POLICY CONSISTENCY IN MULTI-AGENT REINFORCE MENT LEARNING WITH MIXED REWARD

Anonymous authors

004

010 011

012

013

014

015

016

017

018

019

021

025

026

Paper under double-blind review

ABSTRACT

The sparsity of team rewards poses a significant challenge that hinders the effective learning of optimal team policies in cooperative multi-agent reinforcement learning. One common approach to mitigate this issue involves augmenting sparse rewards with individual rewards to guide policy training. However, a significant drawback of such approaches is that modifying the reward function can potentially alter the optimal policy. To tackle this challenge, we propose a novel multi-agent policy optimization approach that ensures consistency between the mixed policy (learned from a combination of individual and team rewards) and the team policy (based solely on team rewards), through a new policy consistency constraint that aligns the returns of both policies in policy optimization model. We further develop an iterated policy optimization procedure to solve the formulated problem, deriving an approximate optimization objective for each iteration of the mixed and team policies. Experimental evaluation conducted in the StarCraft II Multi-Agent Challenge Environment (SMAC), Multi-Agent Particle Environment (MPE), and Google Research Football (GRF) environments demonstrate that our proposed approach effectively addresses the policy inconsistency problem, *i.e.*, it evenly outperforms strong baseline methods.

028 029 1 INTRODUCTION

Cooperative multi-agent reinforcement learning (MARL) has attracted significant interest due to its potential in solving complex decision-making tasks (Yan & Xu, 2020; Chen et al., 2024). Despite advancements in MARL algorithms, the issue of sparse team rewards remains a major obstacle, limiting the practical application of these algorithms in real-world scenarios such as power grids (Tittaferrante & Yassine, 2021), aerial vehicles (Du et al., 2021), and robotics (Sun et al., 2020).

Previous research addressing the issue of sparse rewards commonly relies on additional individual dense rewards. These approaches can be categorized into three main types: utilizing expert knowl-037 edge (Kurach et al., 2020; Lowe et al., 2017; Huang et al., 2022; Zhu & Zhao, 2021), exploring the state space (Strehl & Littman, 2008; Bellemare et al., 2016; Liu et al., 2023; Jeon et al., 2022; Liu et al., 2021; Xu et al., 2024), and action exploration (Li et al., 2021; Xu et al., 2023a). While 040 incorporating individual rewards have shown promise in addressing sparse rewards, recent studies 041 highlight a critical issue: learned policies may deviate from optimal policies due to modifications 042 in the reward function, especially in multi-agent environments (Wang et al., 2022). For instance, in a cooperative battle simulation, agents incentivized by individual rewards may prioritize individual 043 skills (such as shooting or escaping) over the collective goal of winning the battle. 044

The motivation of the our research is illustrated in Figure 1. While team rewards are expected to guide agents towards the optimal policy (orange dashed line), the sparsity of the reward function often hinders the policy learning process (orange solid line). Shaping mixed rewards (through a combination of team and individual rewards) can facilitate more efficient policy learning but may lead to suboptimal policies due to alterations in the reward function (green line). This highlights the need to maintain consistency between the learned policy and the optimal team policy when incorporating mixed rewards.

While IRAT (Wang et al., 2022) mitigates the policy inconsistency issue through improving policy similarity between learned and team policies, our approach completely eliminates this issue by deriving exact policy objectives from a constrained Lagrangian dual optimization model. By maximizing

054 mixed rewards subject to consistency constraints between 055 learned and team policies' cumulative rewards, we derive an op-056 timization objective with an extended TD error. Unlike IRAT's 057 standard TD error using only individual rewards, CMT incor-058 porates team rewards with a Lagrangian multiplier λ . Team rewards provide a more comprehensive metric for policy evaluation, while λ enforces policy consistency constraints in the dual 060 optimization problem. These innovations yield policies with 061 higher team rewards and reduced variance. Moreover, unlike 062 IRAT's focus on individual rewards, our approach incorporates 063 mixed rewards during training, better balancing individual skill 064 execution and group collaboration. 065

More specifically, our approach begins with the presented con-066 strained policy optimization problem, which is transformed into 067 its Lagrangian dual form, allowing us to solve it with the un-068 known optimal team policy. Furthermore, we establish the 069 equivalence between the solutions of the original problem and its dual counterpart. 071

We propose the Consistency between Mixed and Team policies 072 (CMT) algorithm, which iteratively updates both policies for 073 each agent to solve the Lagrangian dual problem. Using perfor-074 mance difference lemma (Kakade & Langford, 2002) and pol-075 icy approximation techniques, we simplify the objective func-076 tion with an extended TD error, while avoiding data inefficiency 077



Training steps

Figure 1: Sparse team rewards often hinders the policy learning process (orange solid line) despite the expectation of guiding agents towards the optimal team policy (orange dashed line). Shaping mixed rewards allows for more efficient policy learning, but it may lead to suboptimal policy due to changes in the reward function (green line). Our approach introduces mixed rewards to efficiently develop policy while ensuring consistency with the optimal team policy (red line). Detailed test results are provided in Section 5.

from simultaneous sampling of both mixed and team policies. Further, we reconstruct the objective 078 with KL terms between policies, maintaining objective equivalence while constraining policy gaps. 079

Extensive experiments across SMAC, MPE, and GRF environments demonstrate the effectiveness of 080 our proposed approach. Specifically, the proposed approach achieves a 28.5 percentage point higher 081 winning rate and 4.2 percentage point lower standard deviation compared to IRAT across 11 maps 082 of SMAC environments. Furthermore, our approach outperforms other state-of-the-art baselines, 083 including MAPPO, QMIX, MASER, and LAIES, across nearly all tasks. Overall, CMT achieved 084 the best performance in 20 out of 21 tasks across all benchmarks. 085

087

880

2 **RELATED WORK**

Individual Rewards in MARL: Introducing individual rewards has become one of the most preva-090 lent and effective strategies to mitigate the sparse reward issue in MARL. Existing research in this 091 area can be broadly categorized into three groups: expert knowledge-based, state space explorationbased, and action space exploration-based approaches. 092

093 External expert knowledge plays a crucial role in formulating individual rewards by leveraging prior 094 understanding of environmental dynamics (Kurach et al., 2020; Lowe et al., 2017; Huang et al., 095 2022; Zhu & Zhao, 2021). For example, a simple design rewards the elimination of enemies and 096 the health of teammates in SMAC (Samvelyan et al., 2019). Further, MAPPER algorithm utilizes 097 expert knowledge to decompose tasks and construct sub-tasks with dense rewards (Liu et al., 2020). 098 However, relying solely on external knowledge can be impractical, as obtaining such knowledge is prohibitively expensive in real-world environments (Zhang et al., 2021a; Ryu et al., 2022). 099

100 To address the aforementioned challenge, some works incorporate individual rewards based on ac-101 quired transactional information. One straightforward approach is to explore novel states by count-102 ing visited states (Strehl & Littman, 2008; Bellemare et al., 2016). However, this approach faces 103 difficulties in complex environments with vast state spaces. Recently, methods such as EDTI/EITI 104 have been developed to promote the exploration of novel states that significantly influence agents' 105 actions (Wang et al., 2019). LAIES partitions the state space into internal and external states, constructing individual rewards to promote exploration of external states (Liu et al., 2023). MASER 106 formulates individual rewards based on the distance between the current state and a target state 107 chosen by the Q-value of visited states and actions (Jeon et al., 2022). Furthermore, DIFFER decomposes team experience into individual experience for constructing individual rewards (Hu et al., 2023). Additionally, Liu et al. (2021) and Xu et al. (2024; 2023b) exploit prior structural knowledge to encourage agents to explore subsets of the state space.

In contrast to state exploration, another line of research focuses on constructing individual rewards by exploring the action space. One example is the CDS approach (Li et al., 2021), which aims to maximize the mutual information between agent identities and trajectories, thereby encouraging more diverse actions. Xu et al. (2023a) introduce individual rewards based on the concept of joint policy diversity, which quantifies the disparity between the current policy and previous policies.

Policy consistency in RL: While previous studies have shown promising results in addressing reward sparsity, the policy inconsistency issue (resulting from the introduction of individual rewards) is often overlooked. To the best of our knowledge, there are two works focusing on this inconsistency issue. The first work proposes a constrained policy optimization method within a single-agent environment (Chen et al., 2022). This method iteratively updates the extrinsic policy and the actual policy using lower bounds of dual objectives as optimization objectives. In contrast, our approach leverages transformed dual optimization objectives to train policies directly, thereby avoiding the bias introduced by lower bounds of the objective function.

124 The work closest to ours is IRAT (Wang et al., 2022). While IRAT focuses on reshaping the opti-125 mization objective to enhance the *policy similarity* between individual and team policies, Figure 4 126 in Appendix A shows that merely maximizing policy similarity does not necessarily lead the learned 127 policy to converge towards an optimal team policy. Besides, the oversight of mixed rewards prevents 128 IRAT from achieving a balance between individual skill execution and collaborative team objec-129 tives. This work therefore derives the optimization objective precisely by solving the Lagrangian 130 dual optimization problem under a policy consistency constraint, which ensures the equivalence of 131 cumulative team rewards obtained by the learned mixed policy and the optimal team policy.

132 133

134

3 PRELIMINARIES

135 In a cooperative multi-agent decision task, the decentralized partially observable Markov decision 136 process (Dec-POMDP) framework defined as $G = \langle \mathcal{N}, \mathcal{S}, \{\mathcal{A}^i\}_{i \in \mathcal{N}}, P, R, \{\mathcal{Z}^i\}_{i \in \mathcal{N}}, \mathcal{O}, \gamma \rangle$ is commonly used to model the problem. Herein, $\mathcal{N} = \{1, 2, \dots, n\}$ denotes the set of agents. S 137 represents the global state space, with $s \in S$ denoting the environmental state. \mathcal{A}^i is the action 138 space of agent i, and $a^i \in \{\mathcal{A}^i\}$ denotes the action taken by agent i. P denotes the transition 139 probability, specifying the probability of transitioning from state s to state s' under a joint action 140 $a = (a_t^1, a_t^2, \dots, a_t^n)$. R represents the reward function, with $r^E \in R$ denoting the shared team 141 rewards received by all agents. $o^i \in \mathcal{Z}^i$ represents the local observation by agent i based on the 142 observation function $\mathcal{O}: \mathcal{S} \times \mathcal{N} \to \mathcal{Z}^i$, where \mathcal{Z}^i denotes the observation space of agent *i*. Given 143 the observation-action history $\tau^i \in T^i = (\{Z^i\} \times \{A^i\})$, agent *i* learns a team policy $\pi^i_E(a^i | \tau^i)$ 144 with the aim to maximize the following **cumulative team reward**: 145

$$\max J_E\left(\pi_E^i\right) = \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^t r_t^E\right],\tag{1}$$

where $\gamma \in [0, 1]$ is the discount factor used to weigh the importance of future rewards. The optimal team policy is defined as $\pi_E^{i^*} = \arg \max_{\pi_E^i} J_E(\pi_E^i)$.

151 Sparse reward presents a common challenge in MARL. To mitigate this issue, individual reward 152 functions are often introduced into policy learning process. In this setting, at each time step t, agents 153 select a joint action a, and the environment returns a reward $\mathbf{r} = (r_t^1, r_t^2, \dots, r_t^n, r_t^E)$, consisting of 154 individual reward r_t^i and a shared team reward r_t^E for each agent *i*. Consequently, agent *i* learns a 155 mixed policy π_{E+i}^i with the aim to maximize the following **cumulative mixed reward**:

156

146

147

157 158

$$\max \hat{J}_{E+i}\left(\pi_{E+i}^{i}\right) = \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^{t}\left(r_{t}^{E} + r_{t}^{i}\right)\right],\tag{2}$$

where $\hat{r}_t^i = r_t^E + r_t^i$ is the mixed reward for each agent *i*. Note that reward shaping, including the determination of the importance of individual reward in mixed reward, plays a crucial role in exploration of RL. This is not the primary focus of this paper. For further insights on this topic, readers are referred to Chen et al. (2022) and Yuan et al. (2023).

¹⁶² 4 METHOD

With a team reward-oriented objective, the optimal team policy for each agent *i* is denoted as $\pi_E^{i*} = \arg \max_{\pi_E^i} J_E(\pi_E^i)$, as defined in Equation 1. However, when introducing a mixed reward-oriented objective, the optimal mixed policy shifts to $\pi_{E+i}^{i*} = \arg \max_{\pi_{E+i}^i} \hat{J}_{E+i}(\pi_{E+i}^i)$ according to Equation 2. It is evident that, at convergence, the optimal mixed policy must deviate from the optimal team policy due to the change in the reward function.

To resolve this inconsistency issue, we first present a consistency constrained policy optimization
model, which defines the research target of our work. Then, we propose a dual policy optimization
procedure to solve the optimization model. Finally, we integrate our approach into the centralized
training with decentralized execution (CTDE) framework and outline the implementation of our
CMT algorithm.

175 176

4.1 CONSISTENCY CONSTRAINED POLICY OPTIMIZATION MODEL

In the environment incorporating mixed rewards, the policy optimization problem with policy consistency constraint for each agent i is given by:

$$\max_{\pi_{E+i}^{i}} \hat{J}_{E+i} \left(\pi_{E+i}^{i}\right) \quad \text{subject to} \quad J_{E} \left(\pi_{E+i}^{i}\right) - \max_{\pi_{E}^{i}} J_{E} \left(\pi_{E}^{i}\right) = 0.$$
(3)

As the performance is often evaluated by the cumulative team rewards J_E , the difference between the cumulative team rewards achieved by the learned mixed policy and that learned by the optimal team policy is constrained by Eq. 3. Moreover, the objective function in Eq. 3 remains consistent with the optimization objective introduced in Eq. 2 evaluated by the cumulative mixed rewards \hat{J}_{E+i} . As such, the formulated policy optimization problem aims to find a mixed policy that maximizes the mixed rewards while maintaining consistency with the optimal team policy.

To tackle the intractability of directly solving the policy optimization problem with the unknown term $\max_{\pi_E^i} J_E(\pi_E^i)$ in the consistency constraint, we transform the problem into its Lagrangian dual. The Lagrangian dual problem is given by:

$$\min_{\lambda} \left[\max_{\pi_{E+i}^{i}} \hat{J}_{E+i} \left(\pi_{E+i}^{i} \right) + \lambda \left(J_{E} \left(\pi_{E+i}^{i} \right) - \max_{\pi_{E}^{i}} J_{E} \left(\pi_{E}^{i} \right) \right) \right], \tag{4}$$

where λ represents the Lagrangian multiplier associated with the consistency constraint.

To establish the equivalence between the original problem and its Lagrangian dual, we make the following assumption.

197 Assumption 1. There exists a policy π_{E+i}^i such that $J_E\left(\pi_{E+i}^i\right) - \max_{\pi_E^i} J_E\left(\pi_E^i\right) = 0.$

Assumption 1 requires that there exists a mixed policy π_{E+i}^i (developed by the mixed reward in Eq. 2), the performance of which can match that of optimal team policy π_E^{i*} (defined in Eq. 1) concerning the cumulative team rewards J_E . This assumption is commonly observed in RL (Sun & Xu, 2023; Wang et al., 2022). During the initial stage of policy training process, the policy π_{E+i}^{i*} yields lower cumulative team reward J_E compared to the optimal team policy π_E^{i*} . However, guided by individual rewards, π_{E+i}^i is expected to improve its performance as the training proceeds, and finally approaches or even matches π_E^{i*} . Given this continuous learning process, we can reasonably assume the existence of a policy π_{E+i}^i that satisfies $J_E(\pi_{E+i}^i) - \max_{\pi_E^i} J_E(\pi_E^i) = 0$.

Under Assumption 1, the Slater's condition (Ding et al., 2020; Zhang et al., 2021b) holds. Consequently, we conclude that the solution to Eq. 3 is equivalent to the solution of its Lagrangian dual problem.

211 212

4.2 MIN-MAX DUAL POLICY OPTIMIZATION

To find the optimal solution of Eq. 4, we first rewrite it for each agent i as follows:

$$\min_{\lambda} \left[\min_{\pi_{E}^{i}} \max_{\pi_{E+i}^{i}} \hat{J}_{E+i}^{\lambda} \left(\pi_{E+i}^{i} \right) - \lambda J_{E} \left(\pi_{E}^{i} \right) \right]$$
(5)

where $\hat{J}_{E+i}^{\lambda}(\pi_{E+i}^{i}) := \hat{J}_{E+i}(\pi_{E+i}^{i}) + \lambda J_{E}(\pi_{E+i}^{i}).$

In Eq. 5, we observe that there is an opposing optimization objective between the mixed policy π_{E+i}^{i} and the team policy π_E^i . For the mixed policy π_{E+i}^i , the objective is to maximize $\hat{J}_{E+i}^{\lambda}(\pi_{E+i}^i)$ – $\lambda J_E(\pi_E^i)$, which is minimized in the optimization objective of the team policy π_E^i . Next, we discuss procedures to optimize the mixed and team policies with the proposed dual objective through an iterated optimization approach.

Optimizing mixed policy: The optimizing objective of mixed policy π_{E+i} for agent *i* is

$$\max_{\substack{\pi_{E+i}^{i}}} \hat{J}_{E+i}^{\lambda} \left(\pi_{E+i}^{i} \right) - \lambda J_{E} \left(\pi_{E}^{i} \right).$$
(6)

We first expand the objective in Eq. 6 based on the performance difference lemma (Kakade & Langford, 2002) (details are provided in Appendix B.1):

$$\hat{J}_{E+i}^{\lambda}\left(\pi_{E+i}^{i}\right) - \lambda J_{E}\left(\pi_{E}^{i}\right) = -\mathbb{E}_{\pi_{E}^{i}}\left[\sum_{t=0}^{\infty} \lambda r_{t}^{E} - V_{E+i}^{\pi_{E+i}^{i},\lambda}\left(\tau_{t}^{i}\right) + \gamma V_{E+i}^{\pi_{E+i}^{i},\lambda}\left(\tau_{t+1}^{i}\right)\right] = -\sum_{i\in\mathcal{T}^{i}} d_{\rho_{0}}^{\pi_{E}^{i},\gamma}\left(\tau^{i}\right) \sum_{e\in\mathcal{A}} \pi_{E}^{i}\left(a_{t}^{i}|\tau_{t}^{i}\right) U_{E+i}\left(\tau_{t}^{i},a_{t}^{i}\right),$$
(7)

where $V_{E+i}^{\pi_{E+i}^i,\lambda} = \mathbb{E}_{\pi_{E+i}^i} \left[\sum_{t=0}^{\infty} \gamma^t \left[r^i + (1+\lambda) \, r^E \right] \right]$ and $d_{\rho_0}^{\pi_E^i,\gamma} \left(\tau^i \right) = \sum_{t=0}^{\infty} \gamma^t P \left(\tau_t^i = \tau^i | \rho_0, \pi_E^i \right).$

 U_{E+i} is a extended TD error for evaluating mixed policy under the environment with policy consistency requirement. It is defined as:

$$U_{E+i} := \lambda r_t^E - V_{E+i}^{\pi_{E+i}^i,\lambda} \left(\tau_t^i \right) + \gamma V_{E+i}^{\pi_{E+i}^i,\lambda} \left(\tau_{t+1}^i \right).$$
(8)

A significant challenge arises from data inefficiency when directly optimizing policies according to Eq. 7, primarily due to the impracticality of sampling from both policies π_{E+i} and π_E simultaneously. To overcome this challenge, we draw inspiration from the approach proposed in Chen et al. (2022), which leverages trajectories from one policy to approximate another similar policy. Specifically, we approximate the team policy by utilizing trajectories generated by the mixed policy, based on the assumption that mixed and team policies are similar (*similarity assumption*) (Schulman et al., 2017; Kakade & Langford, 2002; Schulman et al., 2015). As a result, the optimization objective can be approximated as follows by changing team policy into mixed policy:

$$-\mathbb{E}_{\pi_{E}^{i}}\left[\sum_{t=0}^{\infty} U_{E+i}\left(\tau_{t}, a_{t}\right)\right] = -\sum_{\tau^{i} \in T^{i}} d_{\rho_{0}}^{\pi_{E}^{i}, \gamma}\left(\tau^{i}\right) \sum_{a \in A} \pi_{E+i}^{i}\left(a_{t}^{i}|\tau_{t}^{i}\right) U_{E+i}\left(\tau_{t}^{i}, a_{t}^{i}\right) \\ = -\sum_{\tau^{i} \in T^{i}} d_{\rho_{0}}^{\pi_{E}^{i}, \gamma}\left(\tau^{i}\right) \sum_{a \in A} \pi_{E}^{i}\left(a_{t}^{i}|\tau_{t}^{i}\right) \frac{\pi_{E+i}^{i}\left(a_{t}^{i}|\tau_{t}^{i}\right)}{\pi_{E}^{i}\left(a_{t}^{i}|\tau_{t}^{i}\right)} U_{E+i}\left(\tau_{t}^{i}, a_{t}^{i}\right).$$

$$\tag{9}$$

It is worth noting that the *similarity assumption* holds true particularly when the mixed and team policies networks share the same initialized parameters, leading to the minimal disparity between the two types of policies.

To further simplify the computation process while preventing the mixed policy from deviating from the team policy, we introduce a transformation of the objective by incorporating the KL divergence between the mixed and team policies. This transformation maintains the equality between the original objective and the transformed objective. The derivation process for this transformation can be found in Appendix B.1. The final optimization objective for mixed policy is expressed as follows:

$$-\mathbb{E}_{\pi_{E}^{i}}\left[\min\left\{\frac{\pi_{E+i}^{i}\left(a_{t}^{i}|\tau_{t}^{i}\right)}{\pi_{E}^{i}\left(a_{t}^{i}|\tau_{t}^{i}\right)}U_{E+i}, clip\left(\frac{\pi_{E+i}^{i}\left(a_{t}^{i}|\tau_{t}^{i}\right)}{\pi_{E}^{i}\left(a_{t}^{i}|\tau_{t}^{i}\right)}, 1-\epsilon, 1+\epsilon\right)U_{E+i}\right\}\right] - D_{KL}\left(\pi_{E+i}^{i}||\pi_{E}^{i}\right)$$

$$(10)$$

Optimizing team policy: The optimizing objective of team policy π_E^i in Eq. 5 can be rewritten as

$$\max_{\pi_E^i} \lambda J_E\left(\pi_E^i\right) - \hat{J}_{E+i}^\lambda\left(\pi_{E+i}^i\right). \tag{11}$$

270 Algorithm 1 CMT algorithm 271 Input: 272 For each agent *i*, initialize parameters θ_{k+i}^i for actor network of mixed policy, ϕ_{k+i}^i for critic network of 273 mixed policy, θ_E^i for actor network of team policy, ϕ_E^i for critic network of team policy. 274 Initialize the Lagrangian multiplier λ and learning rate α . 275 1: Initialize empty data buffer \mathcal{D}_m and \mathcal{D}_t for mixed policy and team policy, respectively. 2: while training step \leq step_{max} do 276 for $step = 1, 2, \ldots, step_{max}$ do 3: 277 4: Initialize empty trajectories lists \mathcal{D}_m and \mathcal{D}_t for mixed policy and team policy, respectively. 278 5: Generate mixed action, team action from mixed policy and team policy, respectively. 279 6: Interact with environment with mixed action. 280 7: Compute the extended TD error U_{E+i} and U_E . 8: Store state, mixed action, reward, termination information and U_{E+i} into \mathcal{T}_m . 281 9: Store state, team action, reward, termination information and U_E into \mathcal{T}_t . 282 10: Incorporate \mathcal{T}_m and \mathcal{T}_t into \mathcal{D}_m and \mathcal{D}_t , respectively. 283 11: end for 284 12: Sample training data from \mathcal{D}_m . 285 Update actor network of mixed policy according to Eq. 10 and Eq. 29. 13: 14: Update critic network of mixed policy according to Eq. 31. 286 15: Sample training data from \mathcal{D}_t . 287 Update actor network of team policy according to the Eq. 13 and Eq. 30. 16: 288 17: Update critic network of team policy according to Eq. 32. 289 18: Update Lagrangian multiplier according to Eq. 28. 290 19: end while **Output:** 291 Learned policy π_{E+i} 292

Similar to the optimization of mixed policy, we transform the optimizing objective as follows (details can be found in Appendix B.2):

$$\lambda J_E\left(\pi_E^i\right) - \hat{J}_{E+i}^{\lambda}\left(\pi_{E+i}^i\right) = \mathbb{E}_{\tau_0^i}\left[V_E^{\pi_E^i}\left(\tau_0^i\right)\right] - \lambda \mathbb{E}_{\pi_{E+i}^i}\left[\sum_{t=0}^{\infty} \gamma^t r_t^{E+i}\right]$$
$$= -\mathbb{E}_{\pi_{E+i}^i}\left[\sum_{t=0}^{\infty} \gamma^t \left(\left(1+\lambda\right)r_t^E + r_t^i - \lambda V_E^{\pi_E^i}\left(\tau_t^i\right) + \gamma \lambda V_E^{\pi_E^i}\left(\tau_{t+1}^i\right)\right)\right]$$
$$:= -\mathbb{E}_{\pi_{E+i}^i}\left[\sum_{t=0}^{\infty} \gamma^t U_E\left(\tau_t^i, a_t^i\right)\right].$$
(12)

Based on the same technique of introducing KL divergence term and policy ratio clip during mixed policy optimization, the final optimization objective of team policy is given by:

$$-\mathbb{E}_{\pi_{E+i}^{i}}\left[\min\left\{\frac{\pi_{E}^{i}\left(a_{t}^{i}|\tau_{t}^{i}\right)}{\pi_{E+i}^{i}\left(a_{t}^{i}|\tau_{t}^{i}\right)}U_{E}, clip\left(\frac{\pi_{E}^{i}\left(a_{t}^{i}|\tau_{t}^{i}\right)}{\pi_{E+i}^{i}\left(a_{t}^{i}|\tau_{t}^{i}\right)}, 1-\epsilon, 1+\epsilon\right)U_{E}\right\}\right] - D_{KL}\left(\pi_{E}^{i}||\pi_{E+i}^{i}\right)\right).$$

$$(13)$$

Optimizing Lagrangian multiplier: To update the Lagrangian multiplier λ , we employ the gradient descent method, considering the optimization objective defined in Eq. 5. The gradient objective for updating λ is as follows (cf. Appendix B.3 for the deriving process):

$$\lambda \leftarrow \lambda - \alpha \mathbb{E}_{\pi_E^i} \left[\sum_{t=0}^{\infty} \gamma^t \min \left\{ \begin{array}{l} \frac{\pi_{E+i}^i(a_t^i | \tau_t^i)}{\pi_E^i(a_t^i | \tau_t^i)} A^{\pi_E^i}\left(\tau_t^i, a_t^i\right), \\ clip\left(\frac{\pi_{E+i}^i(a_t^i | \tau_t^i)}{\pi_E^i(a_t^i | \tau_t^i)}, 1 - \varepsilon, 1 + \varepsilon\right) A^{\pi_E^i}\left(\tau_t^i, a_t^i\right) \end{array} \right\} \right], \quad (14)$$

319 320 321

322 323

293 294

295

305 306

307

314

315

316 317 318

where α is the step size, and $A^{\pi_E^i}(\tau_t^i, a_t^i)$ denotes the TD error for team policy:

$$A^{\pi_{E}^{i}}\left(\tau_{t}^{i}, a_{t}^{i}\right) = r_{t}^{E} + \gamma V_{E}^{\pi_{E}^{i}}\left(\tau_{t+1}^{i}\right) - V_{E}^{\pi_{E}^{i}}\left(\tau_{t}^{i}\right).$$
(15)

326	Mon	Diff cultur	Rule-based Individual Reward				
327	Map Difficulty		IRAT	MAPPO	QMIX	Ours	MAPPO (Sparse)
328	2m_vs_1z	Easy	100.0(0.0)	100.0(0.0)	100.0(0.0)	100.0(0.0)	0.0(0.0)
220	2s3z	Easy	98.9(1.1)	97.2(2.8)	97.2(2.8)	100.0(0.0)	0.0(0.0)
329	3m	Easy	100.0(0.0)	97.3(2.3)	92.6(5.1)	100.0(0.0)	0.0(0.0)
330	1c3s5z	Easy	60.9(13.4)	59.3(13.2)	98.4(1.6)	100.0(0.0)	0.0(0.0)
331	3s_vs_5z	Hard	89.5(4.9)	87.3(4.5)	0.0(0.0)	95.3(1.3)	0.0(0.0)
332	8m_vs_9m	Hard	56.2(15.6)	83.2(3.1)	58.3(24.0)	93.7(6.2)	0.0(0.0)
000	5m_vs_6m	Hard	52.6(20.0)	57.8(8.7)	62.5(6.3)	65(6.7)	0.0(0.0)
333	3s5z	Hard	24.4(22.6)	64.6(8.8)	81.0(12.2)	89.0(9.1)	0.0(0.0)
334	MMM2	Super Hard	19.1(15.3)	3.6(2.5)	71.8(9.9)	62.5(22.0)	0.0(0.0)
335	6h_vs_8z	Super Hard	37.5(35.0)	10.9(6.6)	49.8(36.1)	89.0(2.2)	0.0(0.0)
336	3s5z_vs_3s6z	Super Hard	10.0(0.8)	34.3(20.6)	42.4(49.1)	64.1(35.9)	0.0(0.0)

Table 1: Median winning rate (%) and standard deviation (%) of five MARL algorithms in more than 10 maps of SMAC environment under rule-based reward setting, using 5 random seeds and at most 10M training steps. 325

4.3 IMPLEMENTATION

324

337

338 339

340

341

343

344

345

346

347

348

349 350

351 352

353

354

355

356

357

358

360

361 362

364

365

366

367

368

369

370

CMT is implemented within the CTDE framework, with MAPPO algorithm as the backbone. The implementation involves two types of policies: mixed policy π_{E+i}^i and team policy π_E^i . Each policy is parameterized by separate networks. During policy execution phase, each agent utilizes the actor 342 network of mixed policy to interact with environment, while the policy information for mixed policy and team policy are stored into data buffer \mathcal{D}_m and \mathcal{D}_t , respectively.

When training, CMT algorithm iteratively updates the policies and the Lagrange multiplier λ . In each iteration, it optimizes the mixed policy π_{E+i}^i while keeping the team policy π_E^i fixed. It then updates the team policy π_E^i based on the optimized mixed policies of all agents. The pseudocode of CMT algorithm is summarized in Algorithm 1. For more details about the CMT algorithm implementation, please refer to Appendix C.

5 **EXPERIMENTS**

To evaluate the effectiveness of the proposed CMT algorithm, we benchmark it against state-ofthe-art baselines across three widely recognized multi-agent benchmarks: SMAC (Samvelyan et al., 2019), MPE (Lowe et al., 2017), and GRF (Kurach et al., 2020). Our approach demonstrates superior performance compared to existing SOTA MARL methods such as IRAT, MAPPO, QMIX, MASER, and LAIES, excelling in 25 out of 27 tasks. Comprehensive experimental details and hyperparameter settings are provided in Appendix D.

5.1 EXPERIMENTS ON SMAC

The experiments on SMAC are conducted using two types of individual reward settings:

- **Rule-based Individual Reward**: A sparse reward of 20 is awarded for winning a battle, while a reward of 0 is given otherwise. Additionally, dense individual rewards are allocated to each agent based on the health of team members and enemies. Specifically, an individual reward of 10 is given for each defeated enemy. Meanwhile, a scaled reward is provided according to the agent's remaining health state. Under this reward setting, the CMT algorithm is compared with IRAT (Wang et al., 2022), MAPPO, and the QMIX algorithm. Additionally, MAPPO (Sparse), which trains MAPPO without any individual reward, is included in the experiments to demonstrate the impact of introducing individual rewards. The implementations of IRAT and MAPPO are consistent with the source code in Wang et al. (2022).
- Heuristic Individual Reward: The team reward is set at 20 for a battle win. LAIES (Liu et al., 372 2023) and MASER (Jeon et al., 2022) implement their respective individual rewards, as introduced 373 in Section 2. Since the individual reward in MASER relies on the mixing network of QMIX, 374 which is incompatible with other types of MARL algorithms, our approach employs the same reward setting as LAIES. Under this reward setting, both the proposed and the LAIES algorithms 375 376 are implemented with IPPO, ensuring all experimental details align with the source code in Liu et al. (2023). The MASER algorithm is implemented with QMIX, using the original source code 377 from Jeon et al. (2022).

379

380

381

382

384

386

387

388

389

390

391

392

394

395

397

398

399 400 401



(a) Training curves of five algorithms on four maps of SMAC environments with rule-based individual reward.



(b) Training curves of three algorithms on four maps of SMAC environments with heuristic individual reward.

Figure 2: Partial training curves on four maps of SMAC environments.

Experimental results for the rule-based individual reward setting are presented in Table 1. As il-402 lustrated in Table 1, our approach is the only method that achieves a 100% win rate across all easy 403 maps, and attains the highest win rates on almost all hard and super hard maps. IRAT secures the 404 second-best performance in 6 out of 11 maps. QMIX and MAPPO exhibit varied performances 405 across different maps; In contrast, MAPPO (Sparse), which does not take into account the addi-406 tional individual rewards, fails to develop an effective policy on any map. These results highlight the 407 significance of introducing effective individual rewards and addressing policy inconsistency when 408 individual rewards are incorporated in sparse reward environments. 409

We further select four representative maps to display the training curves of five algorithms in Figure 2a. The figures demonstrate that the proposed CMT algorithm exhibits superior sample efficiency compared to the other four algorithms. CMT requires fewer training steps to converge, and its converged win rate is either higher or comparable to other algorithms. Training curves on all 11 maps are provided in Appendix D.3.

Experimental results for the heuristic individ-415 ual reward setting are presented in Table 2. We 416 evaluated the CMT, MASER, and LAIES al-417 gorithms on two easy maps, two hard maps, 418 and one super hard map. As shown in Table 2, 419 our proposed approach achieves the highest win 420 rate across five maps. This result demonstrates 421 that our approach enhances algorithm perfor-422 mance in both rule-based and heuristic individual reward settings. LAIES outperforms 423 or matches MASER in four out of five maps, 424 which is consistent with the results reported in 425 Liu et al. (2023). 426

Table 2: Median winning rate (%) and standard deviation (%) of three MARL algorithms in 5 maps of SMAC environment under heuristic reward setting, using 5 random seeds and at most 5M training steps.

Man	Heuristic Individual Reward					
wiap	MASER	Ours	LAIES			
2m_vs_1z	0.0(0.0)	100.0(0.0)	90.6(5.6)			
3m	63.8(35.3)	96.8(3.2)	82.8(2.2)			
1c3s5z	0.0(0.0)	72.4(9.9)	49.8(43.2)			
5m_vs_6m	15.2(2.8)	27.1(8.9)	0.0(0.0)			
MMM2	0.0(0.0)	28.2 (15.1)	15.6(6.8)			

We provide training curves for the three algorithms on four selected maps in Figure 2b. It can be
observed that CMT exhibits superior sample efficiency compared to LAIES and MASER. When
comparing algorithm performance between heuristic and rule-based individual reward settings, it
is evident that the algorithm performance under rule-based individual rewards setting significantly
better than that under heuristic individual rewards. This underscores the importance of leveraging
environment knowledge and understanding for effective reward shaping.



Figure 3: Experimental results on 3 scenarios of MPE and 2 scenarios of GRF

5.2 TEST ON MPE

441

442 443

444

In the MPE environment, we compare our approach with IRAT, MAPPO, QMIX, and MAPPO (Sparse) across three types of environments: Spread, Attack, and Predator-Prey.

445 In the Attack environment, three agents collaborate to reach and attack a single landmark, earning 446 a positive reward 20 upon successful completion. Each agent incurs a penalty based on their distance 447 to the landmark, and a penalty -1 when colliding with other agents. As shown in the first column 448 of Figure 3, CMT outperforms all other algorithms, achieving a team reward exceeding 20. IRAT 449 ranks second with a team reward around 7, while MAPPO, QMIX, and MAPPO (Sparse) perform 450 worse, with team rewards below 5. These results indicate that algorithms assisted by individual 451 rewards tend to overlook policy inconsistency, which adversely affects their performance. Although IRAT partially addresses this issue, it fails to achieve consistency between the mixed policy learned 452 from mixed rewards and the team policy derived from team rewards. In contrast, CMT successfully 453 establishes this consistency. 454

455 In the **Spread environment**, four agents collaborate to locate two landmarks, receiving a sparse 456 positive reward 10 when multiple agents simultaneously discover a landmark. Additionally, each 457 agent earns an individual reward based on its minimum distance to undiscovered landmarks. Experimental results are presented in the second column of Figure 3. These results demonstrate that 458 CMT once again achieves the highest team rewards, around 10. IRAT and MAPPO follow closely 459 with a team reward of 8. MAPPO (Sparse) and QMIX perform the worst among all algorithms. 460 These findings further show that CMT can identify the sub-optimal policy trap caused by additional 461 individual rewards. 462

463 In the **Predator-Prey environment**, five predators work together to capture two prey, receiving a sparse positive reward 20 when multiple agents successfully capture the same prey. Each agent 464 also earns an individual reward 5 when it successfully hits a prey. The results in the Predator-465 Prey environment differ slightly from those in the previous two environments. As shown in the 466 third column of Figure 3, the CMT and IRAT algorithms perform similarly, with CMT achieving 467 a slightly higher final team reward. This is because in the Predator-Prey environment, individual 468 rewards play a more significant role than in the other two environments. The IRAT algorithm focuses 469 on improving policy performance by fully utilizing individual rewards, whereas the CMT algorithm 470 also considers the consistency between the team policy and the mixed policy.

471 472

473

5.3 TEST ON GRF

474 We also benchmark CMT in the widely-used GRF environment. We employ a common setting where 475 the team reward is defined as +1 when the team scores and -1 when the team is scored against. The 476 individual reward is based on the default checkpoint, where an agent receives a reward of 0.1 for 477 possessing the football in a region near the goal. We present the results for the two popular academy tasks, 3_{vs_1} with the keeper and academy counter-attack (easy) (cf. the experimental settings in 478 Appendix D.4). In the GRF environment, we compare the CMT algorithm with IRAT, MAPPO, 479 and MAPPO (Sparse) algorithms. We did not test QMIX because the GRF environment does not 480 provide global status information. 481

As shown in the fourth and fifth column of Figure 3, the CMT algorithm achieves the highest win
 rate, exceeding 60% on both maps. This result demonstrates that CMT can deliver superior per formance in environments requiring full collaboration among agents, such as passing and off-ball
 movement. Even in sparse reward environments, CMT demonstrates excellent performance in de veloping effective strategies through the utilization of mixed rewards.

Table 3: Average team reward of CMT and five variants, i.e., CMT-TD(w), CMT-SP(w), CMT-KL(w), CMT($3\times$), CMT($5\times$) and CMT(RD) on MPE environment using 5 random seeds and at most 2M training steps.

-	Scenario	CMT	CMT-TD(w)	CMT-SP(w)	CMT-KL(w)	$CMT(3\times)$	$CMT(5\times)$	CMT(RD)
	Spread	9.8(3.9)	7.2(2.3)	9.7(2.5)	7.5(2.5)	9.5(2.0)	12.2(1.7)	0.0(0.0)
	Attack	23.7(3.8)	8.5(4.0)	13.0(6.0)	11.0(5.5)	18.5(6.5)	17.2(7.5)	0.0(0.0)
	Predator-Pray	112.5(17.5)	14.0(6.0)	106.2(26.5)	22.5(8.0)	80.5(22.5)	107.5(20.0)	7.8(6.0)

5.4 ABLATION STUDIES

495 It is noteworthy that CMT incorporates three crucial design elements: the extended TD error, policy 496 approximation based on the *similarity assumption*, and optimization objective reconstruction with 497 KL divergence. To assess the impact of all components on CMT's performance, we conduct an 498 ablation study in MPE environments. We conduct three ablation studies: i) CMT-TD(w) as CMT 499 with the extended TD error replaced by a standard TD error (defined as Eq. 15) used by MAPPO 500 and IRAT for policy optimization; ii) CMT-SP(w) is CMT with the policy approximation removed by initially using distinct parameters between mixed and team policy networks; iii) CMT-KL(w) 501 represents CMT without the KL terms-based reconstruction module. 502

503 Table 3 demonstrates that the inclusion of the extended TD error exerts the most significant impact, 504 followed by the KL divergence-based optimization objective reconstruction model, with the policy 505 approximation technique showing the least influence on algorithm performance. This finding aligns 506 our core innovation outlined in Section 1. The extended TD error, which originates from resolving the policy consistency constraint in the Lagrangian dual problem and the incorporation of mixed 507 rewards, notably contribute to CMT's enhanced performance. Moreover, the minimal disparity be-508 tween the outcomes of CMT-SP(w) and CMT suggests that eliminating the similarity assumption 509 does not dramatically affect algorithm performance. 510

The influence of scaling individual rewards on the algorithm's performance is also investigated. In Table 3, $CMT(3\times)$ denotes the CMT developed with individual rewards amplified three times, while CMT(5×) indicates the CMT developed with individual rewards amplified five times. Comparing CMT(3×) and CMT(5×) with baseline CMT across three scenarios, we observe that the performance variation remains within 30%, indicating CMT's robustness to individual reward scaling.

Finally, we examine CMT's robustness to reward design by testing with random individual rewards sampled from [-1, 1] (denoted as CMT(RD) in Table 3). The results show that with inappropriate individual rewards, CMT(RD) performs similarly to MAPPO(Sparse), both showing limited effectiveness. This indicates that even with poorly designed individual rewards, our approach maintains a performance floor equivalent to methods without additional rewards. Additional ablation studies investigating the impact of the initial Lagrangian multiplier λ on CMT's performance can be found in Appendix D.3.

523 524

494

6 CONCLUSION AND FUTURE WORK

525 526

527 In this paper, we focus on addressing the challenge of deviation in optimal policies in MARL due 528 to the introduction of individual rewards. To tackle this problem, we propose a novel multi-agent 529 constrained policy optimization procedure, which maximizes the cumulative rewards while ensur-530 ing the consistency between the team policy and the mixed policy learned from the sum of team and 531 individual rewards. Leveraging the min-max dual objective presented in the constrained policy optimization model, our approach iteratively updates the mixed and the team policies under the proposed 532 policy consistency constraint. Experimental results validate that CMT successfully overcomes the 533 policy inconsistency issue, gaining superior performance across 27 tasks on MPE, SMAC, and GRF 534 environments, compared to SOTA MARL algorithms. 535

536 One limitation of the proposed approach lies that we utilize the *similarity assumption* when simpli-537 fying the optimization objectives. Although we can develop policy networks with the same parame-538 ters to make this assumption hold during implementation, there may still be divergence between the 539 mixed and team policies during the policy learning process. Relaxing this assumption leads to an 539 interesting direction that is worth further exploration in future work.

540 REFERENCES

548

569

570

571

- Marc Bellemare, Sriram Srinivasan, Georg Ostrovski, Tom Schaul, David Saxton, and Remi Munos.
 Unifying count-based exploration and intrinsic motivation. *Advances in Neural Information Processing Systems (NeurIPS)*, 29, 2016.
- Eric Chen, Zhang-Wei Hong, Joni Pajarinen, and Pulkit Agrawal. Redeeming intrinsic rewards via
 constrained optimization. *Advances in Neural Information Processing Systems (NeurIPS)*, 35:
 4996–5008, 2022.
- Sirui Chen, Zhaowei Zhang, Yaodong Yang, and Yali Du. Stas: Spatial-temporal return decompo sition for solving sparse rewards problems in multi-agent reinforcement learning. In *Proceedings* of the AAAI Conference on Artificial Intelligence, volume 38, pp. 17337–17345, 2024.
- ⁵⁵² Dongsheng Ding, Kaiqing Zhang, Tamer Basar, and Mihailo Jovanovic. Natural policy gradient
 ⁵⁵³ primal-dual method for constrained markov decision processes. *Advances in Neural Information* ⁵⁵⁴ *Processing Systems (NeurIPS)*, 33:8378–8390, 2020.
- Wenbo Du, Tong Guo, Jun Chen, Biyue Li, Guangxiang Zhu, and Xianbin Cao. Cooperative pursuit of unauthorized uavs in urban airspace via multi-agent reinforcement learning. *Transportation Research Part C: Emerging Technologies*, 128:103122, 2021.
- Xunhan Hu, Jian Zhao, Wengang Zhou, Ruili Feng, and Houqiang Li. Differ: Decomposing individ ual reward for fair experience replay in multi-agent reinforcement learning. *Advances in Neural Information Processing Systems (NeurIPS)*, 36, 2023.
- 562
 563
 564
 564
 565
 2022.
- Jeewon Jeon, Woojun Kim, Whiyoung Jung, and Youngchul Sung. Maser: Multi-agent reinforce ment learning with subgoals generated from experience replay buffer. In *International Conference on Machine Learning (ICML)*, pp. 10041–10052. PMLR, 2022.
 - Sham Kakade and John Langford. Approximately optimal approximate reinforcement learning. In *Proceedings of the Nineteenth International Conference on Machine Learning (ICML)*, pp. 267–274, 2002.
- Karol Kurach, Anton Raichuk, Piotr Stańczyk, Michał Zając, Olivier Bachem, Lasse Espeholt, Carlos Riquelme, Damien Vincent, Marcin Michalski, Olivier Bousquet, et al. Google research football: A novel reinforcement learning environment. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pp. 4501–4510, 2020.
- 577 Chenghao Li, Tonghan Wang, Chengjie Wu, Qianchuan Zhao, Jun Yang, and Chongjie Zhang. Cel 578 ebrating diversity in shared multi-agent reinforcement learning. *Advances in Neural Information* 579 *Processing Systems (NeurIPS)*, 34:3991–4002, 2021.
- Boyin Liu, Zhiqiang Pu, Yi Pan, Jianqiang Yi, Yanyan Liang, and Du Zhang. Lazy agents: a new perspective on solving sparse reward problem in multi-agent reinforcement learning. In *International Conference on Machine Learning (ICML)*, pp. 21937–21950. PMLR, 2023.
- Iou-Jen Liu, Unnat Jain, Raymond A Yeh, and Alexander Schwing. Cooperative exploration for
 multi-agent deep reinforcement learning. In *International Conference on Machine Learning* (*ICML*), pp. 6826–6836. PMLR, 2021.
- Zuxin Liu, Baiming Chen, Hongyi Zhou, Guru Koushik, Martial Hebert, and Ding Zhao. Mapper: Multi-agent path planning with evolutionary reinforcement learning in mixed dynamic environments. In 2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), pp. 11748–11754. IEEE, 2020.
- Ryan Lowe, Yi I Wu, Aviv Tamar, Jean Harb, OpenAI Pieter Abbeel, and Igor Mordatch. Multiagent actor-critic for mixed cooperative-competitive environments. *Advances in Neural Information Processing Systems (NeurIPS)*, 30, 2017.

594 Heechang Ryu, Hayong Shin, and Jinkyoo Park. Remax: Relational representation for multi-agent 595 exploration. In Proceedings of the 21st International Conference on Autonomous Agents and 596 Multiagent Systems (AAMAS), pp. 1137–1145, 2022. 597 Mikayel Samvelyan, Tabish Rashid, Christian Schroeder De Witt, Gregory Farquhar, Nantas 598 Nardelli, Tim GJ Rudner, Chia-Man Hung, Philip HS Torr, Jakob Foerster, and Shimon Whiteson. The starcraft multi-agent challenge. arXiv preprint arXiv:1902.04043, 2019. 600 601 John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region policy optimization. In International Conference on Machine Learning (ICML), pp. 1889–1897. 602 PMLR, 2015. 603 604 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy 605 optimization algorithms. arXiv preprint arXiv:1707.06347, 2017. 606 Alexander L Strehl and Michael L Littman. An analysis of model-based interval estimation for 607 markov decision processes. Journal of Computer and System Sciences, 74(8):1309–1331, 2008. 608 609 Chuangchuang Sun, Macheng Shen, and Jonathan P How. Scaling up multiagent reinforcement 610 learning for robotic systems: Learn an adaptive sparse communication graph. In 2020 IEEE/RSJ 611 International Conference on Intelligent Robots and Systems (IROS), pp. 11755–11762. IEEE, 612 2020. 613 Shaoqi Sun and Kele Xu. Temporal inconsistency-based intrinsic reward for multi-agent reinforce-614 ment learning. In 2023 International Joint Conference on Neural Networks (IJCNN), pp. 1–7. 615 IEEE, 2023. 616 Andrew Tittaferrante and Abdulsalam Yassine. Multiadvisor reinforcement learning for multiagent 617 multiobjective smart home energy control. IEEE Transactions on Artificial Intelligence, 3(4): 618 581-594, 2021. 619 620 Li Wang, Yupeng Zhang, Yujing Hu, Weixun Wang, Chongjie Zhang, Yang Gao, Jianye Hao, Tangjie 621 Lv, and Changjie Fan. Individual reward assisted multi-agent reinforcement learning. In Interna-622 tional Conference on Machine Learning (ICML), pp. 23417–23432. PMLR, 2022. 623 Tonghan Wang, Jianhao Wang, Yi Wu, and Chongjie Zhang. Influence-based multi-agent explo-624 ration. arXiv preprint arXiv:1910.05512, 2019. 625 Pei Xu, Junge Zhang, and Kaiqi Huang. Exploration via joint policy diversity for sparse-reward 626 multi-agent tasks. In Proceedings of the Thirty-Second International Joint Conference on Artifi-627 cial Intelligence (IJCAI), pp. 326–334, 2023a. 628 629 Pei Xu, Junge Zhang, Qiyue Yin, Chao Yu, Yaodong Yang, and Kaiqi Huang. Subspace-aware ex-630 ploration for sparse-reward multi-agent tasks. In Proceedings of the AAAI Conference on Artificial 631 Intelligence, volume 37, pp. 11717–11725, 2023b. 632 Pei Xu, Junge Zhang, and Kaiqi Huang. Population-based diverse exploration for sparse-reward 633 multi-agent tasks. In Proceedings of the Thirty-Third International Joint Conference on Artificial 634 Intelligence, pp. 283-291, 2024. 635 Ziming Yan and Yan Xu. A multi-agent deep reinforcement learning method for cooperative load 636 frequency control of a multi-area power system. *IEEE Transactions on Power Systems*, 35(6): 637 4599-4608, 2020. 638 639 Mingqi Yuan, Bo Li, Xin Jin, and Wenjun Zeng. Automatic intrinsic reward shaping for exploration 640 in deep reinforcement learning. In International Conference on Machine Learning (ICML), pp. 641 40531-40554. PMLR, 2023. 642 Qizhen Zhang, Chris Lu, Animesh Garg, and Jakob Foerster. Centralized model and exploration 643 policy for multi-agent rl. arXiv preprint arXiv:2107.06434, 2021a. 644 645 Yang Zhang, Bo Tang, Qingyu Yang, Dou An, Hongyin Tang, Chenyang Xi, Xueying Li, and Feiyu Xiong. Bcorle (λ): An offline reinforcement learning and evaluation framework for coupons al-646 location in e-commerce market. Advances in Neural Information Processing Systems, 34:20410-647 20422, 2021b.

648	Zevu Zhu and Huijing Zhao. A survey of deep rl and il for autonomous driving policy learning
649 650	<i>IEEE Transactions on Intelligent Transportation Systems</i> , 23(9):14043–14065, 2021.
651	
652	
653	
654	
655	
656	
657	
658	
659	
660	
661	
662	
663	
664	
665	
666	
667	
668	
669	
670	
671	
672	
673	
674	
676	
677	
678	
679	
680	
681	
682	
683	
684	
685	
686	
687	
688	
689	
690	
691	
692	
693	
694	
695	
696	
697	
698	
700	
700	
/01	



706



A VISUAL COMPARISON OF POLICY LEARNING APPROACHES

722 LAIES (Liu et al., 2023))

Figure 4: Comparison of policy convergence across different approaches. Blue and orange lines represent
 learned and team policies, respectively. While approaches in (a) and (b) cannot guarantee convergence to the
 optimal policy, our approach (c) ensures alignment between learned and optimal policies in terms of team
 rewards.

728

729 Figure 4 comprehensively illustrates that our approach fully resolves the policy inconsistency prob-730 lem compared to existing research. In Figure 4, the learning trajectories of learned and team policies 731 are depicted as blue and orange lines, respectively. In Figure 4a, prior research (e.g., MASER (Jeon 732 et al., 2022) and LAIES (Liu et al., 2023)) incorporates individual rewards into MARL while over-733 looking the policy inconsistency issue. Consequently, the blue line remains parallel to the orange line, and the optimal team policy fails to develop (shown by dashed lines). In Figure 4b, IRAT 734 mitigates the policy inconsistency problem but does not guarantee convergence of the learned pol-735 icy to the optimal team policy. As a result, while blue and orange lines gradually approach each 736 other from their starting points and may eventually intersect, the team policy line deviates from its 737 intended direction (red line). In Figure 4c, our approach is the **first** to introduce the consistency 738 policy constraint, enforcing consistency between the mixed policy and optimal team policy's team 739 rewards. Consequently, the direction of team policy remains unchanged, and the blue and orange 740 lines intersect precisely at the optimal team policy point. 741

742 743

744 745

746 747

748 749

B COMPLETE MATHEMATICAL DERIVATION

B.1 OPTIMIZATION OBJECTIVE OF MIXED POLICY

The original optimization objective for mixed policy defined in Eq. 6 is:

$$\max_{\substack{\pi_{E+i}^{i}}} \hat{J}_{E+i}^{\lambda} \left(\pi_{E+i}^{i}\right) - \lambda J_{E} \left(\pi_{E}^{i}\right).$$
(16)

750 751 752

To address the data inefficiency issue in sampling policies π_{E+i} and π_E simultaneously, we make an approximation of optimizing objective in training process of mixed policy network.

First, we expand the objective function as follows:

$$\hat{J}_{E+i}^{\lambda}\left(\pi_{E+i}^{i}\right) - \lambda J_{E}\left(\pi_{E}^{i}\right) = \mathbb{E}_{\tau_{0}}\left[V_{E+i}^{\pi_{E+i}^{i},\lambda}\left(\tau_{0}^{i}\right)\right] - \lambda \mathbb{E}_{\pi_{E}^{i}}\left[\sum_{t=0}^{\infty}\gamma^{t}r_{t}^{E}\right]$$
$$= -\left[\mathbb{E}_{\tau_{0}^{i}}\left[-V_{E+i}^{\pi_{E+i}^{i},\lambda}\left(\tau_{0}^{i}\right)\right] + \lambda \mathbb{E}_{\pi_{E}^{i}}\left[\sum_{t=0}^{\infty}\gamma^{t}r_{t}^{E}\right]\right]$$
$$= -\mathbb{E}_{\pi_{E}^{i}}\left[-V_{E+i}^{\pi_{E+i}^{i},\lambda}\left(\tau_{0}^{i}\right) + \sum_{t=0}^{\infty}\lambda\gamma^{t}r_{t}^{E}\right].$$
(17)

(18)

 Since the policy π_E^i has no influence on $V_{E+i}^{\pi_{E+i}^i,\lambda}$, we can merge the two components in the second equation in Eq. 17, yielding the last equation in Eq. 17.

For simplicity, we let $V_t = V_{E+i}^{\pi_{E+i}^i,\lambda}(\tau_t^i)$ and $r_t = \lambda r_t^E$. We can expand the expression within the expectation in Eq. 16 as follows:

 $-V_{0} + \sum_{t=0}^{\infty} \gamma^{t} r_{t} = (r_{0} - V_{0} + \gamma V_{1}) + \gamma (r_{1} - V_{1} + \gamma V_{2}) + \gamma^{2} (r_{2} - V_{2} + \gamma V_{3}) + \cdots$ $=\sum_{t=0}^{\infty}\gamma^{t}\left(r_{t}-V_{t}+\gamma V_{t}\right)$

 $=\sum_{t=0}^{\infty}\gamma^{t}\left(\lambda r_{t}^{E}-V_{E+i}^{\pi_{E+i}^{i},\lambda}\left(\tau_{t}^{i}\right)+\gamma V_{E+i}^{\pi_{E+i}^{i},\lambda}\left(\tau_{t+1}^{i}\right)\right)$

$$:=\sum_{t=0}^{\infty}\gamma^{t}U_{E+i}\left(\tau_{t}^{i},a_{t}^{i}\right).$$

Therefore, the optimization objective can be rewritten as:

$$\begin{array}{ll}
\begin{aligned}
\hat{J}_{E+i}^{\lambda}\left(\pi_{E+i}^{i}\right) - \lambda J_{E}\left(\pi_{E}^{i}\right) &= -\mathbb{E}_{\pi_{E}^{i}}\left[\sum_{t=0}^{\infty}\gamma^{t}U_{E+i}\left(\tau_{t}^{i},a_{t}^{i}\right)\right] \\
\\
\begin{array}{l}
= -\sum_{t=0}^{\infty}\sum_{\tau^{i}\in(\{Z^{i}\}\times\{A^{i}\})}\gamma P\left(\tau_{t}=\tau|\rho_{0},\pi_{E}^{i}\right)\sum_{a\in A}\pi_{E}^{i}\left(a_{t}^{i}|\tau_{t}^{i}\right)U_{E+i}\left(\tau_{t}^{i},a_{t}^{i}\right) \\
= -\sum_{\tau^{i}\in(\{Z^{i}\}\times\{A^{i}\})}\eta^{\pi_{E}^{i},\gamma}\left(\tau^{i}\right)\sum_{a\in A}\pi_{E}^{i}\left(a_{t}^{i}|\tau_{t}^{i}\right)U_{E+i}\left(\tau_{t}^{i},a_{t}^{i}\right) \\
= -\sum_{\tau^{i}\in(\{Z^{i}\}\times\{A^{i}\})}\eta^{\pi_{E}^{i},\gamma}\left(\tau^{i}\right)\sum_{a\in A}\pi_{E}^{i}\left(a_{t}^{i}|\tau_{t}^{i}\right)U_{E+i}\left(\tau_{t}^{i},a_{t}^{i}\right) \\
\end{aligned}$$
(19)

 where $d_{\rho_0}^{\pi_E^i,\gamma}(\tau^i) = \sum_{t=0}^{\infty} \gamma^t P\left(\tau_t^i = \tau^i | \rho_0, \pi_E^i\right)$ denotes the discounted observation-action frequency under policy π_E^i with the initial observation-action distribution ρ_0 .

To mitigate the data inefficiency resulting from simultaneously sampling policies π_{E+i} and π_E , as per the similarity assumption, we substitute π_E^i with π_{E+i}^i and approximate the optimization objectives as follows:

$$\hat{J}_{E+i}^{\lambda}\left(\pi_{E+i}^{i}\right) - \lambda J_{E}\left(\pi_{E}^{i}\right) = -\sum_{\tau^{i} \in \left(\{Z^{i}\} \times \{A^{i}\}\right)} d_{\rho_{0}}^{\pi_{E}^{i},\gamma}\left(\tau^{i}\right) \sum_{a \in A} \pi_{E+i}^{i}\left(a_{t}^{i}|\tau_{t}^{i}\right) U_{E+i}\left(\tau_{t}^{i},a_{t}^{i}\right).$$
(20)

To further simplify the computation process, we introduce the KL divergence between the mixed policy and the team policy while preserving the equality between the original and converted opti810 mization objectives:

$$\begin{aligned} \hat{J}_{E+i}^{\lambda} \left(\pi_{E+i}^{i} \right) - \lambda J_{E} \left(\pi_{E}^{i} \right) &= - \left[\sum_{\tau^{i} \in (\{Z^{i}\} \times \{A^{i}\})} d_{\rho_{0}}^{\pi_{E,\gamma}^{i}} \sum_{a \in A} \pi_{E+i}^{i} \left(a_{t}^{i} | \tau_{t}^{i} \right) U_{E+i} \left(\tau_{t}^{i}, a_{t}^{i} \right) \right] \\ &+ D_{KL} \left(\pi_{E+i}^{i} \parallel \pi_{E}^{i} \right) - D_{KL} \left(\pi_{E+i}^{i} \parallel \pi_{E}^{i} \right) \\ &= - \left[\sum_{\tau^{i} \in (\{Z^{i}\} \times \{A^{i}\})} d_{\rho_{0}}^{\pi_{E,\gamma}^{i}} \sum_{a \in A} \pi_{E}^{i} \left(a_{t}^{i} | \tau_{t}^{i} \right) \frac{\pi_{E+i}^{i} \left(a_{t}^{i} | \tau_{t}^{i} \right)}{\pi_{E}^{i} \left(a^{i} | \tau^{i} \right)} U_{E+i} \left(\tau_{t}^{i}, a_{t}^{i} \right) \right] \\ &+ D_{KL} \left(\pi_{E+i}^{i} \parallel \pi_{E}^{i} \right) - D_{KL} \left(\pi_{E+i}^{i} \parallel \pi_{E}^{i} \right) \\ &+ D_{KL} \left(\pi_{E+i}^{i} \parallel \pi_{E}^{i} \right) - D_{KL} \left(\pi_{E+i}^{i} \parallel \pi_{E}^{i} \right) \\ &= - \left[\mathbb{E}_{\pi_{E}} \left[\frac{\pi_{E+i}^{i} \left(a_{t}^{i} | \tau_{t}^{i} \right)}{\pi_{E}^{i} \left(a_{t}^{i} | \tau_{t}^{i} \right)} \right] - D_{KL} \left(\pi_{E+i}^{i} \parallel \pi_{E}^{i} \right) \right] \\ &- D_{KL} \left(\pi_{E+i}^{i} \parallel \pi_{E}^{i} \right). \end{aligned}$$

Leveraging the clip technique presented in Schulman et al. (2017), we finally convert the objective to

$$\hat{J}_{E+i}^{\lambda}\left(\pi_{E+i}^{i}\right) - \lambda J_{E}\left(\pi_{E}^{i}\right) = -\mathbb{E}_{\pi_{E}^{i}}\left[\min\left\{\begin{array}{l} \min\left\{\begin{array}{c} \frac{\pi_{E+i}^{i}\left(a_{t}^{i}|\tau_{t}^{i}\right)}{\pi_{E}^{i}\left(a_{t}^{i}|\tau_{t}^{i}\right)}U_{E+i}\left(\tau_{t}^{i},a_{t}^{i}\right),\\ clip\left(\frac{\pi_{E+i}^{i}\left(a_{t}^{i}|\tau_{t}^{i}\right)}{\pi_{E}^{i}\left(a_{t}^{i}|\tau_{t}^{i}\right)},1-\epsilon,1+\epsilon\right)U_{E+i}\left(\tau_{t}^{i},a_{t}^{i}\right)\right\}\right] \\ -D_{KL}\left(\pi_{E+i}^{i}||\pi_{E}^{i}\right).$$
(22)

The introduction of the KL-divergence term between the mixed policy and team policy in Eq. 22 not only assists in acquiring the entire optimization objective but also helps to make the *similarity assumption* valid.

B.2 OPTIMIZATION OBJECTIVE OF TEAM POLICY

The original optimization objective for the team policy, as defined in Eq. 6, is:

$$\max_{\pi_E^i} \lambda J_E\left(\pi_E^i\right) - \hat{J}_{E+i}^\lambda\left(\pi_{E+i}^i\right).$$
(23)

Similar to optimizing the mixed policy, we rewrite the optimization objective for the team policy:

$$\lambda J_E\left(\pi_E^i\right) - \hat{J}_{E+i}^{\lambda}\left(\pi_{E+i}^i\right) = \mathbb{E}_{\tau_0^i}\left[V_E^{\pi_E^i}\left(\tau_0^i\right)\right] - \lambda \mathbb{E}_{\pi_{E+i}^i}\left[\sum_{t=0}^{\infty} \gamma^t r_t^{E+i}\right]$$
$$= -\left[\mathbb{E}_{\tau_0^i}\left[-V_E^{\pi_E^i}\left(\tau_0^i\right)\right] + \lambda \mathbb{E}_{\pi_{E+i}^i}\left[\sum_{t=0}^{\infty} \gamma^t r_t^{E+i}\right]\right]$$
$$= -\mathbb{E}_{\pi_{E+i}^i}\left[-V_E^{\pi_E^i}\left(\tau_0^i\right) + \sum_{t=0}^{\infty} \lambda \gamma^t r_t^{E+i}\right]$$
$$= -\mathbb{E}_{\pi_{E+i}^i}\left[\sum_{t=0}^{\infty} \gamma^t \left((1+\lambda) r_t^E + r_t^i - \lambda V_E^{\pi_E^i}\left(\tau_t^i\right) + \gamma \lambda V_E^{\pi_E^i}\left(\tau_{t+1}^i\right)\right)\right]$$
$$:= -\mathbb{E}_{\pi_{E+i}^i}\left[\sum_{t=0}^{\infty} \gamma^t U_E\left(\tau_t^i, a_t^i\right)\right].$$
(24)

By applying the same technique used in optimizing the mixed policy, we obtain the optimization objective for the team policy:

$$\lambda J_{E}(\pi_{E}^{i}) - \hat{J}_{E+i}^{\lambda}(\pi_{E+i}^{i}) = -\mathbb{E}_{\pi_{E+i}^{i}} \left[\min \left\{ \begin{array}{c} \frac{\pi_{E}^{i}(a_{t}^{i}|\tau_{t}^{i})}{\pi_{E+i}^{i}(a_{t}^{i}|\tau_{t}^{i})} U_{E}(\tau_{t}^{i}, a_{t}^{i}), \\ clip\left(\frac{\pi_{E}^{i}(a_{t}^{i}|\tau_{t}^{i})}{\pi_{E+i}^{i}(a_{t}^{i}|\tau_{t}^{i})}, 1 - \epsilon, 1 + \epsilon\right) U_{E}(\tau_{t}^{i}, a_{t}^{i}) \end{array} \right\} \right] \\ - D_{KL}(\pi_{E}^{i}||\pi_{E+i}^{i}).$$

B.3 Optimizing the Lagrangian multiplier λ

To update Lagrangian multiplier λ , we start by deriving the objective with respect to λ based on the optimization objective defined in Equation 5:

$$\nabla_{\lambda} \left(\hat{J}_{E+i}^{\lambda} \left(\pi_{E+i}^{i} \right) - \lambda J_{E} \left(\pi_{E}^{i} \right) \right) = J_{E} \left(\pi_{E+i}^{i} \right) - J_{E} \left(\pi_{E}^{i} \right).$$
(26)

(25)

、 ¬

To approximate the gradient, we employ a technique presented in the PPO algorithm (Schulman et al., 2017). The gradient can be lower bounded as follows:

$$J_E\left(\pi_{E+i}^{i}\right) - J_E\left(\pi_{E}^{i}\right) \ge \mathbb{E}_{\pi_E^{i}} \left[\sum_{t=0}^{\infty} \gamma^t \min \left\{ \begin{array}{l} \frac{\pi_{E+i}^{i}\left(a_t^{i}|\tau_t^{i}\right)}{\pi_E^{i}\left(a_t^{i}|\tau_t^{i}\right)} A^{\pi_E^{i}}\left(\tau_t^{i}, a_t^{i}\right), \\ clip\left(\frac{\pi_{E+i}^{i}\left(a_t^{i}|\tau_t^{i}\right)}{\pi_E^{i}\left(a_t^{i}|\tau_t^{i}\right)}, 1-\varepsilon, 1+\varepsilon\right) A^{\pi_E^{i}}\left(\tau_t^{i}, a_t^{i}\right) \end{array} \right\} \right],$$

$$(27)$$

where the advantage function $A^{\pi_E^i}(\tau_t^i, a_t^i)$ is defined as $A^{\pi_E^i}(\tau_t^i, a_t^i) = r_t^E + \gamma V_E^{\pi_E^i}(\tau_{t+1}^i) - V_E^{\pi_E^i}(\tau_{t+1}^i)$ $V_E^{\pi_E^i}\left(\tau_t^i\right).$

Finally, the approximate gradient update step for the Lagrangian multiplier is given by:

$$\lambda \leftarrow \lambda - \alpha \mathbb{E}_{\pi_E^i} \left[\sum_{t=0}^{\infty} \gamma^t \min \left\{ \begin{array}{c} \frac{\pi_{E+i}^i(a_t^i | \tau_t^i)}{\pi_E^i(a_t^i | \tau_t)} A^{\pi_E^i}\left(\tau_t^i, a_t^i\right), \\ clip\left(\frac{\pi_{E+i}^i(a_t^i | \tau_t^i)}{\pi_E^i(a_t^i | \tau_t^i)}, 1 - \varepsilon, 1 + \varepsilon\right) A^{\pi_E^i}\left(\tau_t^i, a_t^i\right) \end{array} \right\} \right], \quad (28)$$

with α denoting the step size.

С IMPLEMENTATION DETAILS

C.1 AUXILIARY OBJECTIVES

With the MAPPO algorithm as the backbone, CMT adds an auxiliary objective from MAPPO into the entire optimization objectives of mixed policy and policy, respectively.

When updating the mixed policy, an auxiliary objective from MAPPO is added to maximize $J_{E+i}\left(\pi_{E+i}^{i}\right)$ for agent *i* as follows:

$$\max \mathbb{E}\left[\min\left\{\frac{\pi_{E+i}^{i}\left(a_{\tau}^{i}|o_{\tau}^{i}\right)}{\pi_{E+i}^{i,old}\left(a_{\tau}^{i}|o_{\tau}^{i}\right)}A_{E+i}^{i,old}\left(a_{\tau}^{i}|o_{\tau}^{i}\right), clip\left(\frac{\pi_{E+i}^{i}\left(a_{\tau}^{i}|o_{\tau}^{i}\right)}{\pi_{E+i}^{i,old}\left(a_{\tau}^{i}|o_{\tau}^{i}\right)}, 1-\epsilon, 1+\epsilon\right)A_{E+i}^{i,old}\left(a_{\tau}^{i}|o_{\tau}^{i}\right)\right\}\right]$$

$$(29)$$

where the advantage function for mixed policy $A^{\pi_{E+i}^i}(\tau_t^i, a_t^i)$ is defined as $A^{\pi_{E+i}^i}(\tau_t^i, a_t^i) = \hat{r}_t^i + \hat{r}_t^i$ $\gamma V_{E+i}^{\pi_{E+i}^{i}} \left(\tau_{t+1}^{i} \right) - V_{E+i}^{\pi_{E+i}^{i}} \left(\tau_{t}^{i} \right).$

Similarly, when updating the team policy, an auxiliary objective is added to maximize $\hat{J}_E(\pi_E^i)$ for agent i as follows:

920 921

923

924 925

926

932 933

934

935

940 941

942 943

944

954

919
920
$$\max \mathbb{E}\left[\min\left\{\frac{\pi_{E}^{i}\left(a_{\tau}^{i}|o_{\tau}^{i}\right)}{\pi_{E}^{i,old}\left(a_{\tau}^{i}|o_{\tau}^{i}\right)}A_{E}^{i,old}\left(a_{\tau}^{i}|o_{\tau}^{i}\right), clip\left(\frac{\pi_{E}^{i}\left(a_{\tau}^{i}|o_{\tau}^{i}\right)}{\pi_{E}^{i,old}\left(a_{\tau}^{i}|o_{\tau}^{i}\right)}, 1-\epsilon, 1+\epsilon\right)A_{E}^{i,old}\left(a_{\tau}^{i}|o_{\tau}^{i}\right)\right\}\right]$$
921
922
(30)

C.2 THE UPDATE OF CRITIC NETWORK

The critic network of the mixed policy for agent *i* is updated in the direction of minimizing the loss function,

$$L\left(\phi_{E+i}^{i}\right) = \mathbb{E}\left[\max\left[\begin{array}{c}\left(V_{\phi_{E+i}^{i}}^{i}\left(\tau_{t}\right) - R_{t}^{E+i}\right)^{2}, \\ \left(clip\left(V_{\phi_{E+i}^{i}}^{i}\left(\tau_{t}\right), V_{\phi_{E+i}^{i}}^{i}\left(\tau_{t}\right) - \epsilon, V_{\phi_{E+i}^{i}}^{i}\left(\tau_{t}\right) + \epsilon\right) - R_{t}^{E+i}\right)^{2}\right]\right],$$
(31)

where R_t^{E+i} denotes the discounted reward-to-go of agent i's mixed reward, which is the cumulative reward obtained by agent *i* in the mixed policy.

The critic network of team policy for agent i is updated in the direction of minimizing the loss function,

$$L\left(\phi_{E}^{i}\right) = \mathbb{E}\left[\max\left[\left(V_{\phi_{E}^{i}}^{i}\left(\tau_{t}\right) - R_{t}^{E}\right)^{2}, \left(clip\left(V_{\phi_{E}^{i}}^{i}\left(\tau_{t}\right), V_{\phi_{E}^{i}}^{i}\left(\tau_{t}\right) - \epsilon, V_{\phi_{E}^{i}}^{i}\left(\tau_{t}\right) + \epsilon\right) - R_{t}^{E}\right)^{2}\right]\right],\tag{32}$$

where R_t^E denotes the discounted reward-to-go of agent *i*'s team reward.

EXPERIMENT DETAILS D

D.1 CODE BASE

945 Our approach is implemented with two versions: CMT with MAPPO and CMT with IPPO. When 946 comparing our method with IRAT, MAPPO and QMIX algorithms, since IRAT is implemented with 947 MAPPO algorithm, we implement our approach with MAPPO, and all implementing details are kept 948 consistent with IRAT (Wang et al., 2022). When comparing our method with LAIES and MASER 949 algorithms, since LAIES is implemented with IPPO algorithm, we implement our approach with IPPO, and all implementing details are kept consistent with LAIES (Liu et al., 2023). We sincerely 950 thanks the authors of IRAT and LAIES research for their excellent work producing the codebase. 951

- 952 953
 - D.2 TEST ON MPE ENVIRONMENTS

The experiments were conducted on a computing platform equipped with an AMD Ryzen 9 7950X 955 CPU and a Nvidia 4090 GPU with 96GB of memory. The common parameters for the CMT algo-956 rithm and baselines in the MPE environment are summarized in Table 4. The parameters for the 957 IRAT baseline are identical to those reported in Wang et al. (2022). The specific parameters used 958 for the CMT algorithm in the MPE environment are listed in Table 5. 959

To deeply investigate the CMT algorithm, we conduct more ablation studies to examine the impact 960 of the initial Lagrangian multiplier value selection. Using the Predator-Prey environment of MPE 961 as examples, Figure 5 reveals that the selection of initial values for the Lagrangian multiplier has a 962 significant impact on the algorithm's performance. Notably, positive values of λ tend to yield better 963 policy performance compared to negative values of λ , providing valuable guidance for applying the 964 CMT algorithm in real-world scenarios.

965

966 D.3 TEST ON SMAC ENVIRONMENTS 967

968 We employ the SC2.4.10 version as the benchmark to evaluate the performance of all algorithms 969 in the SMAC environment. The parameters of the CMT algorithm on the SMAC environment are provided in Table 6. Full experimental training curves of five algorithms on rule-based individual 970 reward setting on 11 maps are provided in Figure 6, and the training curves of three algorithms on 971 heuristic individual reward setting on 5 maps are provided in Figure 7.

973 974	Table 4: The common hyper-par	meters of all algorithms on MPE environment.
075	Huper peremeter	Value
976	Number of fully connected layers	
970	Dim of fully-connected layer	2
977	Number of GRU layers	1
976	Dim of RNN hidden layer	64
979	Optimizer	Adam
980	Value loss	huber loss
981	Huber delta	10
982	Batch Size	Jumber of Envs*Number of Agents*Buffer Length
983	Discount factor γ	0.99
984	Activation	ReLU
985	Use reward normalization	True
986	Use feature normalization	True
987	Learning rate of Actor Network	5e-4
988	Learning rate of Critic Network	5e-4
989		
990		
991		
992	Table 5: The hyper-paramete	rs of CMT algorithm on MPE environment.
993	Hyper-parameter	Value
994	Initial Policy Clipping Ratio it	mixed policy 3.0
995	Final Policy Clipping Ratio in	mixed policy 0.5
996	Decaying time range of Policy	Clipping Ratio 2.0 million training steps
997	Policy Clipping Ratio in team	policy 0.2
998	Learning Rate of Lagrangian r	nultiplier 0.01
999		
1000		
1001		
1002		$02 \longrightarrow \lambda = -02$
1003		$\lambda = -0.5$
1004		0
1005	<u>د</u> ۲۵	
1006	san co	for a
1007	P≝ 50	
1008	Ö a	
1009	a 25	
1010	×	
1011	0.0	0.5 1.0
1012		Training Steps 1e6
1013		
1014	Figure 5:	Initial λ value ablation
1015		
1016		
1017		
1018	Table 6: The hyper-parameter	s of CMT algorithm on SMAC environment.
1019	Hyper parameter	Volue
1020	Initial Policy Clipping Patio in	mixed policy 1.0
1021	Final Policy Clipping Ratio in	mixed policy 0.5
1022	Decaying time range of Policy	Clipping Ratio 0.4 million training steps
1023	Policy Clipping Ratio in team	policy 0.05
1024	Learning Rate of Lagrangian r	nultiplier 0.01
4005		





Figure 7: Training curves of three algorithms on heuristic individual reward setting evaluated on 5 maps of SMAC.

D.4 TEST ON FOOTBALL ENVIRONMENTS

1082	Table 7: The hyper-parameters of CMT algorith	nm on GRF environment.
1083		
1084	Hyper-parameter	Value
1085	Initial Policy Clipping Ratio in mixed policy	0.5
1086	Final Policy Clipping Ratio in mixed policy	0.2
1087	Decaying time range of Policy Clipping Ratio	0.5 million training steps
1088	Policy Clipping Ratio in team policy	0.2
1089	Learning Rate of Lagrangian multiplier	0.01

Under GRF environment, we benchmark the proposed approach on the academy counterattack and academy 3_vs_1 with keeper scenarios. The parameters of the CMT algorithm on the GRF environ-ment are listed in Table 6.

Overall, our experimental results on MPE, SMAC, and GRF environments demonstrate that al-though the CMT algorithm may not perform the best in the early stages of policy learning, it keeps improving performance as learning progresses. Ultimately, the CMT algorithm achieves the best performance in most benchmarks. These findings demonstrate that the proposed approach can ef-fectively align the learned policy with the optimal policy.

Ε **BROADER IMPACT**

Our approach, which enables the development of policies that align with the optimal team policy in multi-agent environments with mixed rewards, has far-reaching implications for various realworld applications. For example, in power grid management, unmanned aerial vehicles control, and robotics, where sparse reward functions are prevalent, our approach can be leveraged to efficiently develop policies that ensure consistency with the optimal policy. As such, the proposed approach has the potential to significantly enhance the efficiency and safety performance of MARL algorithms in real-world scenarios.

Furthermore, we are confident that our work has no negative societal implications. Our proposed approach is designed to be benign, with no potential for malicious or unintended uses, and does not raise any concerns related to fairness or privacy.