2 \mathbf{I}_d)$ , and  $m_i = p_i + u_i$  is the privatized message to be sent by agent  $i$ . Here we drop the dependency of message on the index of the target agent (i.e., we abbreviate the message  $m_{i,j}$  sent from agent  $i$  to agent  $j$  as  $m_i$ ) when it is clear from the context. We slightly abuse the notation by writing  $f_i^s = f_i^s(o_i, a_i; \theta_i^s)$ . To bound the sensitivity of  $f_i^s$ , one can perform the norm clipping to restrict the  $\ell_2$  norm of  $f_i^s$  by replacing  $f_i^s$  with  $f_i^s / \max(1, \|f_i^s\|_2 / C)$ , which ensures that  $\|f_i^s\|_2 \leq C$ . At each time  $t$ , for agent  $i$ , each message function  $f_i^s$  is applied to a subset of transitions in local trajectory  $\tau_i$  of agent  $i$  using uniform sampling without replacement with sampling rate  $\gamma_1$ , and agent  $i$  samples a subset of target agents to send messages by sampling without replacement with sampling rate  $\gamma_2$ .

*Proof of Theorem 5.1.* Due to the norm clipping of  $f_i^s$  and the triangle inequality, the  $\ell_2$ -sensitivity of  $f_i^s$  could be bounded as

$$\begin{aligned} \Delta_2(f_i^s) &= \sup_{D, D'} \|f_i^s(D) - f_i^s(D')\|_2 \\ &\leq 2C, \end{aligned} \tag{5}$$

where  $D, D' \in \mathcal{D}$  are any two neighbouring datasets differing by one element. By Equation (5) and Lemma A.1, the privatized message  $m_i$  satisfies  $(\alpha, 14\gamma_1^2 C^2 \alpha / \sigma^2)$ -RDP. Since agent  $i$  samples  $\gamma_2 N$  target agents to communicate, all the messages sent by agent  $i$  at time  $t$  are actually sent by a composite message function  $M_i = \{m_{i,k_j}\}_{j=1}^{\gamma_2 N}$ , which satisfies  $(\alpha, 14\gamma_2 \gamma_1^2 N C^2 \alpha / \sigma^2)$ -RDP by

**Proposition A.1.** Substituting  $\sigma_i^2 = \frac{14\gamma_2\gamma_1^2NC^2\alpha}{\beta\epsilon_i} = \frac{14\gamma_2\gamma_1^2NC^2\alpha}{\epsilon_i + \frac{\log \delta}{\alpha-1}}$  shows that  $M_i$  satisfies  $(\alpha, \epsilon_i + \frac{\log \delta}{\alpha-1})$ -RDP. Applying Proposition A.2 shows that  $M_i$  satisfies  $(\epsilon_i, \delta)$ -DP, which concludes the proof.  $\square$

## A.2 ANALYSIS OF EPISODE-LEVEL $(\epsilon_i, \delta)$ -DP

We would like to emphasize that Theorem 5.1 also provides an episode-level privacy guarantee. To see this, consider a multi-step game with finite episode length of  $T$ . Since Theorem 5.1 guarantees the  $(\epsilon_i, \delta)$ -DP for the communication mechanism  $M_i^t$  in each step  $t$ , the communication mechanism of the whole episode  $M_i = \{M_i^t\}_{t=1}^T$  satisfies  $(T\epsilon_i + (T-1)\frac{\log \delta}{\alpha-1}, \delta)$ -DP based on Proposition A.1 and Proposition A.2. Further, we can attain the  $(\epsilon_i, \delta)$ -DP of the communication of DPMAC in the whole episode by adjusting the noise variance in Theorem 5.1, detailed in the following corollary.

**Corollary A.1** (Episode-level  $(\epsilon_i, \delta)$ -DP guarantee for DPMAC). *Consider an episode with finite length  $T$ . Let  $\gamma_1, \gamma_2 \in (0, 1)$  having the same definitions as in Theorem 5.1, and let  $C$  be the sensitivity of the message functions. For any  $\delta > 0$  and privacy budget  $\epsilon_i$ , the communication of agent  $i$  in the whole episode satisfies  $(\epsilon_i, \delta)$ -DP when  $\sigma_i^2 = \frac{14\gamma_2\gamma_1^2NC^2\alpha T}{\beta\epsilon_i}$ , if we have  $\alpha = \frac{\log \delta^{-1}}{\epsilon_i(1-\beta)} + 1 \leq 2\sigma'^2 \log(1/\gamma_1\alpha(1+\sigma'^2))/3 + 1$  with  $\beta \in (0, 1)$  and  $\sigma'^2 = \sigma_i^2/(4C^2) \geq 0.7$ .*

*Proof.* Analogous to the proof of Theorem 5.1, by Equation (5) and Lemma A.1, the privatized message  $m_i^t$  at time  $t$  satisfies  $(\alpha, 14\gamma_1^2C^2\alpha T/\sigma^2)$ -RDP. Similarly, since agent  $i$  samples  $\gamma_2N$  target agents to communicate, all the messages sent by agent  $i$  at time  $t$  are through a composite message function  $\mathcal{M}_i^t = \{m_{i,k_j}^t\}_{j=1}^{\gamma_2N}$ . By Proposition A.1,  $\mathcal{M}_i^t$  satisfies  $(\alpha, 14\gamma_2\gamma_1^2NC^2\alpha T/\sigma^2)$ -RDP. Substituting  $\sigma_i^2 = \frac{14\gamma_2\gamma_1^2NC^2\alpha T}{\beta\epsilon_i} = \frac{14\gamma_2\gamma_1^2NC^2\alpha T}{\epsilon_i + \frac{\log \delta}{\alpha-1}}$ , we have  $\mathcal{M}_i^t$  satisfies  $(\alpha, (\epsilon_i + \frac{\log \delta}{\alpha-1})/T)$ -RDP. Then applying Proposition A.1 again shows that the composite message mechanism  $\mathcal{M}_i = \{\mathcal{M}_i^t\}_{t=1}^T$  is  $(\alpha, \epsilon_i + \frac{\log \delta}{\alpha-1})$ -RDP. As all the messages sent by agent  $i$  in the whole episode are through message mechanism  $\mathcal{M}_i$ , the privacy of these messages is strictly protected. Lastly, the proof is completed by applying Proposition A.2 to translate  $(\alpha, \epsilon_i + \frac{\log \delta}{\alpha-1})$ -RDP to  $(\epsilon_i, \delta)$ -DP for mechanism  $\mathcal{M}_i$ .  $\square$

## B EQUILIBRIUM ANALYSIS

### B.1 SINGLE-STEP GAME

In this section, we give the detailed analysis of our single-step game. We first give the notion of the potential game (PG) as follows.

**Definition B.1** (Potential game, Monderer & Shapley (1996)). *A two-player game  $G$  is a potential game if the mixed second order partial derivative of the utility functions are equal:*

$$\partial_{p_1}\partial_{p_2}u_1 = \partial_{p_1}\partial_{p_2}u_2.$$

The intuition behind the PG is that it tracks the changes in the payoff when some player deviates, without taking into account which one. Thus the PG usually helps with the analysis of the cooperative game, where the players might have the similar potential to act. To analyze the games involving multi-agent coordination with state dependence, the Markov potential game (MPG) is recently studied (Macua et al., 2018; Leonardos et al., 2021), where the action potential of all agents is described by a potential function. The solution concept of the PG relies on Nash equilibrium (NE) (Nash et al., 1950), the existence of which guarantees that all the agents could act in the best response to others.

We present the definition of the privacy loss function  $c^{\mathcal{M}}(p_n)$  in Definition B.2. Recall  $\mathcal{M}$  is the privacy preserving mechanism used in the game.

**Definition B.2** (Privacy loss function, Pejo et al. (2019)). *The privacy loss function  $c : [0, 1] \rightarrow [0, 1]$  is a continuous and twice differentiable function with  $c(0) = 1$ ,  $c(1) = 0$  and  $\partial_{p_n}c < 0$ .*

Then the benefit function  $b(V_n, V_n^{\mathcal{M}}(p_1, p_2))$  is given in Definition B.3. [Pejo et al. \(2019\)](#) consider the negative training error as the benefit while ours is different, as introduced here.  $V_n$  is the value that agent  $n$  receives when agent  $n$  acts alone without any cooperation.  $V_n(p_1, p_2)$  is the value that agent  $n$  receives when agent  $n$  acts cooperatively under the privacy mechanism  $\mathcal{M}$  and the actions  $p_1$  and  $p_2$ . Intuitively, we want that the benefit function portrays the benefit of the cooperative actions over the non-cooperative ones, thus in the meaningless case  $V_n \geq V_n^{\mathcal{M}}$ ,  $b(V_n, V_n^{\mathcal{M}}(p_1, p_2)) = 0$ . In addition,  $\partial_{p_n} b \leq 0$  since the value of the benefit function should decrease as the level of privacy protection increases.

**Definition B.3** (Benefit function). *The benefit function  $b : \mathbb{R}^+ \times \mathbb{R}^+ \rightarrow \mathbb{R}_0^+$  is a continuous and twice differentiable function, with  $\partial_{p_n} b \leq 0$  and  $b(V_n, V_n^{\mathcal{M}}(p_1, p_2)) = 0$  if  $V_n \geq V_n^{\mathcal{M}}$ .*

Finally, we give the definition of the value function  $V_n$  and  $V_n^{\mathcal{M}}(p_1, p_2)$ .  $V_n$  can be viewed as the upper bound of  $V_n^{\mathcal{M}}(p_1, p_2)$ , which intuitively means the value for the pure cooperation without privacy protection is always larger than that with some privacy protection. Note that the value function defined here is not the one which is commonly defined in RL literature.

**Definition B.4** (Value function).  *$V_n$  is the value function for agent  $n$  acting alone and without any cooperation.  $V_n^{\mathcal{M}} : [0, 1] \times [0, 1] \rightarrow \mathbb{R}^+$  is continuous, twice differentiable and:*

- $\exists m \in \mathcal{N} : p_m = 1 \Rightarrow \forall n \in \mathcal{N} : V_n^{\mathcal{M}}(p_1, p_2) \leq V_n$ ,
- $\forall n, m \in \mathcal{N} : \partial_{p_m} V_n^{\mathcal{M}} < 0$ ,
- $\forall n \in \mathcal{N} : V_n < V_n^{\mathcal{M}}(0, 0)$ .

For  $V_n^{\mathcal{M}}$ , the first rule ensures that if one agent protects the privacy entirely, the total value will be no greater than the value agents act alone. The second rule constrains the value function  $V_n^{\mathcal{M}}$  should be negatively monotonic wrt  $p_m$  since another agent's privacy gets stronger thus leading to less corporation. The third rule tells that the total value without any privacy protection (i.e., pure corporation) should be certainly larger than the value agents act alone.

The following theorem shows that the PG enjoys a good property that the NE exists. This motivates us to translate CGP into a potential game, which will be the key to prove Theorem 6.1.

**Theorem B.1** (Monderer and Shapley, [Monderer & Shapley \(1996\)](#)). *The potential game has at least one pure-strategy NE.*

We are not ready to prove Theorem 6.1.

*Proof of Theorem 6.1.* Starting from the definition of the two-player potential game (Definition B.1) and the utility function (Definition 6.1), we have

$$\begin{aligned} & \partial_{p_1} \partial_{p_2} u_1 = \partial_{p_1} \partial_{p_2} u_2, \\ \iff & \partial_{p_1} \partial_{p_2} b(V_1, V_1^{\mathcal{M}}(p_1, p_2)) = \partial_{p_1} \partial_{p_2} b(V_2, V_2^{\mathcal{M}}(p_1, p_2)), \\ \stackrel{(*)}{\iff} & (\partial_{V_n}^2 b) \cdot (\partial_{p_1} V_1^{\mathcal{M}} - \partial_{p_2} V_2^{\mathcal{M}}) = (\partial_{V_n}^2 b) \cdot (\partial_{p_1} \partial_{p_2} V_2^{\mathcal{M}} - \partial_{p_1} \partial_{p_2} V_1^{\mathcal{M}}), \end{aligned} \quad (6)$$

where  $(*)$  comes from the chain rule.

If the value function is with the property  $\partial_{p_1}^i V_1 = \partial_{p_2}^i V_2, \forall i \in \{1, 2\}$ , CGP will satisfy Equation 6. Further, by Theorem B.3, CGP is a potential game and has at least one non-trivial pure-strategy Nash equilibrium, thus concluding the proof.  $\square$

Utilizing Theorem 6.1 and designing the necessary functions for CGP, we show that single round binary sums could be proved with the existence of NE.

**Theorem B.2** (NE guarantee in single round binary sums). *For the single round binary sums game, let  $N = 2$  be the number of players,  $p_n$  be the probability of the random,  $c^{\mathcal{M}}(p_n) = 1 - p_n$  be the privacy loss function,  $V_n = -\frac{1}{2}$  and  $V_n^{\mathcal{M}}(p_1, p_2) = -\frac{1}{2}p_1^2 - \frac{1}{2}p_2^2 - p_1p_2 = -\frac{1}{2}(p_1 + p_2)^2$  be the value function,  $b(V_n, V_n^{\mathcal{M}}(p_1, p_2)) = V_n^{\mathcal{M}}(p_1, p_2) - V_n = -\frac{1}{2}p_1^2 - \frac{1}{2}p_2^2 - p_1p_2 + \frac{1}{2}$  be the benefit function,  $u_n(p_1, p_2) = B_n \cdot b(V_n, V_n^{\mathcal{M}}(p_1, p_2)) - C_n^{\mathcal{M}} \cdot c^{\mathcal{M}}(p_n)$  be the utility function, then*

Theorem 6.1 holds, which is to say, single round binary sums can be formulated into a CGP, further leading to the existence of one non-trivial pure-strategy NE. Further, the utility function is

$$u_n = -\frac{B_n}{2}(p_1 + p_2)^2 + C_n p_n + \frac{B_n}{2} - C_n, \quad \forall n \in \{1, 2\}. \quad (7)$$

The strategy taken by the agent is

$$p_n = \arg \max_{p_n} u_n = \frac{C_n}{B_n} - p_{-n}, \quad \forall n \in \{1, 2\}.$$

## B.2 MULTI-STEP GAME

In this section, we present the equilibrium analysis of the multi-step game. We start by presenting the definition of MPG.

**Definition B.5** (Markov potential game, Leonardos et al. (2021)). A Markov decision process (MDP),  $\mathcal{M}$ , is called a Markov Potential Game (MPG) if there exists a state-dependent function  $\Phi_s : \Pi \rightarrow \mathbb{R}$  for  $s \in \mathcal{S}$  such that

$$\Phi_s(\pi_i, \pi_{-i}) - \Phi_s(\pi'_i, \pi_{-i}) = V_s^i(\pi_i, \pi_{-i}) - V_s^i(\pi'_i, \pi_{-i})$$

holds for all agents  $i \in \mathcal{N}$ , all states  $s \in \mathcal{S}$  and all policies  $\pi_i, \pi'_i \in \Pi_i, \pi_{-i} \in \Pi_{-i}$ . By linearity of expectation, it follows that  $\Phi_\rho(\pi_i, \pi_{-i}) - \Phi_\rho(\pi'_i, \pi_{-i}) = V_\rho^i(\pi_i, \pi_{-i}) - V_\rho^i(\pi'_i, \pi_{-i})$ , where  $\Phi_\rho(\pi) := \mathbb{E}_{s \sim \rho} [\Phi_s(\pi)]$ .

The solution concept of the PG and the MPG relies on NE (Nash et al., 1950), the existence of which guarantees that all the agents could act best response of agent  $i$  with respect to the opponent's policy  $\pi_{\theta-i}^{-i}$  indicates the policy  $\pi_*^i$  such that  $V^i(s; \pi_*^i, \pi_{\theta-i}^{-i}) \geq V^i(s; \pi_{\theta^i}^i, \pi_{\theta-i}^{-i})$  for all feasible  $\pi_{\theta^i}^i$ . While it is still not clear how privacy will influence the equilibrium in the MARL setting, Kumari & Chakravarthy (2016) study the cooperative game with privacy and Pejo et al. (2019) model the game with privacy to solve the private learning problem.

Then we introduce the following assumptions and Theorem B.3 (Macua et al., 2018) to support the main theorems.

**Assumption 1.** The state and parameter sets,  $\mathbb{Z}$  and  $\mathbb{W}$ , are nonempty and convex.

We slightly abuse the notation in Assumption 2 and use  $\sigma_{k,i}$  as a given random variable with distribution  $p_{\sigma_k}(\cdot | x_i, a_i)$  instead of the standard deviation of the noise.

**Assumption 2.** The reward functions  $r_k(x_i, a_i, \sigma_{k,i})$  are twice continuously differentiable in  $\mathbb{Z} \times \mathbb{W}, \forall k \in \mathcal{N}$ .

**Assumption 3.** The state-transition function,  $f$ , and constraints,  $g$ , are continuously differentiable in  $\mathbb{Z} \times \mathbb{W}$ , and satisfy some regularity conditions (e.g., Mangasarian-Fromovitz).

**Assumption 4.** The reward functions  $r_k$  are proper, and there exists a scalar  $B$  such that the level sets  $\{a_0 \in \mathbb{C}_0, (x_i, a_i) \in \mathbb{C}_i : \mathbb{E}[r_k(x_i, a_i, \sigma_{k,i})] \geq B\}_{i=0}^\infty$  are nonempty and bounded  $\forall k \in \mathcal{N}$ .

**Theorem B.3** ((Macua et al., 2018)). Let Assumptions 1, 2, 3, 4 hold. Then, the game in Equation (8) of Macua et al. (2018) is an MPG if and only if: i) the reward function of every agent can be expressed as the sum of a term common to all agents plus another term that depends neither on its own state-component vector, nor on its policy parameter:

$$\begin{aligned} & r_k(x_{k,i}^r, \pi_k(x_{k,i}^\pi, w_k), \pi_{-k}(x_{-k,i}^\pi, w_{-k}), \sigma_{k,i}) \\ & = J(x_i, \pi(x_i, w), \sigma_i) + \Theta_k(x_{-k,i}^r, \pi_{-k}(x_{-k,i}^\pi, w_{-k}), \sigma_i), \quad \forall k \in \mathcal{N}; \end{aligned} \quad (8)$$

and ii) the following condition on the non-common term holds:

$$\mathbb{E} \left[ \nabla_{x_{k,i}^\pi} \Theta_k(x_{-k,i}^r, \pi_{-k}(x_{-k,i}^\pi, w_{-k}), \sigma_i) \right] = 0. \quad (9)$$

Moreover, if Equation (9) holds, then the common term in Equation (8),  $J$ , equals the potential function.

To make multiple round sums one MPG, the following detailed settings including the state space, the state transition, the action space, and the reward function are given. By such a specific design, we could achieve Theorem B.4, which shows the existence of a NE.

**State** The state  $x_{i,t} \in \mathbb{R}$  represents the remaining saving of agent  $i$  at time  $t$  (initially  $x_{i,0}$ ).

**State transition** The transition function is deterministic, which is  $x_{i,t+1} = x_{i,t} - b_{i,t}$ .

**Action** For  $i \in [N]$ , at time  $t$ , agent  $i$  with policy  $\pi_i$  assigns  $b_{i,t}$  savings according to the current state and the message sent by other agents. Then agent  $i$  performs the randomized perturbation on the assignment with privacy level  $p_{i,t}$ , with the message sender protocol  $f_i^s$ . After that the message  $m_{i,t}$  is generated and sent to other agents. Formally,

$$\begin{aligned} b_{i,t}, p_{i,t} &= \pi_i^b(x_{i,t}, m_{-i,t-1}), \pi_i^p(x_{i,t}, m_{-i,t-1}), \\ m_{i,t} &= f_i^s(x_{i,t}, b_{i,t}, p_{i,t}). \end{aligned}$$

**Reward function** The reward function for agent  $i$  is designed with the trade-off between privacy and utility, as shown in Equation (10). Privacy reward is performed to praise the privacy preserving. Formally,  $\forall i \in N$ ,

$$r_i(x_{i,t}, \pi_i(x_{i,t}), \pi_{-i}(x_{-i,t})) = \sum_{j \in N} (1 - p_{j,t}) b_{j,t} + \alpha x_{i,t} + \beta p_{i,t}. \quad (10)$$

The reward function given in Equation (10) could be written in two terms respectively, where  $\sum_{j \in N} (1 - p_{j,t}) b_{j,t}$  represents the first term in Equation (8) and  $\alpha x_{i,t} + \beta p_{i,t}$  represents the second term. Theorem B.4 is thus given.

**Theorem B.4** (NE guarantee in multiple round sums). *For the multiple round sums game, the reward function in Equation (10) will satisfy Theorem B.3 if we make the gradient chain between the message  $m_{-i}$  and the message function  $f_i$  cut down after sending. Further, if Assumptions 1, 2, 3, 4 (see Appendix B.2) are satisfied, our MPG could find a NE with potential function  $J$  as shown below,*

$$J(x_t, \pi(x_t)) = \sum_{j \in [N]} ((1 - p_{j,t}) b_{j,t} + \alpha x_{j,t} + \beta p_{j,t}). \quad (11)$$

## C IMPLEMENTATION DETAILS

**Agent** We use recurrent neural networks (RNN) for the actors of agents to approximate policies, which is shown with effectively in partially observable environments (Hausknecht & Stone, 2015), and build our communication method upon MADDPG (Lowe et al., 2017). However, we note that our communication method is general enough to be built on the top of any MARL algorithm with the CTDE paradigm.

**Communication protocol** The sender utilizes the multi-layer perceptron (MLP) and the receiver is an attention-based network (Vaswani et al., 2017). For the sender, a shared linear layer is first applied and two linear layer follows to output the mean the logarithm of the standard deviation. For the receiver, three linear layers corresponding to  $K, Q, V$  is directly applied with no shared layer. One can see that the network structure of our communication protocol is simple.

**Whether to encode the agent index** At each time  $t$ , the messages sent from agent  $i$  to different target agents could be sampled from different message distributions by including the index of target agent (say agent  $j$ ) into the input of  $f_i^s$  as  $p_{i,j} = f_i^s(o_i, a_i, j; \theta_i^s)$ , or just from the same message distribution by making  $p_{i,j} = f_i^s(o_i, a_i; \theta_i^s)$ . In our implementation, we do not include the target agent index into the input of  $f_i^s$  due to that including the target agent index into the input of  $f_i^s$  may lead to performance degeneration in some scenarios.

**Reparameterization trick** The stochastic message sender adopts the multivariate Gaussian distribution, which is implemented via the reparameterization trick (Kingma & Welling, 2014). Specifically, in our implementation, the stochastic message  $p_i$  is

$$p_i = \mu_i + \tilde{\sigma}_i \odot \xi,$$

where  $\mu_i$  and  $\tilde{\sigma}_i$  are the outputs of two neural networks of the sender,  $\xi \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ , and  $\odot$  denotes the element-wise product.



## D TRAINING DETAILS

**Optimization process of message senders and receivers** The message receiver of agent  $i$  encodes all the received messages into the encoded messages, which together with the observation of agent  $i$  will be fed into the actor of agent  $i$ . Hence the message receiver of agent  $i$  serves as a component in the actor of agent  $i$ , which will be updated by backpropagating the actor loss using policy gradient. By the chain rule, since the messages received by the receiver of agent  $i$  are sent from the message senders of all other agents, the message senders of all other agents will then be updated by backpropagation from the message receiver of agent  $i$ .

**Parameter setting** All the experiments are conducted on a server with 4 NVIDIA GeForce RTX 3090 GPUs. The hyperparameters of all algorithms are shown in Table 1. Since all the algorithms are built on the top of MADDPG, we tune the hyperparameters such that MADDPG has the best performance, which is also adopted in Kim et al. (2021). The hidden dimension for actor, critic, and the message protocol is selected by grid searching in  $\{32, 64, 128, 256\}$ , and the message dimension is selected by conducting grid search among values  $\{4, 8, 16, 32\}$ .

Table 1: Hyperparameters of all algorithms.

	MADDPG	TarMAC	I2C	DPMAC
Discount Factor	0.99	0.99	0.99	0.99
Batch Size (CCN&CC)	128	128	128	128
Batch Size (PP)	256	256	256	256
Buffer Size	$1 \times 10^4$	$1 \times 10^4$	$1 \times 10^4$	$1 \times 10^4$
Optimizer	Adam	Adam	Adam	Adam
Activation Function	ReLU	ReLU	ReLU	ReLU
Learning Rate	$7 \times 10^{-4}$	$7 \times 10^{-4}$	$7 \times 10^{-4}$	$7 \times 10^{-4}$
Hidden Dimension for Actor/Critic	128	128	128	128
Hidden Dimension for Message Protocol	—	32	32	32
Message Dimension	—	8	8	8

## E ENVIRONMENT DETAILS

**Cooperative navigation (CN)** Cooperative navigation is a standard task for multi-agent systems, introduced in Lowe et al. (2017), where agents target at reaching their own landmarks while avoiding collision. There are in total  $N = 3$  agents for our experiment setting.

**Cooperative communication and navigation (CCN)** The cooperative communication task is introduced in Mordatch & Abbeel (2017), where agents’ goal is to reach their target landmark while each agent only knows the location of the other agent’s target. The agents do not communicate with the channels embedded in the task and only share personal information with a learned protocol.

**Predator prey (PP)** The predator prey task is a standard task and we use the same setting and the same evaluation way as in Ding et al. (2020). Specifically, there are  $N = 3$  predators and  $M = 2$  preys in this task, whose initial positions are randomized initialized. Each predator is controlled by a agent, and each prey moves in the closest predator’s opposite direction. Since the speed of preys is higher than that of predators, cooperation is required for agents to capture a prey. The team reward is the negative sum of physical distances from all the predators to their closet preys. Besides, the predators will be penalized for each time any two predators collide with each other and we set the collision penalty as  $r_{\text{collision}} = -1$ . Each episode has 40 timesteps in this task.

## F ADDITIONAL EXPERIMENT RESULTS

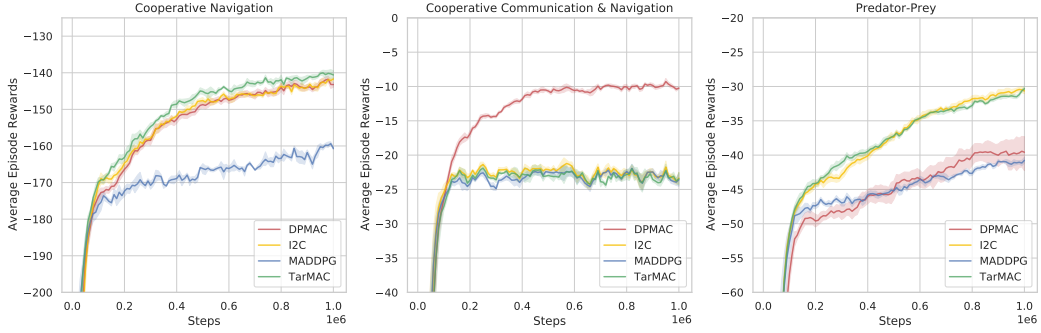


Figure 3: Learning curves where we compare DPMAC (non-private) with other baselines.

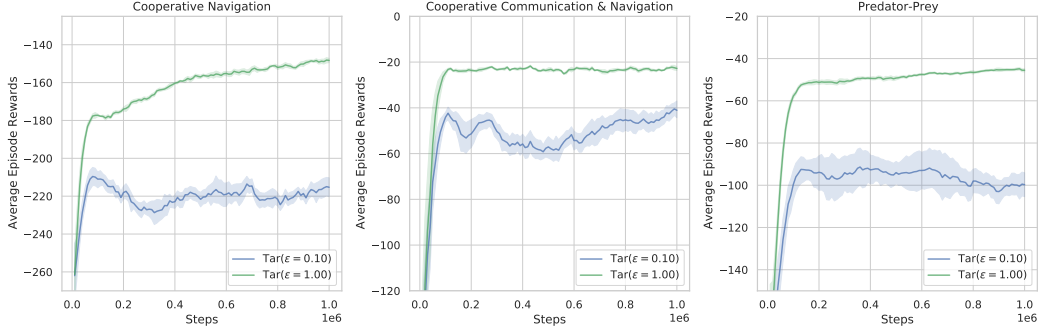


Figure 4: Learning curves of TarMAC with  $\epsilon = 0.1$  and  $\epsilon = 1.0$ .

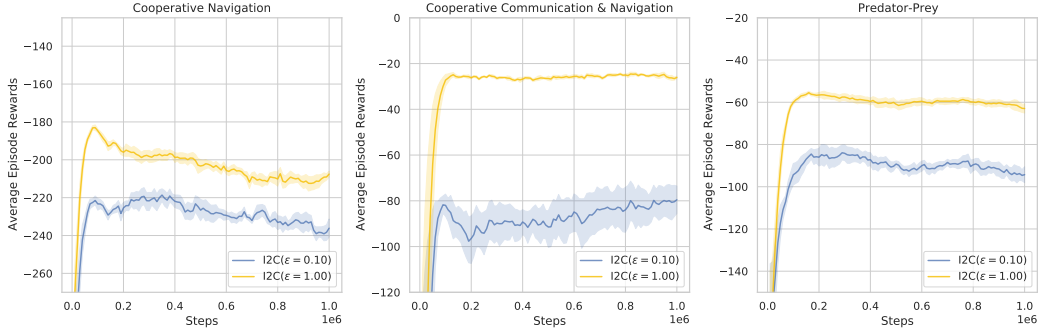


Figure 5: Learning curves of I2C with  $\epsilon = 0.1$  and  $\epsilon = 1.0$ .

## G ADDITIONAL EXPERIMENT ANALYSES

In this section, we present the additional experiment analyses regarding Figure 2.

### G.1 EXPERIMENT RESULTS WITHOUT PRIVACY

We first analyze our baseline methods, TarMAC and I2C, which are effective communication-based algorithm and can achieve superior performance than methods without communication such as

MADDPG. As shown in Figure 2a, TarMAC and I2C outperform MADDPG on cooperative navigation (CN) and predator-prey (PP) tasks and achieve similar performance as MADDPG on cooperative communication and navigation (CCN) task. However, DPMAC without privacy constraint outperforms TarMAC and I2C by a large margin on CCN task and is comparable with TarMAC and I2C on CN task. Under the privacy constraints, DPMAC with  $\epsilon = 0.10$  outperforms TarMAC and I2C on PP task. These results demonstrate the effectiveness of the message sender and receiver of DPMAC, which also has competitive performance under no privacy constraints.

## G.2 EXPERIMENT RESULTS WITH PRIVACY

We further investigate how the DP noise affects our DPMAC and other baseline methods. From Figure 2b, one can observe a serious performance degradation of TarMAC and I2C with a privacy budget of  $\epsilon = 0.10$ , which makes these two methods even worse than non-communicative MADDPG. In contrast, our DPMAC could still significantly outperform MADDPG on CN and PP tasks, and is comparable with MADDPG on CCN task. From Figure 2c, which is with a privacy budget of  $\epsilon = 1.00$  (note that  $\epsilon = 1.00$  enjoys weaker privacy guarantee than  $\epsilon = 0.10$  due to less privacy noises injected), our DPMAC still clearly outperforms all other methods on CN and PP tasks and is comparable with MADDPG. In sharp contrast to DPMAC, TarMAC still has clear degeneration on all tasks and is even worse than MADDPG on PP task, and I2C fails across all three tasks under such privacy constraints. These results demonstrate that the design principle of DPMAC that adjusting the learned message distribution via incorporating the distribution of privacy noise into the optimization process of the message sender and receiver can well alleviate the negative impacts of the privacy noises.

## G.3 DPMAC UNDER DIFFERENT PRIVACY BUDGETS

In Figure 2d, we further show the performance of DPMAC under different privacy budget  $\epsilon$ . On CN task, the performance of DPMAC with different privacy budgets are very close except  $\epsilon = 0.01$ , which is a very stringent privacy level. On CCN task, DPMAC without privacy constraint achieves the best performance, and DPMAC with different  $\epsilon$  could also reach meaningful performance. On PP task, DPMAC with  $\epsilon = 0.10$  and  $\epsilon = 1.00$  even outperform non-private DPMAC. As we have analyzed in Section 5.1, this is because DPMAC can learn the parameter  $\theta_i^s$  of the stochastic message sender such that  $D_{\text{KL}}(\mathcal{N}(\mu_i, \Sigma_i + \sigma_i^2 \mathbf{I}_d) \parallel \mathcal{N}(\mu_i^*, \Sigma_i^*)) \leq D_{\text{KL}}(\mathcal{N}(\mu_i', \Sigma_i') \parallel \mathcal{N}(\mu_i^*, \Sigma_i^*))$ , which means that the private messages can even encourage better team cooperation to gain higher team rewards than the non-private messages. Overall, one can see that the performance of DPMAC drops clearly only when  $\epsilon = 0.01$  on PP task. The above results also clearly demonstrate that DPMAC can learn to adjust the message distribution to alleviate the potential negative impacts of privacy noises, which ensures meaningful performance even under small privacy budgets and stabilizes the learning process.

## H EXPERIMENTS WITH EPISODE-LEVEL $(\epsilon_i, \delta)$ -DP

In this section, we present the experiment results under episode-level  $(\epsilon_i, \delta)$ -DP constraints with the privacy noise specified in Theorem A.1 and corresponding experiment analyses.

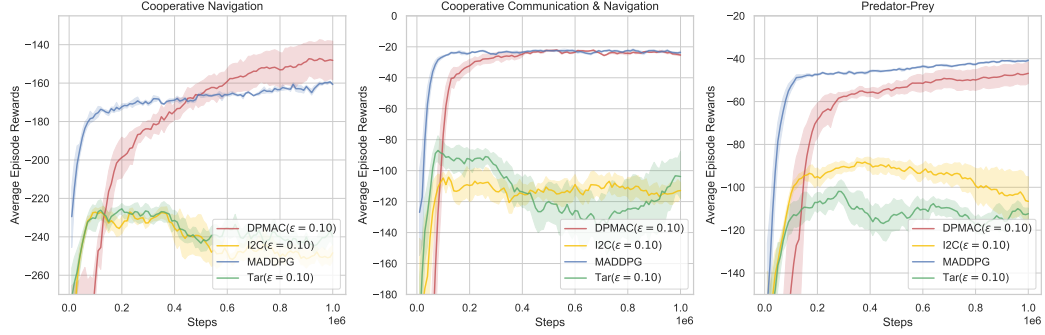
### H.1 EXPERIMENT RESULTS WITH EPISODE-LEVEL $(\epsilon_i, \delta)$ -DP

Comparing Figure 2b and 6a, under episode-level privacy constraint of  $\epsilon = 0.10$ , both the performance of TarMAC and I2C has a further degeneration across all the tasks. However, DPMAC could still clearly outperform MADDPG on CN task and is comparable with MADDPG on CCN task. On PP task, DPMAC is slightly outperformed by MADDPG but still exceeds TarMAC and I2C by a large margin. Comparing Figure 2c and 6b, I2C still fails across all tasks and TarMAC also has a clear performance drop. Specifically, TarMAC exceeded MADDPG on CN task and was comparable with MADDPG on CCN task in Figure 2c, but now TarMAC turns out to be outperformed by MADDPG on these tasks in Figure 6b. Besides, the performance gap between MADDPG and TarMAC now also becomes larger under the episode-level privacy constraint on the PP task. However, even though under episode-level privacy constraint of  $\epsilon = 1.00$ , DPMAC still outperforms MADDPG on CN and PP tasks and is comparable with MADDPG on the CCN task. The above results indeed

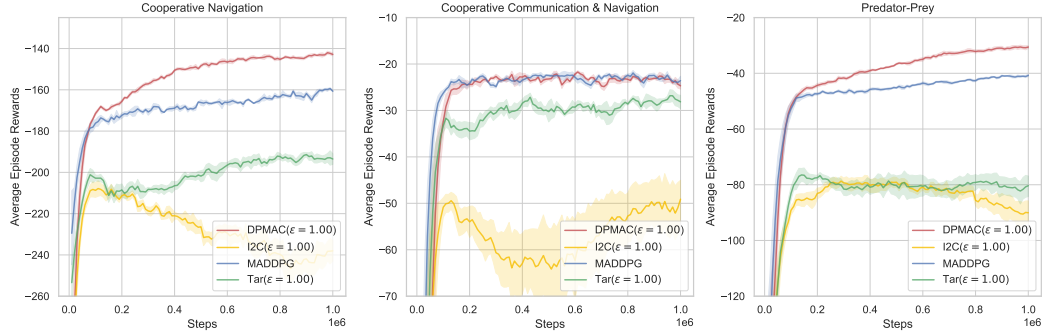
demonstrate that it is harder to learn under the episode-level privacy constraints since larger privacy noises are injected into the messages, but also validate the effectiveness of DPMAC, which still has good performance under episode-level privacy constraints.

## H.2 DPMAC UNDER DIFFERENT EPISODE-LEVEL PRIVACY BUDGETS

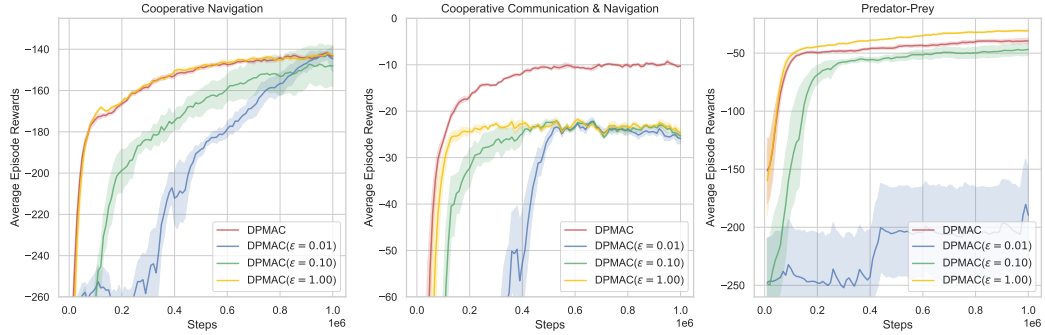
We now analyze the performance of DPMAC under different episode-level privacy constraints. Comparing Figure 2d and 6c, the performance of DPMAC under the episode-level privacy constraint of  $\epsilon = 1.00$  nearly remains the same as before. Under episode-level privacy constraint of  $\epsilon = 0.10$ , DPMAC still has the similar performance on CN and CCN tasks though converges slightly slower. But on the PP task, DPMAC also has a slight performance drop. Under episode-level privacy constraint of  $\epsilon = 0.01$ , though DPMAC still has the similar performance on CN and CCN tasks, it nearly fails to learn on the PP task. We leave the investigation of improving the performance of DPMAC under episode-level privacy constraint of  $\epsilon = 0.01$  as our future work. However, we note that the privacy level of  $\epsilon = 0.01$  is a rather stringent privacy constraint requiring injecting privacy noise with particularly large variance. Overall, the above results show that DPMAC can still learn to adjust the message distributions to alleviate the negative impacts of privacy noises even under episode-level privacy constraints.



(a) Learning curves of different algorithms under the episode-level privacy budget  $\epsilon = 0.10$ . MADDPG (non-private) is also displayed for comparison.



(b) Learning curves of different algorithms under the episode-level privacy budget  $\epsilon = 1.0$ . MADDPG (non-private) is also displayed for comparison.



(c) Learning curves of different episode-level privacy budgets ( $\epsilon = 0.01, 0.10, 1.00$ ) for DPMAC.

Figure 6: Learning curves of DPMAC and baseline algorithms. The curves are averaged over 5 seeds. Shaded areas denote 1 standard deviation.