Robust Coordination under Misaligned Communication via Power Regularization

Anonymous authors

Paper under double-blind review

Abstract

1	Effective communication in Multi-Agent Reinforcement Learning (MARL) can signif-
2	icantly enhance coordination and collaborative performance in complex and partially
3	observable environments. However, reliance on communication can also introduce vul-
4	nerabilities when agents are misaligned, potentially leading to adversarial interactions
5	that exploit implicit assumptions of cooperative intent. Prior work has addressed ad-
6	versarial behavior through power regularization by controlling the influence one agent
7	exerts over another, but has largely overlooked the role of communication in these dy-
8	namics. This paper introduces communicative power regularization (CPR), which ex-
9	tends power regularization specifically to communication channels. By explicitly quan-
10	tifying and constraining agents' communicative influence during training, CPR actively
11	mitigates vulnerabilities arising from misaligned or adversarial communications. Eval-
12	uations in the Grid Coverage benchmark environment demonstrate that our approach
13	significantly enhances robustness to adversarial communication while preserving coop-
14	erative performance, offering a practical framework for secure and resilient cooperative
15	MARL systems.

16 1 Introduction

Effective coordination among agents in Multi-Agent Reinforcement Learning (MARL) is crucial 17 for achieving collective goals. Communication, as an explicit exchange of information, is often em-18 19 ployed to facilitate this coordination, particularly in Cooperative MARL (CoMARL), where agents collaborate. However, common CoMARL approaches emphasizing parameter sharing for training 20 21 efficiency can lead to joint policies vulnerable to issues such as agent free-riding (Ueshima et al., 22 2023) or over-reliance on learned conventions (Köster et al., 2020). Such vulnerabilities are exac-23 erbated when agents are misaligned or face adversarial interactions, especially if policies implicitly 24 assume cooperative intent.

Objective misalignment, where agents pursue self-interested goals, makes public communication channels susceptible to sabotage, particularly against cooperative agents. Resilience against such misalignment is critical for deploying autonomous agent teams, requiring evaluation under nonstandard conditions, considering both team and individual contexts.

29 Communication remains a key research area in MARL (OroojlooyJadid & Hajinezhad, 2021), often 30 modeled with protocol controllers like CommNet (Sukhbaatar et al., 2016) and IC3Net (Singh et al., 2018). When agents learn communication and environment policies concurrently, they may develop 31 32 uncontrolled regularization against misaligned messages, whether from co-learning errors or inten-33 tional adversarial actions. Such self-learned communication, however, can create vulnerabilities if 34 naive agents are targeted. We define misaligned communication as any message negatively affecting 35 a recipient's performance, irrespective of explicit adversarial intent. Fostering resilience requires policies robust enough for mixed settings, differing from adversarial attacks that inject malicious 36 37 payloads (Tu et al., 2021; Dong et al., 2022).

38 Power, one agent's influence over another's utility and decisions, offers a mechanism to enhance pol-

39 icy robustness when incorporated into training. While agents typically do not explicitly optimize for

40 power, decomposing utility functions, akin to intrinsic rewards (Du et al., 2019), could offer greater

41 control. This paper introduces Communicative Power Regularization (CPR), extending the concept

42 of power regularization (Li & Dennis, 2024) specifically to communication channels. By quan-43 tifying and constraining communicative influence, CPR mitigates vulnerabilities from misaligned

45 trying and constraining communicative innuence, CFK innugates vulnerability 44 communication.

45 Our contributions are: (1) We propose CPR as a technique to control power dynamics within learned

46 communication policies. (2) We evaluate CPR in the Grid Coverage benchmark, demonstrating sig-

47 nificantly enhanced robustness to adversarial communication while preserving cooperative perfor-

48 mance. CPR thus offers a practical framework for more secure and resilient cooperative MARL.

49 This paper is organized as follows: Section 2 reviews related work. Section 3 covers preliminaries.

50 Section 4 details our CPR approach. Section 5 presents experimental results, followed by conclu-

51 sions in Section 6.

52 2 Related Work

Adversarial communication in MARL settings is often highlighted by its emergence in noncooperative settings (Blumenkamp & Prorok, 2020), which may be the product of misaligned agents. Adversarial attacks in MARL settings are diverse in their methodology, ranging from sparse targeted attacks (Hu & Zhang, 2022) to attacks that exploit vulnerabilities in mechanism design, such as consensus-based mechanisms (Figura et al., 2021) and adversarial minority influence (Li et al., 2024).

Adversarial training, an approach to mitigating against adversarial interests, is an umbrella-term for incorporating adversarial interactions into training for hardening and better resilience against adversarial opponents. In support, there are works on the robustness of CoMARL, such as Lin et al. (2020) and Guo et al. (2022). There are diverse defenses against adversarial communication, including works that consider test-time settings with theory of mind inspired mechanisms (Piazza & Behzadan, 2023). In our work, adversarial training is used to address misaligned communication.

65 Many CoMARL works that address credit assignment between global reward and local reward can be viewed as a means for regularizing agent behaviors and dynamics. For example, a reward-shaping 66 67 mechanism was proposed by Ibrahim et al. (2020) to portion out the team reward based on individual 68 contributions in order to address free-riders. Foerster et al. (2018) proposed COMA, Counterfactual 69 Multi Agent Policy Gradients, which marginalizes out single agent actions and also addresses credit 70 assignment with the incentive that agents will maximize their contribution to the global reward. The 71 motivation behind COMA is complementary in the sense that it quantifies how much influence an 72 individual agent's action has on the joint action, whereas this work quantifies how much influence 73 other agents' actions have upon an individual agent's policy.

74 Additionally, investigative work by Jaques et al. (2019), for example, explores causal relationships

among agents in MARL through counterfactual reasoning. While promoted for more efficient com-

76 munication and coordination, this approach can also quantify the contribution of other agents to a 77 self-agent's return.

78 Some other related works on explicit regularization in MARL originate from maximum-entropy 79 MARL, which reconstructs the return as the reward and the entropy of the policy distribution, 80 weighed by a temperature parameter. An example would be FOP (Zhang et al., 2021), an actor-81 critic method that factorizes the optimal joint policy from maximum-entropy MARL. The existing 82 work on quantifying power in MARL by Li & Dennis (2024) studies adversarial power, defined as 83 power associated with an adversarial opponent. The authors discuss various fine-tune parameters 84 for implementing power and measuring power in multi-opponent settings. This work investigates 85 power in settings with communication.

86 **3** Preliminaries

87 3.1 Communicative MARL

Multi-Agent Reinforcement Learning (MARL) provides a framework for sequential decisionmaking problems involving multiple interacting agents. Cooperative MARL scenarios are often formalized as Decentralized Partially Observable Markov Decision Processes (Dec-POMDPs), representable by the tuple:

$$\langle \mathcal{N}, \mathcal{S}, \{\mathcal{A}^i\}_{i \in \mathcal{N}}, \{\mathcal{O}^i\}_{i \in \mathcal{N}}, P, \{\mathcal{R}^i\}_{i \in \mathcal{N}}, \gamma \rangle$$

Here, \mathcal{N} denotes the set of N agents, \mathcal{S} is the global state space, \mathcal{A}^i is the action space for agent *i*, and \mathcal{O}^i is its observation space. The function $P(s'|s, \mathbf{a})$ defines the state transition dynamics, where $\mathbf{a} = \{a^i\}_{i \in \mathcal{N}}$ is the joint action assembled from individual agent actions sampled from their policies, $a^i \sim \pi^i(\cdot|o^i)$. Each agent *i* receives a local observation o^i and a reward $r^i = \mathcal{R}^i(s, \mathbf{a})$. In fully cooperative settings, agents typically share a common team reward, $r^i = R(s, \mathbf{a})$. The collective goal is to learn policies that maximize the expected discounted return $G = \sum_{t=0}^{T} \gamma^t r_t$ where $\gamma \in [0, 1]$ is the discount factor.

A prevalent paradigm for training MARL agents is Centralized Training with Decentralized Execution (CTDE). During the training phase, CTDE algorithms utilize global information, such as the full state or the actions of all agents, to facilitate learning. However, during execution, each agent must operate solely based on its local observation history. For instance, Value Decomposition Networks (VDN) (Sunehag et al., 2017) learn decentralized policies by decomposing the global team Q-function $Q_{tot}(s, a)$ into a sum of individual agent Q-functions $Q^i(o^i, a^i)$, as shown in Equation 1:

$$Q_{tot}(s,a) = \sum_{i \in \mathcal{N}} Q^i(o^i, a^i) \tag{1}$$

105 Such methods often employ parameter sharing across agent networks to improve learning efficiency.

106 To further enhance coordination, particularly under partial observability, agents can utilize explicit 107 communication. Communication channels allow agents to exchange information directly. Typically, 108 agent j generates a message m^j based on its internal state or history h^j via a learned communication policy $m^j \sim C^j(\cdot|h^j)$. Agent i's effective input s^i_{input} for decision-making can then incorporate 109 received messages m^{-i} from other agents alongside its own local observation $s_{input}^i = f(o^i, m^{-i})$. 110 Various architectures facilitate this exchange, such as those employing shared network modules for 111 message processing (e.g., CommNet) or learning selective communication via gating mechanisms 112 (e.g., IC3Net). The structure of communication, dictating which agents can exchange messages, is 113 114 frequently modeled using graph representations.

115 3.2 Implicit Communication via Graph Neural Networks

116 Communication can also be learned implicitly through structured feature aggregation. As ex-117 plored in Li et al. (2020) and utilized in our experiments, Graph Neural Networks (GNNs) offer 118 a powerful mechanism for this. Instead of learning explicit messages, agents first process their 119 local observations o_t^i to generate feature embeddings $x_t^i \in \mathbb{R}^F$. These are stacked into a matrix 120 $X_t = [x_t^1, \ldots, x_t^N]^T \in \mathbb{R}^{N \times F}$. The GNN operates over a dynamic communication graph repre-121 sented by an adjacency matrix (or Graph Shift Operator) $S_t \in \mathbb{R}^{N \times N}$, where $[S_t]_{ij} = 1$ if agent j122 can transmit to agent i at time t. (typically based on proximity).

123 The core mechanism is a graph convolution layer. For a single layer GNN transforming input fea-124 tures $X_{in} \in \mathbb{R}^{N \times F}$ to output features $X_{out} \in \mathbb{R}^{N \times G}$, the operation can be defined as shown in 125 Equation 2:

$$X_{out} = \sigma \left(\sum_{k=0}^{K-1} S_t^k X_{in} A_k \right)$$
⁽²⁾

Here, $S_t^k X_{in}$ represents features aggregated from the *k*-hop neighborhood, effectively requiring *k* rounds of message passing or communication exchanges. *K* is the maximum communication hop count (filter size) defining the spatial receptive field of the graph convolution. $A_k \in \mathbb{R}^{F \times G}$ are learnable weight matrices specific to hop *k*, transforming and combining features across hops and dimensions. $\sigma(\cdot)$ is a non-linear activation function applied element-wise. For multi-layer GNNs (*L*

$$X_l = \sigma[\mathcal{A}_l(X_{l-1}; S_t)] \quad \text{for } l = 1, \dots, L$$
(3)

132 where X_0 is the input to the first layer and A_l denotes the graph convolution operation at layer l.

133 This GNN architecture learns what information is relevant to share and aggregate from the local 134 neighborhood defined by S_t and K. The resulting aggregated feature vector for agent *i*, denoted 135 $[X_L]_i \in \mathbb{R}^{G_L}$ (the *i*-th row of the final layer's output), captures context from its communicative 136 neighbors and serves as the input to its decentralized policy network $\pi^i(a^i|[X_L]_i)$.

137 3.3 Power

The concept of power refers to the influence one agent has over another agent's decision-making 138 139 and utility. In shared environments, power dynamics play a critical role in determining how agents interact and coordinate. Li & Dennis (2024) introduced power as a formal measure, redefining the 140 141 optimization criterion as a combination of expected task return and power utility. By incorporating 142 power regularization into the training process, agents can learn policies that are more resilient to 143 states where power imbalances make them vulnerable. Specifically, power quantifies the expected 144 difference between the current joint policy and a hypothetical joint policy where other agents take 145 adversarial actions over k steps.

Power is closely related to the concept of game security, which evaluates the expected return when facing adversarial opponents. In sequential games, the minimax strategy is often used to select the best action among the worst possible outcomes. When power dynamics naturally emerge or are necessary for completing team tasks, power regularization provides designers with a mechanism to control the autonomous behaviors of agents, ensuring they remain robust to adversarial influences.

The original formulation of power by Li and Dennis, referred to as standard power in this paper, estimates the influence agent j has over agent i as follows (Equation 4):

$$\rho_{i:j}^{\text{standard}}(\pi^{i}, \pi^{j}, s) = Q_{i}^{\pi^{i}, \pi^{j}}(s, a^{i}) - \min_{a^{j} \in A^{j}} Q_{i}^{\pi^{i}, \pi^{j}}(s, a^{j})$$
(4)

Here, $Q_i^{\pi^i,\pi^j}(s,a^i)$ represents the expected return for agent *i* under the joint policy π^i,π^j , while min_{$a^j \in A^j$} $Q_i^{\pi^i,\pi^j}(s,a^j)$ represents the worst-case return if agent *j* takes an adversarial action. In cooperative settings, the joint policy π differs from adversarial policies, stabilizing the estimation of power.

To incorporate power into the learning process, the state-value function for agent i is modified to include a power regularization term, as shown in Equation 5:

$$V_{i}(s,a) = V_{i}^{\pi}(s,a) + \lambda V_{i}^{\pi,\rho_{i:j}}(s,a)$$
(5)

Here, $V_i^{\pi}(s, a)$ is the original state-value function of agent *i*, and $V_i^{\pi, \rho_{i:j}}(s, a)$ represents the power component, which penalizes states where agent *j* exerts high influence over agent *i*. This regularization can be viewed as a form of reward shaping, where the power measure is used to guide agents toward policies that are less vulnerable to adversarial influence.

163 4 Power Regularization Over Communication

164 Learning communication in cooperative settings can lead to more efficient coordination and strong, 165 mutually dependent relationships among agents. However, misaligned agents can exploit these de-166 pendencies through sensory manipulation over the communication medium. Given the potential 167 misuse of the communication medium, it is important to address how much dependency an agent 168 delegates to other agents through the communication channel or protocol. Furthermore, it is im-169 perative to ask how much dependency an agent should delegate to the communication medium, 170 regardless of who uses the communication medium. In our work, we train policies to be more 171 resilient to misaligned communication and miscommunication through adversarial training. Adver-172 sarial training is the practice of incorporating a variety of adversarial experiences and adversarial 173 communication into training. We propose Communicative Power Regularization (CPR) to improve 174 robustness against adverse impacts from misaligned communication by incorporating adversarial 175 messages during training. Unlike standard power, which keeps the agent state constant, the com-176 munication message affects the perceived agent state and, therefore, can be viewed as a form of 177 state regularization where the proximity of other states consisting of adversarial messages affects its 178 perceived utility, similar to that of stochastic transitions by an environment.

179 4.1 Communicative Power Regularization (CPR)

180 We define communicative power as the decomposition of power into two components: standard 181 power and power of communication. Standard power is the power delegated to other agents without 182 leveraging the communication channel or protocol. In contrast, the power of communication is the 183 power delegated to other agents over the communication channel or protocol. The total power ρ_{ij}^{CPR} 184 that agent *j* has over agent *i* is defined as the sum of these two components (Equation 6):

$$\rho_{ij}^{\text{CPR}}(\pi^{i}, \pi^{j}, s^{i}, m^{j}) = \rho_{ij}^{\text{Standard}}(\pi^{i}, \pi^{j}, s^{i}) + \rho_{ij}^{\text{Communication}}(\pi^{i}, \pi^{j}, s^{i}, m^{j})$$

$$(6)$$

Here, $\rho_{ij}^{\text{Standard}}(\pi^i, \pi^j, s^i)$ represents the standard power, which quantifies the influence agent *j* has over agent *i* through actions alone, and $\rho_{ij}^{\text{Communication}}(\pi^i, \pi^j, s^i, m^j)$ represents the power of communication, which quantifies the influence agent *j* has over agent *i* through communication.

188 The communicative power ρ_{ij}^{CPR} is defined as shown in Equation 7:

$$\rho_{ij}^{\text{CPR}}(\pi^{i}, \pi^{j}, s^{i}, m^{j}) = Q_{i}^{\pi^{i}, \pi^{j}}(s, m^{j}, a^{i}, s^{i'}, m^{j'}) - \min_{a^{j} \in A^{j}} Q_{i}^{\pi^{i}, \pi^{j}}(s^{i}, m_{\text{adv}}^{j}, a^{i}, s^{i'}, m^{j'})$$

$$(7)$$

This measures the difference in agent *i*'s expected return when agent *j* takes an adversarial action and sends an adversarial message m_{adv}^{j} compared to when agent *j* follows the joint policy.

Standard power can directly regulate misaligned communication in scenarios where communications are considered actions in the action space. However, its effectiveness in regularizing state-related misaligned communication assumes that appropriate variance is introduced into training, such as simultaneously learning a communication policy with an environment policy. Communicative power incorporates adversarial messages, whether they are individually sent or aggregated. This is particularly important in cases where individual messages are not misaligned but the aggregation of messages is misaligned.

198 In the traditional approach, the methodology does not provide direct defense mechanisms against 199 adversarial attacks that specifically exploit model parameterization. Adversaries perform adversar-200 ial attacks to craft and inject adversarial samples, which usually target a model's parameterization 201 (e.g., a neural network's decision boundary). However, some state-actions performed by certain 202 agent roles contribute more to the environment's expected utility than others, which finite-budget 203 adversaries often consider. To address these challenges, we propose a framework that explicitly 204 accounts for the influence of communication on power dynamics, ensuring robustness against both misaligned communication and adversarial exploitation of model parameterization. 205

The subsequent expressions are adapted from Li & Dennis (2024) to align with our proposed setting. To incorporate power into the learning process, the state-value function for agent i is modified to 208 include a power regularization term, as shown in Equation 8:

$$V_i(s,a) = V_i^{\pi}(s,a) + \lambda V_i^{\pi,\rho_{ij}}(s,a)$$
(8)

Here, $V_i^{\pi}(s, a)$ is the original state-value function for the task, and $V_i^{\pi, \rho_{ij}}(s, a)$ represents the power component, which penalizes states where other agents exert influence over agent *i* through both actions and communication. The parameter λ is a scalar that controls the degree of power regularization over the expected return.

Another approach involves applying both standard power and power of communication separately to
 enable individualized penalization of states where other agents exert greater control and penalization
 for states where there is excessive reliance on communicated messages.

The power regularization term $V_i^{\pi,\rho_{ij}}(s,a)$ is defined as the sum of power rewards $R_i^{\text{power}}(s_t,\pi)$ over states starting from *s* reached by unrolling the policy π , as given by Equation 9:

$$V_i^{\pi,\rho_{ij}}(s,a) = \sum_{t=0}^T R_i^{\text{power}}(s_t,\pi)$$
(9)

In the 2-agent setting, the power reward $R_i^{\text{power}}(s, \pi)$ is defined as shown in Equation 10:

$$R_{i}^{\text{power}}(s,\pi) = -\rho_{ij}^{\text{CPR}}(\pi^{i},\pi^{j},s^{i},m^{j})$$
(10)

This indicates that the power reward penalizes the influence agent j has over agent i through both actions and communication in the state s. In settings with more than two agents, the power reward

captures the strongest individual influence that any other agent j exerts on agent i, as defined in Equation 11:

$$R_{i}^{\text{power}}(s,\pi) = -\max_{j \neq i} \rho_{ij}^{\text{CPR}}(\pi^{i},\pi^{j},s^{i},m^{j})$$
(11)

By applying a maximization, this formulation emphasizes the worst-case dependency, making the regularization more sensitive to the most dominant external influence.

Our definition of communicative power is to further specify how power is allocated in the presence of a communication medium. This is in contrast to standard power, which makes no distinction over how power is distributed over coordinating devices or mechanisms. It is within the designer's discretion whether communication is appropriate for a cooperative task.

229 5 Experiment Results

We evaluate CPR in the Grid Coverage (Blumenkamp & Prorok, 2020) (GC) environment. This setting allows us to assess CPR's effectiveness in a scenario where cooperative agents must coordinate effectively while being resilient to misaligned communication from adversarial entities.

233 Grid Coverage. To evaluate the effectiveness of Communicative Power Regularization (CPR) in 234 mitigating the effects of adversarial communication, we conducted experiments within the non-235 Convex Coverage map from the Adversarial Comms repository (Blumenkamp & Prorok, 2020). 236 This environment challenges a team of cooperative agents to maximize area coverage on a grid 237 map while contending with explicit adversarial agents designed to disrupt cooperative performance 238 through the communication channel. Agents operate with limited local observations and commu-239 nication ranges, utilizing a CNN-GNN-MLP architecture to process environmental input, exchange 240 messages via the GNN, and select actions using the MLP. Detailed experimental configuration pa-241 rameters for this environment are provided in Appendix A (Table 2).

Our evaluation directly compares the performance of cooperative MARL agents trained with CPR against baseline agents trained without CPR. Both sets of agents were evaluated over 100 trials in scenarios featuring varying numbers of cooperative and adversarial agents, where all agents, includ-ing adversaries, actively communicated throughout the episodes. This setup isolates the impact of CPR on the robustness of the cooperative strategy against communication-based attacks.

Table 1: Grid Coverage: Cooperative agents' scores (mean (\pm std.) over 100 trials), comparing
performance with and without CPR across various [adversarial, cooperative] agent compositions.
Asterisked (*) configurations denote training and evaluation with the same number of agents; others
are scaled in evaluation.

# of agents	Cooperative agents' scores		Improvement (%)
	With CPR	Without CPR	
$[1, 5]^*$	$257.74 (\pm 45.93)$	218.82 (± 66.40)	18
$[2, 4]^*$	204.92 $(\pm$ 59.84)	93.56 (± 77.53)	119
$[3,3]^*$	$188.85 (\pm 88.50)$	33.53 (± 39.69)	463
[6, 30]	294.65 $(\pm$ 31.21 $)$	285.40 (± 43.90)	3
[4, 8]	$\textbf{242.25}~(\pm~\textbf{53.18})$	116.27 (± 83.21)	108
[8, 16]	$262.29(\pm46.94)$	134.22 (± 86.52)	95
[12, 24]	$273.94 (\pm 39.33)$	$142.86 \ (\pm \ 89.00)$	92
[18, 18]	$232.66 (\pm 77.85)$	61.89 (± 64.19)	276

247 The comprehensive results presented in Table 1 quantitatively demonstrate the significant advan-248 tage conferred by CPR across various team compositions and evaluation scales. We first established 249 baseline performance in configurations where training and evaluation agent counts were identical 250 ([1,5], [2,4], [3,3]). In these scenarios, cooperative agents employing CPR consistently achieved 251 substantially higher mean scores, often with reduced variance as indicated by the standard devia-252 tions, compared to baseline agents without CPR. For instance, with 1 adversary and 5 cooperative 253 agents ([1,5]), CPR-trained agents achieved a mean score of 257.74, while baseline agents scored 254 218.82. This performance gap widened as the proportion of adversarial agents increased; in the chal-255 lenging [3,3] scenario, agents with CPR maintained a strong cooperative score of 188.85, whereas 256 the performance of baseline agents severely degraded to 33.53.

257 To assess the scalability of the learned policies, models trained on these starred configurations were 258 then evaluated on significantly larger teams without retraining. Remarkably, CPR-trained agents 259 consistently maintained robust performance levels even in these scaled-up scenarios. For example, 260 when scaling from the [2,4] training setup, agents with CPR achieved mean scores of 242.25 for 261 [4,8] and 262.29 for [8,16], substantially outperforming baseline agents whose scores were 116.27 262 and 134.22, respectively. This trend continued with further scaling; for instance, in the [12,24] setup 263 (also scaled from [2,4]), CPR agents scored 273.94 against the baseline's 142.86. Even in the highly 264 scaled [18,18] scenario (from [3,3] training), CPR-enabled agents achieved a mean score of 232.66, 265 a stark contrast to the 61.89 achieved by baseline agents. While the score difference in the [6,30] 266 configuration (294.65 with CPR vs. 285.40 without, scaled from [1,5]) was more modest, likely 267 due to the very low adversary-to-cooperative agent ratio, CPR still provided a clear benefit. The 268 overall trend strongly indicates that CPR facilitates the learning of robust coordination strategies 269 that generalize effectively and preserve a high level of performance when deployed in larger, more 270 complex multi-agent systems.

Figure 1 provides a more granular view, illustrating the average coverage percentage achieved by cooperative agents over episode time steps for [1,5], [2,4], and [3,3] configurations, respectively.

In all depicted scenarios, the agents trained with CPR (blue curves) consistently outperform the baseline agents (red curves). They not only reach a higher final coverage percentage but also exhibit faster convergence towards their optimal performance early in the episode. Furthermore, the tighter variance bands (shaded regions) associated with the CPR agents suggest that CPR contributes to



Figure 1: Grid coverage, average cooperative coverage percentage (over 100 trials) across varying team compositions, comparing agents trained with CPR (blue) vs without CPR (red).

more stable and reliable performance across trials, reducing the detrimental impact of adversarialinterference.

Synthesizing these findings, both the aggregate scores and the temporal coverage dynamics unequivocally show that CPR enhances the ability of cooperative agents to maintain effective coordination and achieve superior performance in the presence of adversarial communication. By regularizing the power or influence of messages, CPR fosters more robust communication strategies that are less susceptible to manipulation. This validation in the complex Grid Coverage task, featuring decentralized control and explicit adversaries, underscores the practical value of CPR for developing resilient multi-agent systems.

286 A key concern is that CPR might inadvertently lead agents to avoid communication to prevent penal-287 ties. To investigate this, we conducted an ablation study and evaluated 5 CPR-trained cooperative 288 agents, comparing their performance with active communication versus communication explicitly 289 disabled. The results (detailed in Appendix B, Figure 2) clearly refute this concern. Agents uti-290 lizing communication (blue triangles) significantly outperform the same agents operating without 291 it (orange circles), achieving faster score accumulation and a higher final score. This confirms that 292 CPR-trained agents do not abandon communication but learn to use it robustly, leveraging it for 293 improved coordination and performance, thus demonstrating that CPR encourages resilient commu-294 nication strategies rather than avoidance.

295 6 Conclusion

296 This work tackled the inherent vulnerability of communicating MARL systems where reliance on 297 information exchange can be exploited by misaligned agents. We argued that prior work on power 298 regularization, focused on action-based influence, inadequately captures the risks associated with 299 delegating control through communication protocols. To address this, we introduced Communica-300 tive Power Regularization (CPR), a method to enhance MARL system robustness against misaligned 301 communication by penalizing over-reliance on communication channels and increasing the self-302 autonomy of agents. Evaluations in the Grid Coverage environment demonstrated CPR's ability to 303 significantly improve cooperative performance and resilience in adversarial communication scenar-304 ios, without sacrificing communication when beneficial. While CPR involves a trade-off between 305 robustness and optimal cooperative performance, it provides a practical framework for developing 306 more secure and reliable cooperative MARL systems. Future work could focus on developing adap-307 tive CPR frameworks, for instance by employing an adaptive Power Regularization Factor (λ), and 308 on addressing heterogeneous trust dynamics. CPR offers a valuable step towards deploying resilient 309 multi-agent systems in challenging real-world environments.

Grid Coverage: Experimental Configuration Parameters 310 Α

311 This appendix details the experimental configuration parameters used for the Grid Coverage envi-

ronment, as referenced in Section 5. 312

Description	Value	
Episode max timestep	345	
Communication range	16	
Observation range	8	
Paward	+1 for visiting a previously	
Reward	uncovered cell; 0 otherwise	
Agent actions	up, down, left, right, stay	
World shape	24×24	
Cooperative training	20 million time steps	
Adversarial training	20 million time steps	
Power regularization factor (λ)	0.3	

Table 2: Grid Coverage: Experimental Configuration Parameters

Grid Coverage: Ablation Study on Communication 313 B

314 To address the concern that Communicative Power Regularization (CPR) might inadvertently incen-

tivize agents to cease communication, we conducted an ablation study. This study, referenced in 315

Section 5, compared the performance of five CPR-trained cooperative agents in the Grid Coverage 316

317 environment under two conditions: (1) with their standard learned communication enabled, and (2)

with their ability to send or receive messages explicitly disabled. 318

319 Figure 2 presents the cumulative average scores over 100 trials for these two conditions. The results 320 demonstrate that agents actively utilizing their learned communication protocols achieve signifi-

321 cantly better performance than when communication is unavailable.



Cooperative Agents Cumulative Average Score

Figure 2: Grid coverage, cumulative average score of CPR-trained agents (5 cooperative, 100 trials), comparing performance with (blue triangles) and without (orange circles) communication.

322 **References**

- Jan Blumenkamp and Amanda Prorok. The emergence of adversarial communication in multi-agent reinforcement learning, 2020. URL https://arxiv.org/abs/2008.02616.
- Juncheng Dong, Suya Wu, Mohammadreza Sultani, and Vahid Tarokh. Multi-agent adversarial
 attacks for multi-channel communications, 2022. URL https://arxiv.org/abs/2201.
 09149.
- Yali Du, Lei Han, Meng Fang, Ji Liu, Tianhong Dai, and Dacheng Tao. Liir: Learning individual intrinsic reward in multi-agent reinforcement learning. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (eds.), Advances in Neural Information Processing Systems, volume 32. Curran Associates, Inc., 2019. URL https://proceedings.neurips.cc/paper_files/paper/2019/ file/07a9d3fed4c5ea6b17e80258dee231fa-Paper.pdf.
- Martin Figura, Krishna Chaitanya Kosaraju, and Vijay Gupta. Adversarial attacks in consensus based multi-agent reinforcement learning, 2021. URL https://arxiv.org/abs/2103.
 06967.
- Jakob Foerster, Gregory Farquhar, Triantafyllos Afouras, Nantas Nardelli, and Shimon Whiteson.
 Counterfactual multi-agent policy gradients. *Proceedings of the AAAI Conference on Artificial Intelligence*, 32(1), Apr. 2018. DOI: 10.1609/aaai.v32i1.11794. URL https://ojs.aaai.
 org/index.php/AAAI/article/view/11794.
- Jun Guo, Yonghong Chen, Yihang Hao, Zixin Yin, Yin Yu, and Simin Li. Towards comprehensive
 testing on the robustness of cooperative multi-agent reinforcement learning, 2022. URL https:
 //arxiv.org/abs/2204.07932.
- Yizheng Hu and Zhihua Zhang. Sparse adversarial attack in multi-agent reinforcement learning,
 2022. URL https://arxiv.org/abs/2205.09362.
- Aly Ibrahim, Anirudha Jitani, Daoud Piracha, and Doina Precup. Reward redistribution mechanisms
 in multi-agent reinforcement learning. In *Adaptive Learning Agents Workshop at the International Conference on Autonomous Agents and Multiagent Systems*, 2020.
- Natasha Jaques, Angeliki Lazaridou, Edward Hughes, Caglar Gulcehre, Pedro Ortega, Dj Strouse,
 Joel Z. Leibo, and Nando De Freitas. Social influence as intrinsic motivation for multi-agent deep
 reinforcement learning. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 3040–3049. PMLR, 09–15 Jun 2019. URL https://proceedings.
 mlr.press/v97/jaques19a.html.
- Raphael Köster, Kevin R. McKee, Richard Everett, Laura Weidinger, William S. Isaac, Edward
 Hughes, Edgar A. Duéñez-Guzmán, Thore Graepel, Matthew Botvinick, and Joel Z. Leibo.
 Model-free conventions in multi-agent reinforcement learning with heterogeneous preferences,
 2020. URL https://arxiv.org/abs/2010.09054.
- Michelle Li and Michael Dennis. The benefits of power regularization in cooperative reinforcement
 learning, 2024. URL https://arxiv.org/abs/2406.11240.
- Qingbiao Li, Fernando Gama, Alejandro Ribeiro, and Amanda Prorok. Graph neural networks
 for decentralized multi-robot path planning, 2020. URL https://arxiv.org/abs/1912.
 06095.
- Simin Li, Jun Guo, Jingqiao Xiu, Yuwei Zheng, Pu Feng, Xin Yu, Aishan Liu, Yaodong Yang,
 Bo An, Wenjun Wu, and Xianglong Liu. Attacking cooperative multi-agent reinforcement
 learning by adversarial minority influence, 2024. URL https://arxiv.org/abs/2302.
 03322.

- Jieyu Lin, Kristina Dzeparoska, Sai Qian Zhang, Alberto Leon-Garcia, and Nicolas Papernot. On the
 robustness of cooperative multi-agent reinforcement learning, 2020. URL https://arxiv.
 org/abs/2003.03722.
- Afshin OroojlooyJadid and Davood Hajinezhad. A review of cooperative multi-agent deep rein forcement learning, 2021. URL https://arxiv.org/abs/1908.03963.
- Nancirose Piazza and Vahid Behzadan. A theory of mind approach as test-time mitigation
 against emergent adversarial communication, 2023. URL https://arxiv.org/abs/
 2302.07176.
- Amanpreet Singh, Tushar Jain, and Sainbayar Sukhbaatar. Learning when to communicate at scale
 in multiagent cooperative and competitive tasks, 2018. URL https://arxiv.org/abs/
 1812.09755.
- Sainbayar Sukhbaatar, Arthur Szlam, and Rob Fergus. Learning multiagent communication with
 backpropagation. In *Proceedings of the 30th International Conference on Neural Information Processing Systems*, NIPS'16, pp. 2252–2260, Red Hook, NY, USA, 2016. Curran Associates
 Inc. ISBN 9781510838819.
- Peter Sunehag, Guy Lever, Audrunas Gruslys, Wojciech Marian Czarnecki, Vinicius Zambaldi, Max
 Jaderberg, Marc Lanctot, Nicolas Sonnerat, Joel Z. Leibo, Karl Tuyls, and Thore Graepel. Value decomposition networks for cooperative multi-agent learning, 2017. URL https://arxiv.
 org/abs/1706.05296.
- James Tu, Tsunhsuan Wang, Jingkang Wang, Sivabalan Manivasagam, Mengye Ren, and Raquel
 Urtasun. Adversarial attacks on multi-agent communication, 2021. URL https://arxiv.
 org/abs/2101.06560.
- Atsushi Ueshima, Shayegan Omidshafiei, and Hirokazu Shirado. Deconstructing cooperation and
 ostracism via multi-agent reinforcement learning, 2023. URL https://arxiv.org/abs/
 2310.04623.
- Tianhao Zhang, Yueheng Li, Chen Wang, Guangming Xie, and Zongqing Lu. Fop: Factorizing
 optimal joint policy of maximum-entropy multi-agent reinforcement learning. In Marina Meila
 and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning*,
 volume 139 of *Proceedings of Machine Learning Research*, pp. 12491–12500. PMLR, 18–24 Jul
 2021. URL https://proceedings.mlr.press/v139/zhang21m.html.
- Content that appears after the references are not part of the "main text," have no page limits, are not necessarily reviewed, and should not contain any claims or material central to the paper. If your paper includes supplementary materials, use the
- 401 \beginSupplementaryMaterials
- 402 command as in this example, which produces the title and disclaimer above. If your paper does not403 include supplementary materials, this command can be removed or commented out.