

# EFFICIENT SEQUENTIAL POLICY OPTIMIZATION VIA OFF-POLICY CORRECTION IN MULTI-AGENT REINFORCEMENT LEARNING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Although trust region policy optimization methods have achieved a lot of success in cooperative multi-agent tasks, most of them face a non-stationarity problem during the learning process. Recently, sequential trust region methods that update policies agent-by-agent have shed light on alleviating the non-stationarity problem. However, these methods are still less sample-efficient when compared to their counterparts (i.e., PPO) in a single-agent setting. To narrow this efficiency gap, we propose the Off-Policy-aware Sequential Policy Optimization (OPSPo) method, which explicitly manages the off-policy-awareness that arises from the sequential policy update process among multiple agents. We prove that our OPSPo has the tightness of the monotonic improvement bound compared with other trust region multi-agent learning methods. Finally, we demonstrate that our OPSPo consistently outperforms strong baselines under challenging multi-agent benchmarks, including StarCraftII micromanagement tasks, Multi-agent MuJoCo, and Google Research Football full game scenarios.

## 1 INTRODUCTION

Trust region learning (Kakade & Langford, 2002), as a class of policy gradient methods (Sutton et al., 1999; Silver et al., 2014), have played an important role in recent advances in single-agent reinforcement learning (Schulman et al., 2015; Haarnoja et al., 2018). Trust Region Policy Optimization (TRPO) (Schulman et al., 2015) and its variant Proximal Policy Optimization (PPO) (Schulman et al., 2017) have been widely used in many fields (Mahmood et al., 2018; Baker et al., 2019; Todorov et al., 2012) and achieved impressive experimental performance (Duan et al., 2016; Kurach et al., 2020). The effectiveness of trust region methods mainly stems from their theoretically guaranteed policy optimization process. Specifically, by restricting the policy optimization to a smaller neighborhood of the current policy, trust region learning obtains a guarantee of monotonic performance improvement at every iteration.

Recently, many works that adopt trust region learning to multi-agent reinforcement learning (MARL) have been proposed, such as methods that use trust region learning to update each agent’s policy independently (De Witt et al., 2020; Yu et al., 2022), as well as methods that coordinate policy updates between agents using trust region learning (Wu et al., 2021). However, these methods update the agents simultaneously, that is, all agents perform policy optimization at the same time and can not observe the change of other agents, which leads to the non-stationarity problem (Hernandez-Leal et al., 2017) during training and hurts performance. To this end, recent works (Kuba et al., 2021; Wang et al., 2023) have proposed sequential trust region learning, which uses trust region learning to sequentially execute agent-by-agent policy optimization. Sequential updates allow the later updated agents to use changes made by preceding agents to optimize their own policies (Bertsekas, 2019), thus stabilizing training.

However, the joint monotonic bounds of these trust region learning methods in MARL, both for simultaneous and sequential updates, are much looser than that of their counterpart (i.e., PPO) in the single agent setting, which leads to sample inefficiency (Li et al., 2022; Wang et al., 2023). In this paper, we take a step toward narrowing this gap. We propose the **Off-Policy-aware Sequential Policy Optimization (OPSPo)** method, which enjoys the tightness of the joint mono-

tonic improvement bound compared with other trust region learning methods in MARL (see Tab. 1 and Theorem 2). Our key idea is to explicitly handle the off-policy nature introduced by the sequential policy update process among multiple agents. Specifically, we start with a vanilla extension of TRPO in MARL, and then significantly improve the joint monotonic improvement bound by performing off-policy corrections on the state distribution and advantage estimation, respectively (see Sec. 3.2). Moreover, we also propose clip range correction which corrects the clip range in clipping-based surrogate objective according to the degree of off-policy (see Sec. 3.3), which further improves the performance of our practical algorithm. We test our OPSPPO on three popular cooperative multi-agent benchmarks: StarCraftII (SMAC) (Samvelyan et al., 2019), multi-agent MuJoCo (MA-MuJoCo) (de Witt et al., 2020), and Google Research Football (GRF) full game scenarios (Kurach et al., 2020). On all benchmark tasks, our OPSPPO consistently outperforms strong baselines with a large margin in both performance and sample efficiency.

In summary, we make three contributions: (i) We propose a novel sequential trust region learning method in MARL, which explicitly handles the off-policy nature caused by the sequential policy update process among multiple agents. We further prove that our method enjoys the tightness of the joint monotonic improvement bound compared with other trust region learning methods in MARL; (ii) We propose a practical clipping-based algorithm with clip-range off-policy correction that can further improve performance; and (iii) Our method significantly outperforms the previous trust region learning methods on three challenging multi-agent benchmarks, including SMAC, MA-MuJoCo, and GRF.

The paper is organized as follows: Sec. 2 provides a background. Sec. 3 introduces the derivation process of our OPSPPO. Sec. 4 presents the experimental studies, and Sec. 5 reviews some related works. Finally, Sec. 6 concludes the paper.

## 2 BACKGROUND

In this section, we first introduce the problem formulation and notations for MARL, and then briefly review trust region learning in MARL.

### 2.1 MARL PROBLEM FORMULATION AND NOTATIONS

We consider formulating a cooperative multi-agent task as a decentralized Markov decision process (DEC-MDP) (Bernstein et al., 2002). An  $n$ -agent DEC-MDP can be defined by a tuple  $(\mathcal{S}, \mathcal{A}, \mathcal{N}, \mathbb{P}, \mathcal{R}, \gamma)$ , where  $\mathcal{N} = \{1, \dots, n\}$  is the set of agents.  $\mathcal{S}$  is the state space.  $\mathcal{A} = \mathcal{A}^1 \times \dots \times \mathcal{A}^n$  is the joint action, where  $\mathcal{A}^i$  is the action space of agent  $i$ . The transition function  $\mathbb{P} : \mathcal{S} \times \mathcal{A}^n \rightarrow \Delta(\mathcal{S})$  maps the state  $s_t$  and the joint action  $\mathbf{a}_t \in \mathcal{A}$  at time step  $t$  to a distribution over the next state  $s_{t+1}$ . All agents receive a collective reward  $r_t = \mathcal{R}(s_t, \mathbf{a}_t)$  according to the reward function  $\mathcal{R} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ .  $\gamma \in [0, 1]$  is a discount factor. At each time step  $t$ , each agent  $i$  takes individual action from its policy  $\pi^i(\cdot|s_t)$  according to the state  $s_t$ , forming the joint action  $\mathbf{a}_t = \{a_t^1, \dots, a_t^n\}$ . All agent’s policy  $\pi^i$  form a joint policy  $\pi(\cdot|s_t) = \pi^1 \times \dots \times \pi^n$ . The joint policy  $\pi$  induces a normalized discounted state visitation distribution  $d^\pi(s) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t Pr(s_t = s | \pi)$ , where  $Pr(\cdot | \pi) : \mathcal{S} \rightarrow \Delta(\mathcal{S})$  is the probability function under a joint policy  $\pi$ . We then define the value function:

$$V^\pi(s) = \mathbb{E}_{\tau \sim (\mathbb{P}, \pi)} \left[ \sum_{t=0}^{\infty} \gamma^t r_t | s_0 = s \right], \quad (1)$$

and the advantage function:

$$A^\pi(s, \mathbf{a}) = r_t + \gamma \mathbb{E}_{s' \sim \mathbb{P}(\cdot | s, \mathbf{a})} [V^\pi(s')] - V^\pi(s), \quad (2)$$

where  $\tau$  denotes one sampled trajectory. The agents’ objective is to find an optimal joint policy  $\pi^*$  that can maximize their expected return, denoted as:

$$\pi^* = \arg \max_{\pi} \mathcal{J}(\pi) = \arg \max_{\pi} \mathbb{E}_{\tau \sim (\mathbb{P}, \pi)} \left[ \sum_{t=0}^{\infty} \gamma^t r_t \right], \quad (3)$$

where  $\mathcal{J}(\pi)$  is the performance of joint policy  $\pi$ . In this paper, we follow the standard centralized training with decentralized execution paradigm (Rashid et al., 2018).

Table 1: Comparing the joint monotonic improvement bounds of trust region MARL algorithms. The proofs of the monotonic bounds can be found in Appendix A.3 and A.4. We sort these algorithms by the tightness of their bounds. MAPPO has the loosest bound and our method has the tightest bound.

Algorithm	Update Scheme	Joint Monotonic Improvement Bound ( $\downarrow$ )
MAPPO	Simultaneous	$4\epsilon \sum_{i=1}^n \alpha^i \frac{\alpha^i}{1-\gamma}$
CoPPO	Simultaneous	$4\epsilon \sum_{i=1}^n \alpha^i \left( \frac{1}{1-\gamma} - \frac{1}{1-\gamma(1-\sum_{j=1}^n \alpha^j)} \right)$
HAPPO	Sequential	$4\epsilon \sum_{i=1}^n \alpha^i \left( \frac{1}{1-\gamma} - \frac{1}{1-\gamma(1-\sum_{j=1}^n \alpha^j)} \right)$
A2PPO	Sequential	$4\epsilon \sum_{i=1}^n \alpha^i \left( \frac{1}{1-\gamma} - \frac{1}{1-\gamma(1-\sum_{j=1}^i \alpha^j)} \right) + \frac{\sum_{i=1}^n \xi^i}{1-\gamma}$
OPSPO (Ours)	Sequential	$4\epsilon \sum_{i=1}^n \alpha^i \left( \frac{1}{1-\gamma} - \frac{1}{1-\gamma(1-\alpha^i)} \right) + \frac{\sum_i^n \alpha^i \delta^i \xi^i}{1-\gamma} + \frac{\sum_{i=1}^n \xi^i}{1-\gamma}$

## 2.2 TRUST REGION POLICY OPTIMIZATION

**Trust Region Methods in RL** As a popular trust region policy optimization method, Trust Region Policy Optimization (TRPO) (Schulman et al., 2015) were proposed in single-agent RL and has the guarantee of monotonic performance improvement of  $\mathcal{J}(\pi)$  at every iteration. If we define the surrogate objective:

$$\mathcal{L}_{\pi_{\text{old}}}(\pi_{\text{new}}) = \mathcal{J}(\pi_{\text{old}}) + \frac{1}{1-\gamma} \mathbb{E}_{(s,a) \sim (d^{\pi_{\text{old}}}, \pi_{\text{new}})} [A^{\pi_{\text{old}}}(s,a)], \quad (4)$$

and let  $\alpha = D_{\text{TV}}^{\max}(\pi_{\text{old}}, \pi_{\text{new}}) = \max_s D_{\text{TV}}(\pi_{\text{old}}(\cdot|s), \pi_{\text{new}}(\cdot|s))$  where  $D_{\text{TV}}$  is the total variation distance, TRPO has the following policy monotonic improvement bound:

$$\begin{aligned} |\mathcal{J}(\pi_{\text{new}}) - \mathcal{L}_{\pi_{\text{old}}}(\pi_{\text{new}})| &\leq 4\alpha \max_{s,a} |A^{\pi_{\text{old}}}(s,a)| \left( \frac{1}{1-\gamma} - \frac{1}{1-\gamma(1-\alpha)} \right) \\ &\leq \frac{4\gamma \max_{s,a} |A^{\pi_{\text{old}}}(s,a)|}{(1-\gamma)^2} \alpha^2. \end{aligned} \quad (5)$$

Eq. 5 states that as the distance between the old policy  $\pi_{\text{old}}$  and a new policy  $\pi_{\text{new}}$  decreases, the surrogate objective  $\mathcal{L}_{\pi_{\text{old}}}(\pi_{\text{new}})$  becomes an increasingly accurate estimate of the actual performance metric  $\mathcal{J}(\pi_{\text{new}})$ . This also implies that a tighter bound improves expected performance by optimizing the surrogate objective more effectively (Li et al., 2022). Proximal Policy Optimization (PPO) (Schulman et al., 2017) uses a clipping-based surrogate objective to approximate TRPO, which is defined as:

$$\mathcal{L}_{\pi_{\text{old}}}^{\text{CLIP}}(\pi_{\text{new}}) = \mathbb{E}_{(s,a) \sim (d^{\pi_{\text{old}}}, \pi_{\text{new}})} \left[ \min \left( \frac{\pi_{\text{new}}(a|s)}{\pi_{\text{old}}(a|s)} A^{\pi_{\text{old}}}(s,a), \text{clip} \left( \frac{\pi_{\text{new}}(a|s)}{\pi_{\text{old}}(a|s)}, 1 \pm \epsilon \right) A^{\pi_{\text{old}}}(s,a) \right) \right]. \quad (6)$$

**Trust Region Methods in MARL** Next, we briefly review recent works that extend the trust region method to the MARL setting. (Yu et al., 2022) proposes Multi-Agent PPO (MAPPO) which is a variant of PPO with centralized critics. (Wu et al., 2021) proposes Coordinate PPO (CoPPO), which obtains a tighter monotonic bound than MAPPO by considering the coordinated adaptation of step size. As we mentioned earlier, a tighter bound means that CoPPO has theoretically better sample efficiency than MAPPO (Li et al., 2022; Wang et al., 2023). Heterogeneous PPO (HAPPO) (Kuba et al., 2021) is the first work to combine the sequential update scheme with trust region methods. Although HAPPO does not achieve a tighter bound than COPPO, it is more stable in training than the simultaneous trust region policy optimization methods, such as MAPPO and CoPPO. Agent-by-agent Policy Optimization (A2PO) (Wang et al., 2023) further improves the sample efficiency by considering the update orders of agents.

However, compared with the monotonic bound of TRPO in the single-agent setting, the monotonic bounds of above multi-agent trust region methods are still loose, which lead to sample inefficiency.

In this paper, we narrow this gap by carefully handling the off-policyyness caused by the sequential policy update process among multiple agents. As shown in Tab. 1, our OPSPPO obtains the tightest joint monotonic bound compared to previous multi-agent trust region methods. Moreover, compared to TRPO, our single agent monotonic bound (see Theorem 1) only has two additional estimation error terms, which can theoretically converge to zero under mild assumptions.

### 3 OFF-POLICYNESS-AWARE SEQUENTIAL POLICY OPTIMIZATION

We first give a native extension of TRPO in sequential updating in MARL in Sec. 3.1. Then, we propose a novel method that greatly improves the monotonic bound by carefully handling the off-policyyness of this TRPO's native extension, in Sec. 3.2. Finally, we give a practical algorithm in Sec. 3.3.

#### 3.1 SEQUENTIAL TRUST REGION POLICY OPTIMIZATION

**Sequential Policy Optimization in MARL** We now have an old joint policy  $\pi$  and some data collected by the old policy  $\pi$ . Our goal is to get a new joint policy  $\bar{\pi}$  updated from the old joint policy  $\pi$  by using the collected data. The general policy optimization process can be defined as:

$$\pi \xrightarrow[\text{Update } \pi]{\max_{\bar{\pi}} \mathcal{G}_{\pi}(\bar{\pi})} \bar{\pi},$$

where  $\mathcal{G}_{\pi}(\bar{\pi})$  is the joint surrogate objective of updating all agents. Without loss of generality, we assume agents are updated in the order  $1, 2, \dots, n$ , and define  $\bar{\pi}^i$  as the updated policy of agent  $i$ . We denote the joint policy after updated agent  $i$  as  $\hat{\pi}^i = \bar{\pi}^1 \times \dots \times \bar{\pi}^i \times \pi^{i+1} \times \dots \times \pi^n$ , and define  $\hat{\pi}^0 = \pi$  and  $\hat{\pi}^n = \bar{\pi}$ . A general sequential policy optimization process can be defined as:

$$\pi = \hat{\pi}^0 \xrightarrow[\text{Update } \pi^1]{\max_{\hat{\pi}^1} \mathcal{L}_{\pi}(\hat{\pi}^1)} \hat{\pi}^1 \rightarrow \dots \rightarrow \hat{\pi}^{n-1} \xrightarrow[\text{Update } \pi^n]{\max_{\hat{\pi}^n} \mathcal{L}_{\hat{\pi}^{n-1}}(\hat{\pi}^n)} \hat{\pi}^n = \bar{\pi},$$

where  $\mathcal{L}_{\hat{\pi}^{i-1}}(\hat{\pi}^i) = \mathcal{J}(\hat{\pi}^{i-1}) + \mathcal{C}(\hat{\pi}^i, \pi)$  is the surrogate objective for agent  $i$ , and  $\mathcal{G}_{\pi}(\bar{\pi}) = \mathcal{J}(\pi) + \sum_{i=1}^n \mathcal{C}(\bar{\pi}^i, \pi)$ . The main difference between sequential policy optimization methods (such as HAPPO and A2PO) lies in the specific design of  $\mathcal{C}(\hat{\pi}^i, \pi)$ .

In this paper, we focus on how to design  $\mathcal{L}_{\hat{\pi}^{i-1}}(\hat{\pi}^i)$  for each agent, so that the monotonic improvement bound  $|\mathcal{J}(\bar{\pi}) - \mathcal{G}_{\pi}(\bar{\pi})| \leq B$  of the joint surrogate objective  $\mathcal{G}_{\pi}(\bar{\pi})$  is more tighter. That is to make the bound  $B$  as small as possible. A tighter bound improves expected performance by optimizing the surrogate objective more effectively (Li et al., 2022).

Moreover, we note that the following inequality holds:

$$|\mathcal{J}(\bar{\pi}) - \mathcal{G}_{\pi}(\bar{\pi})| = |\mathcal{J}(\hat{\pi}^n) - \mathcal{J}(\hat{\pi}^0) - \sum_{i=1}^n \mathcal{C}(\hat{\pi}^i, \pi)| \leq \sum_{i=1}^n |\mathcal{J}(\hat{\pi}^i) - \mathcal{L}_{\hat{\pi}^{i-1}}(\hat{\pi}^i)|, \quad (7)$$

which means that the tighter single-agent monotonic bound can lead to the tighter joint monotonic bound. Based on this observation, we mainly discuss the single-agent monotonic bound in the following sections.

**Vanilla Single-Agent Surrogate Objective** If we natively extend the objective of TRPO (i.e., Eq. 4) to the multi-agent sequential policy optimization, then a vanilla single-agent surrogate objective  $\mathcal{L}_{\hat{\pi}^{i-1}}^{\text{VAN}}(\hat{\pi}^i)$  is obtained, which can be defined as:

$$\mathcal{L}_{\hat{\pi}^{i-1}}^{\text{VAN}}(\hat{\pi}^i) = \mathcal{J}(\hat{\pi}^{i-1}) + \frac{1}{1-\gamma} \mathbb{E}_{(s, \mathbf{a}) \sim (d^{\pi}, \hat{\pi}^i)} [A^{\pi}(s, \mathbf{a})]. \quad (8)$$

Similar to Eq. 5, we can give a single-agent policy monotonic improvement bound of the vanilla surrogate objective  $\mathcal{L}_{\hat{\pi}^{i-1}}^{\text{VAN}}(\hat{\pi}^i)$ , as elaborated in the following proposition.

**Proposition 1.** For agent  $i$ , let  $\epsilon^i = \max_{s, \mathbf{a}} |A^{\hat{\pi}^{i-1}}(s, \mathbf{a})|$ ,  $\Delta^i = \max_{s, \mathbf{a}} |A^{\hat{\pi}^{i-1}}(s, \mathbf{a}) - A^{\pi}(s, \mathbf{a})|$ ,  $\alpha^i = D_{TV}^{\max}(\pi^i, \bar{\pi}^i)$ , where  $D_{TV}(p, q)$  is the total variation distance between distributions  $p$  and  $q$  and we define  $D_{TV}^{\max}(\pi, \bar{\pi}) = \max_s D_{TV}(\pi(\cdot|s), \bar{\pi}(\cdot|s))$ , then we have:

$$|\mathcal{J}(\hat{\pi}^i) - \mathcal{L}_{\hat{\pi}^{i-1}}^{\text{VAN}}(\hat{\pi}^i)| \leq 4\alpha^i \epsilon^i \left( \frac{1}{1-\gamma} - \frac{1}{1-\gamma(1-\alpha^i - \sum_{j=1}^{i-1} \alpha^j)} \right) + \frac{1}{1-\gamma} \Delta^i. \quad (9)$$

For proof see Appendix A.2. Compared with the bound of TRPO (Eq. 5), the bound of the vanilla single-agent objective  $\mathcal{L}_{\hat{\pi}^{i-1}}^{\text{VAN}}(\hat{\pi}^i)$  (Eq. 9) is too looser because there are two extra terms (the blue term and the orange term). The extra terms appear because we do not strictly follow TRPO’s surrogate objective (Eq. 4) to design  $\mathcal{L}_{\hat{\pi}^{i-1}}^{\text{VAN}}(\hat{\pi}^i)$  (Eq. 8). Specifically, we replace the normalized discounted state visitation distribution  $d^{\hat{\pi}^{i-1}}$  induced by  $\hat{\pi}^{i-1}$  with the distribution  $d^\pi$  induced by  $\pi$ , and this substitution leads to the appearance of the blue term  $\sum_{j=1}^{i-1} \alpha^j$ . We also replace the advantage estimation  $A^{\hat{\pi}^{i-1}}(s, \mathbf{a})$  under  $\hat{\pi}^{i-1}$  with the advantage estimation  $A^\pi(s, \mathbf{a})$  under  $\pi$ , and this substitution leads to the appearance of the orange term  $\Delta^i$ . These substitutions are due to the fact that we do not have the data under the joint policy  $\hat{\pi}^{i-1}$ , only the data collected by the old joint policy  $\pi$ . Overall, the off-policy-ness of vanilla objective  $\mathcal{L}_{\hat{\pi}^{i-1}}^{\text{VAN}}(\hat{\pi}^i)$  leads to its looser monotonic bound.

### 3.2 IMPROVING JOINT MONOTONIC BOUND BY OFF-POLICY CORRECTION

In this section, we introduce **Off-Policy-aware Sequential Policy Optimization (OPSPPO)**, a novel algorithm that greatly improves the single-agent monotonic bound by carefully handling the off-policy-ness of TRPO’s native extension (i.e., Eq. 8). This tighter single-agent monotonic bound naturally leads to a tighter joint monotonic bound, which is what we are looking for. To the best of our knowledge, our OPSPPO obtains by far the tightest joint bound compared to previous methods.

Our core idea is to perform off-policy corrections on the state distribution and advantage estimation, respectively. A similar off-policy correction idea is also used to adjust the clip range of each agent to stabilize training (Sec. 3.3).

**Our Surrogate Objective** Our single-agent surrogate objective  $\mathcal{L}_{\hat{\pi}^{i-1}}^{\text{Our}}(\hat{\pi}^i)$  is defined as:

$$\mathcal{L}_{\hat{\pi}^{i-1}}^{\text{Our}}(\hat{\pi}^i) = \mathcal{J}(\hat{\pi}^{i-1}) + \frac{1}{1-\gamma} \mathbb{E}_{(s, \mathbf{a}) \sim (d^{\pi, \pi^{i-1}}, \hat{\pi}^i)} [A^{\pi, \pi^{i-1}}(s, \mathbf{a})], \quad (10)$$

where  $d^{\pi, \pi^{i-1}}$  is an approximation of  $d^{\pi^{i-1}}$  using data collected by  $\pi$ , and  $A^{\pi, \pi^{i-1}}(s, \mathbf{a})$  is also an approximation of  $A^{\pi^{i-1}}(s, \mathbf{a})$  using data collected by  $\pi$ . Based on  $\mathcal{L}_{\hat{\pi}^{i-1}}^{\text{Our}}(\hat{\pi}^i)$ , our joint surrogate objective  $\mathcal{G}_\pi^{\text{Our}}(\bar{\pi})$  is defined as:

$$\mathcal{G}_\pi^{\text{Our}}(\bar{\pi}) = \mathcal{J}(\pi) + \frac{1}{1-\gamma} \sum_{i=1}^n \mathbb{E}_{(s, \mathbf{a}) \sim (d^{\pi, \pi^{i-1}}, \hat{\pi}^i)} [A^{\pi, \pi^{i-1}}(s, \mathbf{a})]. \quad (11)$$

We first analysis the single-agent/joint monotonic bound of our  $\mathcal{L}_{\hat{\pi}^{i-1}}^{\text{Our}}(\hat{\pi}^i)$ , and then discuss the details of these approximations.

**Theorem 1 (Corrected Single-Agent Monotonic Bound).** *For agent  $i$ , let  $\epsilon^i = \max_{s, \mathbf{a}} |A^{\hat{\pi}^{i-1}}(s, \mathbf{a})|$ ,  $\xi^i = \max_{s, \mathbf{a}} |A^{\pi, \hat{\pi}^{i-1}}(s, \mathbf{a}) - A^{\hat{\pi}^{i-1}}(s, \mathbf{a})|$ ,  $\delta^i = \sum_s |d^{\pi, \hat{\pi}^{i-1}}(s) - d^{\hat{\pi}^{i-1}}(s)|$ ,  $\alpha^i = D_{TV}^{\max}(\pi^i, \bar{\pi}^i)$ , where  $D_{TV}(p, q)$  is the total variation distance between distributions  $p$  and  $q$  and we define  $D_{TV}^{\max}(\pi, \bar{\pi}) = \max_s D_{TV}(\pi(\cdot|s), \bar{\pi}(\cdot|s))$ , then we have:*

$$\begin{aligned} |\mathcal{J}(\hat{\pi}^i) - \mathcal{L}_{\hat{\pi}^{i-1}}^{\text{Our}}(\hat{\pi}^i)| &\leq 4\alpha^i \epsilon^i \left( \frac{1}{1-\gamma} - \frac{1}{1-\gamma(1-\alpha^i)} \right) + \frac{1}{1-\gamma} \alpha^i \delta^i \xi^i + \frac{1}{1-\gamma} \xi^i \\ &\leq \frac{4\gamma \epsilon^i}{(1-\gamma)^2} (\alpha^i)^2 + \frac{1}{1-\gamma} \alpha^i \delta^i \xi^i + \frac{1}{1-\gamma} \xi^i. \end{aligned} \quad (12)$$

For proof see Appendix A.3. Compared to the bound (Eq. 9) of the vanilla objective  $\mathcal{L}_{\hat{\pi}^{i-1}}^{\text{VAN}}(\hat{\pi}^i)$ , although the bound (Eq. 12) of our method also has two extra terms, it has better theoretical properties. First, the first term in our bound is exactly the same as the first term of the original TRPO’s bound (Eq. 5), which suggests that as long as the last two terms are small enough, our  $\mathcal{L}_{\hat{\pi}^{i-1}}^{\text{Our}}(\hat{\pi}^i)$  will get a very tight bound. Second, in our last two terms,  $\delta^i$  is the error of approximating  $d^\pi$  with  $d^{\pi, \hat{\pi}^{i-1}}$ , and  $\xi^i$  is the error of approximating  $A^{\hat{\pi}^{i-1}}$  with  $A^{\pi, \hat{\pi}^{i-1}}$ . As the accuracy of these approximations improve,  $\delta^i$  and  $\xi^i$  become smaller and even converge to zero. In contrast, the extra terms in Eq. 12 do not go to zero. Third, our blue term is doubly robust, which is considered to be a good theoretical

property (Dudík et al., 2011; Tang et al., 2020; Jiang & Li, 2016). We can see that if  $A^{\pi, \hat{\pi}^{i-1}}$  is exact (i.e.,  $A^{\pi, \hat{\pi}^{i-1}} = A^{\hat{\pi}^{i-1}}$ ), we have  $\xi^i = 0$ ; if  $d^{\pi, \hat{\pi}^{i-1}}$  is exact (i.e.,  $d^{\pi, \hat{\pi}^{i-1}} = d^{\hat{\pi}^{i-1}}$ ), we have  $\delta^i = 0$ . Therefore, our blue term goes to zero, if either  $A^{\pi, \hat{\pi}^{i-1}}$  or  $d^{\pi, \hat{\pi}^{i-1}}$  are exact. The blue term is thus doubly robust in this sense. Fourth, as we show later, using advanced approximation methods, both  $\delta^i$  and  $\xi^i$  can theoretically go to zero.

Given the such tight single-agent bound, we can prove that our joint objective has the tightest monotonic improvement bound compared to previous methods, as elaborated in the following theorem. We present the joint monotonic bounds of other algorithms in Tab. 1.

**Theorem 2 (Corrected Joint Monotonic Bound).** *For each agent  $i \in \mathcal{N}$ , let  $\epsilon^i = \max_{s, \mathbf{a}} |A^{\hat{\pi}^{i-1}}(s, \mathbf{a})|$ ,  $\epsilon = \max_i \epsilon^i$ ,  $\xi^i = \max_{s, \mathbf{a}} |A^{\pi, \hat{\pi}^{i-1}}(s, \mathbf{a}) - A^{\hat{\pi}^{i-1}}(s, \mathbf{a})|$ ,  $\delta^i = \sum_s |d^{\pi, \hat{\pi}^{i-1}}(s) - d^{\hat{\pi}^{i-1}}(s)|$ ,  $\alpha^i = D_{TV}^{\max}(\pi^j, \bar{\pi}^j)$ , then we have:*

$$\begin{aligned} |\mathcal{J}(\bar{\pi}) - \mathcal{G}_{\pi}^{\text{Our}}(\bar{\pi})| &\leq 4\epsilon \sum_{i=1}^n \alpha^i \left( \frac{1}{1-\gamma} - \frac{1}{1-\gamma(1-\alpha^i)} \right) + \frac{\sum_i \alpha^i \delta^i \xi^i}{1-\gamma} + \frac{\sum_{i=1}^n \xi^i}{1-\gamma} \\ &\leq \frac{4\gamma\epsilon}{(1-\gamma)^2} \sum_i (\alpha^i)^2 + \frac{\sum_i \alpha^i \delta^i \xi^i}{1-\gamma} + \frac{\sum_{i=1}^n \xi^i}{1-\gamma}. \end{aligned} \quad (13)$$

For proof see Appendix A.3. As shown in Eq. 13 and Tab. 1, since our first term ( $4\epsilon \sum_{i=1}^n \alpha^i (\frac{1}{1-\gamma} - \frac{1}{1-\gamma(1-\alpha^i)})$ ) is smallest compared to first terms of other methods, our method achieves the tightest joint monotonic bound if  $\delta^i$  and  $\xi^i$ ,  $\forall i \in \mathcal{N}$  are small enough. The assumption about  $\delta^i$  and  $\xi^i$  is valid, because both  $\delta^i$  and  $\xi^i$  can theoretically go to zero, when advanced approximation methods are used.

**Advantage Estimation Correction** Recall that we can use temporal difference error to approximate an advantage function (Sutton & Barto, 2018). That is  $A^{\hat{\pi}^{i-1}}(s_t, \mathbf{a}_t) \approx r + \gamma V^{\hat{\pi}^{i-1}}(s_{t+1}) - V^{\hat{\pi}^{i-1}}(s_t)$ . Thus, to approximate  $A^{\hat{\pi}^{i-1}}(s_t, \mathbf{a}_t)$  using data collected by  $\pi$ , we only need to use V-trace operator (Espenholt et al., 2018) to approximate  $V^{\hat{\pi}^{i-1}}(s)$ . Given a trajectory  $(s_k, \mathbf{a}_k, r_k)_{k=t}^T$  collected by  $\pi$ , the V-trace target for our value approximation  $V(s_{t+1})$  can be defined as:

$$V_{\text{Target}}^{\pi, \hat{\pi}^{i-1}}(s_t) = V(s_t) + \sum_{k=t}^{T-1} \gamma^{k-t} \left( \prod_{j=t}^k c_j \right) (r_t + \gamma V(s_{t+1}) - V(s_t)), \quad (14)$$

where  $c_j = \lambda \min(1.0, \frac{\hat{\pi}^{i-1}(\mathbf{a}_j | s_j)}{\pi(\mathbf{a}_j | s_j)})$  and  $\lambda$  is a hyper-parameter. Based on Eq. 14, we can get  $A^{\pi, \hat{\pi}^{i-1}}(s_t, \mathbf{a}_t)$  by:

$$A^{\pi, \hat{\pi}^{i-1}}(s_t, \mathbf{a}_t) = r_t + \gamma V(s_{t+1}) - V(s_t). \quad (15)$$

Since V-trace operator (Espenholt et al., 2018) has been proven to be a  $\gamma$ -contraction mapping, the error between  $V^{\hat{\pi}^{i-1}}(s)$  and  $V(s)$  theoretically converges to zero when using the target  $V_{\text{Target}}^{\pi, \hat{\pi}^{i-1}}(s_t)$  to update value function  $V$ . Obviously,  $\xi^i = \max_{s, \mathbf{a}} |A^{\pi, \hat{\pi}^{i-1}}(s, \mathbf{a}) - A^{\hat{\pi}^{i-1}}(s, \mathbf{a})|$  also theoretically goes to zero.

**State Distribution Correction** To approximate  $d^{\hat{\pi}^{i-1}}(s)$  using data collected by  $\pi$ , we use BCH (Liu et al., 2018) to estimate stationary state density ratio  $\omega^{i-1}(s) = d^{\pi^{i-1}}(s)/d^{\pi}(s)$ . It has been proved that  $\omega^{i-1}(s)$  can be approximated by finding a function  $\omega$  over data which minimizes  $\max_f L(\omega, f)$ , defined as:

$$L(\omega, f) = \gamma \mathbb{E}_{(s, \mathbf{a}, s') \sim d^{\pi}} [\omega_{\Delta}(s, \mathbf{a}, s') f(s')] + (1-\gamma) \mathbb{E}_{s \sim d_0^{\pi}} [(1-\omega(s)) f(s)], \quad (16)$$

where  $\omega_{\Delta}(s, \mathbf{a}, s') = \left( \omega(s') - \omega(s) \frac{\pi^{i-1}(\mathbf{a} | s)}{\pi(\mathbf{a} | s)} \right)$ ,  $d_0^{\pi}$  is an initial state distribution under  $\pi$ . Then  $\omega(s) = \omega^{i-1}(s) = d^{\pi^{i-1}}(s)/d^{\pi}(s)$  if and only if  $L(\omega, f) = 0$  for any measurable test function  $f$ . The above equation can be solved by the following min-max problem (Liu et al., 2018):

$$\min_{\omega} \max_{f \in \mathcal{F}} \gamma \mathbb{E}_{(s, \mathbf{a}, s') \sim d^{\pi}} [\omega_{\Delta}(s, \mathbf{a}, s') f(s')] + (1-\gamma) \mathbb{E}_{s \sim d_0^{\pi}} [(1-\omega(s)) f(s)]^2, \quad (17)$$

where  $\mathcal{F}$  is a test function space. Based on the above discussion, if we define  $d^{\hat{\pi}, \hat{\pi}^{i-1}}(s)$  by:

$$d^{\hat{\pi}, \hat{\pi}^{i-1}}(s) = \omega(s) d^{\pi}(s), \quad (18)$$

then it has been proven that  $\max_s |d^{\pi^{i-1}}(s) - d^{\hat{\pi}, \hat{\pi}^{i-1}}(s)| \leq \max_{f \in \mathcal{F}} L(\omega, f)$  if the test function space  $\mathcal{F}$  is sufficiently rich (Liu et al., 2018). In other words,  $\delta^i$  theoretically will become small enough when using Eq. 18 to learn  $\omega(s)$ . Details of learning  $\omega(s)$  can be found in Appendix B.2 and Alg. 2.

### 3.3 THE FINAL ALGORITHM

In this section, we give a practical implementation for optimizing our joint surrogate objective  $\mathcal{G}_{\pi}^{\text{Our}}(\hat{\pi})$ . We first give the native implementation  $\tilde{\mathcal{L}}_{\hat{\pi}^{i-1}}^{\text{Native}}(\hat{\pi}^i)$ , and then give our improved implementation  $\tilde{\mathcal{L}}_{\hat{\pi}^{i-1}}^{\text{Our}}(\hat{\pi}^i)$ . We summarize our proposed Off-Policy-aware Sequential Policy Optimization (OPSPPO) in Alg. 1.

If we directly follow the implementation of original PPO (Eq. 4) for our surrogate objective, we can get a native implementation  $\tilde{\mathcal{L}}_{\hat{\pi}^{i-1}}^{\text{Native}}(\hat{\pi}^i)$ :

$$\mathbb{E}_{(s, \mathbf{a}) \sim (d^{\pi}, \hat{\pi}^i)} \left[ \min \left( r^i \mathbf{r}^{i-1} \omega^{i-1}(s) A^{\pi, \hat{\pi}^{i-1}}, \text{clip} \left( r^i \mathbf{r}^{i-1}, l, h \right) \omega^{i-1}(s) A^{\pi, \hat{\pi}^{i-1}} \right) \right], \quad (19)$$

where  $r^i = \hat{\pi}^i(a^i | s) / \pi^i(a^i | s)$ ,  $\mathbf{r}^{i-1} = \hat{\pi}^{i-1}(\mathbf{a} | s) / \pi(\mathbf{a} | s)$ ,  $l = 1 - \epsilon$ ,  $h = 1 + \epsilon$ , and  $A^{\pi, \hat{\pi}^{i-1}}$  is short for  $A^{\pi, \hat{\pi}^{i-1}}(s, \mathbf{a})$ .

**Clip Range Correction** The main issue with the native implementation Eq. 19 is that due to the off-policyness, the joint policy ratio  $r^i \mathbf{r}^{i-1}$  is likely to be less than  $l = 1 - \epsilon$  or greater than  $h = 1 + \epsilon$ , which results in some data being unable to provide gradients due to the clipped operation. To more fully utilize the data, we scale the base clip range  $(l, h)$  by a correction factor  $r^{i-1}$  that represents the degree of off-policy between  $\hat{\pi}^{i-1}$  and  $\pi$ . The corrected clip range is  $(lr^{i-1}, hr^{i-1})$ . Note that when  $r^{i-1} = 1$ , then  $\hat{\pi}^{i-1} = \pi$ , the corrected clip range reduces to the base clip range. It is important to highlight that, although we scale the clip range, training with the corrected range does not cause instability, as shown in the following theorem.

**Theorem 3 (Stability of Corrected Clip Range).** *Let  $\Pi_{opt}^i$  as the optimal policy set maximizing  $\tilde{\mathcal{L}}_{\hat{\pi}^{i-1}}^{\text{Native}}(\hat{\pi}^i)$  with corrected clip range  $(lr^{i-1}, hr^{i-1})$ ,  $\hat{\pi}_{*,off}^i \in \Pi_{opt}^i$  denotes the optimal joint policy, which achieves the minimum KL divergence over all optimal joint policies, i.e.,  $D_{KL}(\hat{\pi}^{i-1}(\cdot | s_t), \hat{\pi}_{*,off}^i(\cdot | s_t)) \leq D_{KL}(\hat{\pi}^{i-1}(\cdot | s_t), \hat{\pi}_{opt}^i(\cdot | s_t))$  for  $\hat{\pi}_{opt}^i \in \Pi_{opt}^i$  at any timestep  $t$ , and let  $\hat{\pi}_{*,on}^i$  have the similar definition for PPO with data collected by  $\hat{\pi}^{i-1}$  and clip range  $(l, h)$ , we have  $\max_t D_{KL}(\hat{\pi}^{i-1}(\cdot | s_t), \hat{\pi}_{*,off}^i(\cdot | s_t)) = \max_t D_{KL}(\hat{\pi}^{i-1}(\cdot | s_t), \hat{\pi}_{*,on}^i(\cdot | s_t))$  for all timestep  $t$ .*

For proof see Appendix A.5. Theorem 3 tells us that the degree of the policy update distance in  $\tilde{\mathcal{L}}_{\hat{\pi}^{i-1}}^{\text{Native}}(\hat{\pi}^i)$  with the corrected clip range  $(lr^{i-1}, hr^{i-1})$  is the same as that in PPO with the base clip range  $(l, h)$ . In summary, although we scale the clip range, new joint policy  $\hat{\pi}^i$  will not be far from the old one  $\hat{\pi}^{i-1}$ , so the training is stable.

Although we scale the base clip range  $(l, h)$  by a correction factor, the base clip range should not be the same for each agent, because each agent faces a different degree of off-policy. Intuitively, when the off-policy degree is large, we should use a smaller step size to update the agent, that is, reduce the clip range. Therefore, we dynamically adjust the base clip range of each agent by letting  $l^i = \min [\max(\mathbf{r}^{i-1}, 1/\mathbf{r}^{i-1}) \cdot (1 - \epsilon_1), 1 - \epsilon_2]$  and  $h^i = \max [\min(\mathbf{r}^{i-1}, 1/\mathbf{r}^{i-1}) \cdot (1 + \epsilon_1), 1 + \epsilon_2]$ , where  $\epsilon_2 < \epsilon_1$ . After adjustment, the maximum base clip range is  $(1 \pm \epsilon_1)$ , and the minimum base clip range is  $(1 \pm \epsilon_2)$ .

To avoid gradient expansion caused by off-policyness when  $A^{\pi, \hat{\pi}^{i-1}}(s, \mathbf{a}) < 0$ , we use the clipped joint policy ratio  $\tilde{r}^{i-1} = \text{clip}(\mathbf{r}^{i-1}, 1 - \epsilon_1, 1 + \epsilon_1)$  (Wu et al., 2021). To stabilize training, we also use clip operation on  $\omega^{i-1}(s)$  (Amortila et al., 2024), and use  $\tilde{\omega}^{i-1}(s) = \min[\omega^{i-1}(s), 1.0]$ .

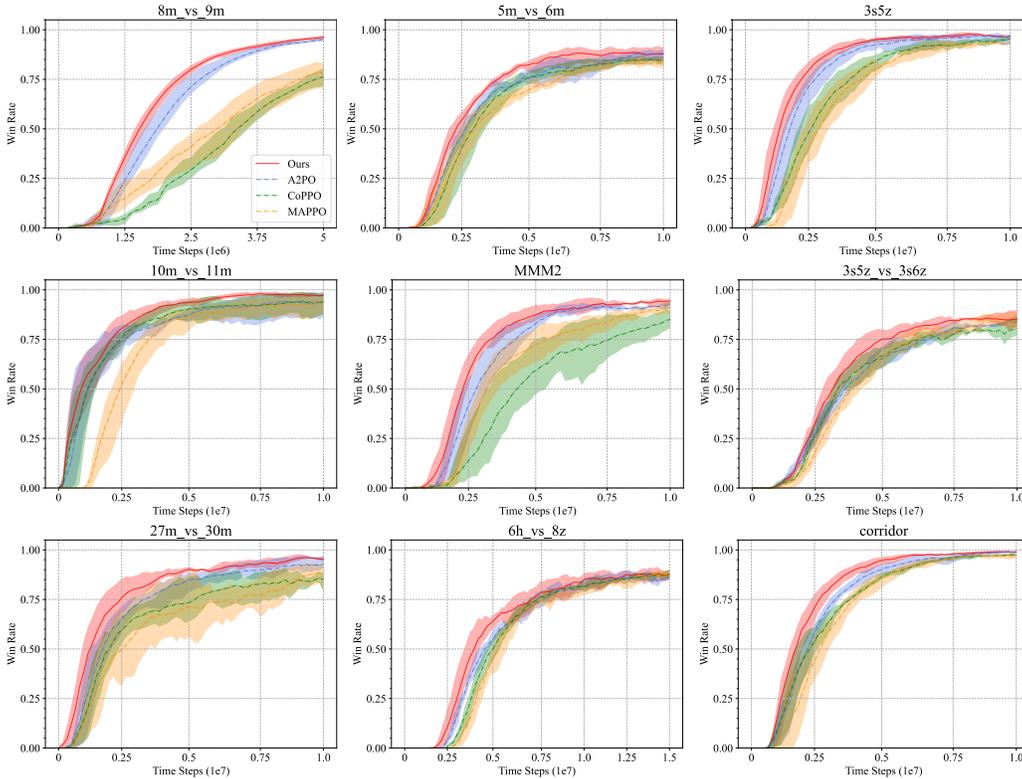


Figure 1: Comparison of our method against baselines over 9 maps of SMAC with various difficulties.

Finally, our practical clipping-based objective  $\tilde{\mathcal{L}}_{\hat{\pi}^{i-1}}^{\text{Our}}(\hat{\pi}^i)$  for updating agent  $i$  is defined as:

$$\mathbb{E}_{(s, \alpha) \sim (d^{\pi}, \hat{\pi}^i)} \left[ \min \left( r^i \tilde{\omega}^{i-1}(s) A^{\pi, \hat{\pi}^{i-1}}, \text{clip} \left( r^i \tilde{\omega}^{i-1}, l^i \tilde{\omega}^{i-1}, h^i \tilde{\omega}^{i-1} \right) \tilde{\omega}^{i-1}(s) A^{\pi, \hat{\pi}^{i-1}} \right) \right]. \quad (20)$$

Compared to the native implementation (Eq. 19), our Eq. 20 introduce a corrected clip range, including a correction factor and a dynamic base clip range, as well as some tricks for stabilizing training.

## 4 EXPERIMENTS

In this section, we empirically evaluate and analyze our OPSPO in the widely-adopted cooperative multi agent benchmarks, including the StarCraftII Multi-agent Challenge (SMAC) (Samvelyan et al., 2019), Multi-agent MuJoCo (MA-MuJoCo) (de Witt et al., 2020), and Google Research Football (GRF) full game scenarios (Kurach et al., 2020).

We compare A2PO with advanced MARL trust-region methods. We first consider MAPPO (Yu et al., 2022) and CoPPO (Wu et al., 2021), which are popular simultaneous trust region learning methods. Then, we consider HAPPO (Kuba et al., 2021) and A2PO (Wang et al., 2023), which are advanced sequential trust region learning methods. Full experimental details can be found in Appendix B.

### 4.1 RESULTS ON CHALLENGING MULTI-AGENT BENCHMARKS

We evaluate our OPSPO in 9 maps of SMAC with various difficulties, 6 scenarios in MA-MuJoCo, and the 5-vs-5 full game scenarios in GRF. As shown in Fig. 1, Fig. 2, and Fig. 3, our method

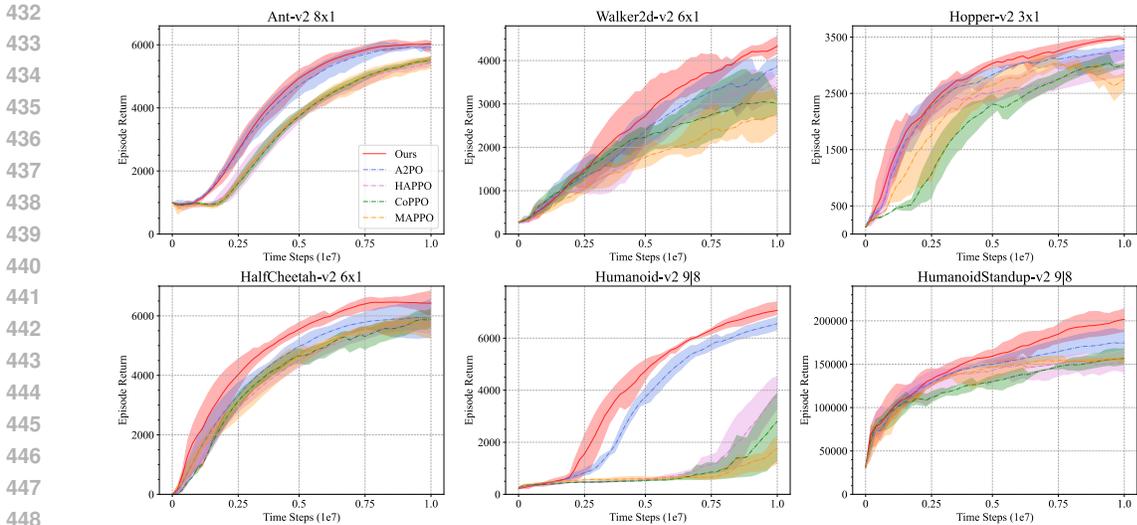


Figure 2: Comparison of our method against baselines over 6 tasks in MA-MuJoCo with different number of robot joints.

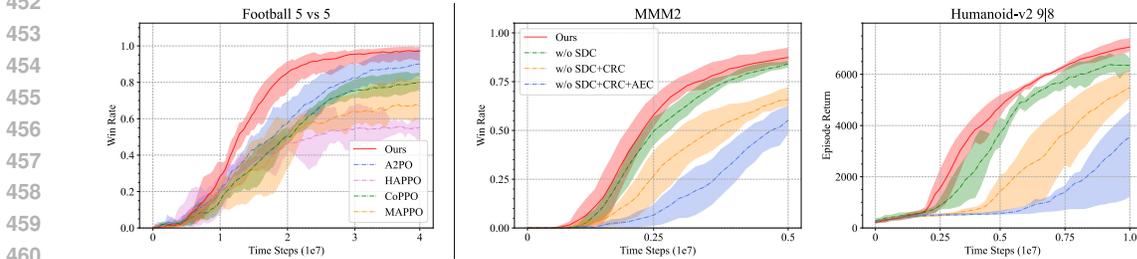


Figure 3: **Left:** Comparison of our method against baselines in 5-vs-5 full game scenarios in Google Research Football. **Right:** Ablation studies for each component in our method on both discrete and continuous action space tasks.

consistently outperforms the strong baselines and achieves better performance and higher sample efficiency in all benchmarks. These experimental results strongly support the theoretical analysis in Sec. 3.

**StarCraftII Multi-agent Challenge (SMAC)** We first evaluate our method on cooperative tasks with discrete action spaces. As shown in Fig. 1, thanks to the better theoretical foundation, i.e. tighter policy monotonic improvement bound, and more flexible policy update step-size adjusted by clip range correction, our method still shows higher sample efficiency even compared with the strong baselines. In addition, our method consistently outperforms other baselines on all tasks, which show the stability of our method.

**Multi-agent MuJoCo environment (MA-MuJoCo)** We then evaluate our method on more complex robotic control multi-agent tasks with continuous action spaces. The experimental results are reported in Fig. 2. As we can see, our method significantly improves the final performance on almost all tasks compared to the baselines. Moreover, we observe that as the complexity of the task increases, our method generally shows increasing advantages over the baselines.

**Google Research Football (GRF)** We also evaluate our method in a full-game GRF scenario with high-dimensional observations, complex action spaces, and long timescales, all of which pose difficulties for agents to discover complex coordination behaviors. As shown in Fig. 3 (left), our method outperforms other methods by a large margin, which once again proves the superiority of our method with a stronger theoretical basis.

## 4.2 ABLATIONS

This section studies how the experimental performance is affected by our proposed components, such as advantage estimation correction (AEC), state distribution correction (SDC), and clip range correction (CRC). We train agents with different components on a super hard task `MMM2` in SMAC, which has a discrete action space, and a complex robotic control task `Humanoid-v2` in MA-MuJoCo, which has a continuous action space. Results are reported in Fig. 3 (right). We make three observations. First, ignoring any component hurts the task performance, which confirms the importance of making off-policy corrections in different aspects. Second, advantage estimation correction brings a very significant performance improvement, which is consistent with our theoretical analysis that it affects both of the last two terms in our monotonic bound (Eq. 13). Third, although we scale the clip range by clip range correction, it does not introduce instability in the training, which is also consistent with the analysis in Sec. 3.3.

## 5 RELATED WORK

Trust Region Policy Optimization (Schulman et al., 2015) and Proximal Policy Optimization (Schulman et al., 2017) are popular trust region methods in the single-agent scenario, which have strong performance mainly due to the guarantee of monotonic policy improvement (Kakade & Langford, 2002). In multi-agent scenarios, De Witt et al. (2020) and Papoudakis et al. (2020) empirically study the performance of directly applying PPO to each agent in multi-agent tasks, and show the inability of native extensions. Their work provides inspiration for subsequent research work. Yu et al. (2022) propose Multi-agent PPO by introducing shared critics and many stable training techniques, and demonstrate strong performance on a large number of multi-agent tasks. Furthermore, Wu et al. (2021) propose Coordinated Proximal Policy Optimization by considering the value decomposition (Sunehag et al., 2017) and coordinated adaptation of step size during the policy update process among agents, and prove the monotonic improvement guarantee. In addition, there are many other works (Wen et al., 2022; Li & He, 2023; Sun et al., 2022) that also discuss trust region methods in MARL scenarios. However, these MARL algorithms suffer from non-stationarity issues due to the simultaneous updating of agents. From the perspective of one agent, the environment dynamics change because other agents also change their policies. As a result, the agent suffers from high variance in gradients and requires more samples to converge (Hernandez-Leal et al., 2017). To this end, sequential updating rather than simultaneous updating has received increasing attention from researchers. Sequential update allows the later updated agents to use changes made by preceding agents to optimize their own policies (Bertsekas, 2019), which makes the environment faced by later agents more stable. Kuba et al. (2021) propose Heterogeneous PPO which combines the sequential update scheme (Bertsekas, 2019) with trust region methods, and demonstrates experimentally and theoretically the advantages offered by sequential updating. Wang et al. (2023) further propose Agent-by-agent Policy Optimization (A2PO) which systematically studies the impact of agent update order on performance and improves the theoretical basis of previous work (Kuba et al., 2021).

The most relevant work to ours is A2PO, which also improves TRPO under sequential policy optimization. A2PO mainly focuses on the impact of agent update order on policy optimization. In contrast, our work is entirely from an off-policy perspective and achieves the tightest bound to date by making off-policy corrections in three aspects: state distribution, advantage estimation, and clip range.

## 6 CONCLUSION

In this paper, we focus on trust region learning in the sequential policy optimization for cooperative multi-agent tasks. We introduce OPSPO, a sequential policy optimization method that explicitly handles the off-policyness caused by the sequential policy update process among agents. We prove that the joint monotonic bound achieved by our OPSPO is the tightest compared to existing trust region MARL methods. Experiments in various benchmarks demonstrate that OPSPO consistently outperforms several strong baselines in performance and sample efficiency in complex tasks. For future work, we plan to continue along the key idea of off-policy correction to improve the broader MARL methods.

## REFERENCES

- 540  
541  
542 Alekh Agarwal, Nan Jiang, Sham M Kakade, and Wen Sun. Reinforcement learning: Theory and  
543 algorithms. *CS Dept., UW Seattle, Seattle, WA, USA, Tech. Rep.*, 32:96, 2019.
- 544 Philip Amortila, Dylan J Foster, Nan Jiang, Ayush Sekhari, and Tengyang Xie. Harnessing density  
545 ratios for online reinforcement learning. *arXiv preprint arXiv:2401.09681*, 2024.
- 546  
547 Bowen Baker, Ingmar Kanitscheider, Todor Markov, Yi Wu, Glenn Powell, Bob McGrew, and Igor  
548 Mordatch. Emergent tool use from multi-agent autocurricula. *arXiv preprint arXiv:1909.07528*,  
549 2019.
- 550 Daniel S Bernstein, Robert Givan, Neil Immerman, and Shlomo Zilberstein. The complexity of  
551 decentralized control of markov decision processes. *Mathematics of operations research*, 27(4):  
552 819–840, 2002.
- 553 Dimitri Bertsekas. Multiagent rollout algorithms and reinforcement learning. *arXiv preprint*  
554 *arXiv:1910.00120*, 2019.
- 555  
556 Xinyun Chen, Lu Wang, Yizhe Hang, Heng Ge, and Hongyuan Zha. Infinite-horizon off-policy  
557 policy evaluation with multiple behavior policies. *arXiv preprint arXiv:1910.04849*, 2019.
- 558  
559 Christian Schroeder De Witt, Tarun Gupta, Denys Makoviichuk, Viktor Makoviychuk, Philip HS  
560 Torr, Mingfei Sun, and Shimon Whiteson. Is independent learning all you need in the starcraft  
561 multi-agent challenge? *arXiv preprint arXiv:2011.09533*, 2020.
- 562  
563 Christian Schroeder de Witt, Bei Peng, Pierre-Alexandre Kamienny, Philip Torr, Wendelin Böhmer,  
564 and Shimon Whiteson. Deep multi-agent reinforcement learning for decentralized continuous  
565 cooperative control. *arXiv preprint arXiv:2003.06709*, 19, 2020.
- 566  
567 Yan Duan, Xi Chen, Rein Houthoofd, John Schulman, and Pieter Abbeel. Benchmarking deep  
568 reinforcement learning for continuous control. In *International conference on machine learning*,  
569 pp. 1329–1338. PMLR, 2016.
- 570  
571 Miroslav Dudík, John Langford, and Lihong Li. Doubly robust policy evaluation and learning.  
572 In *Proceedings of the 28th International Conference on International Conference on Machine*  
573 *Learning*, pp. 1097–1104, 2011.
- 574  
575 Lasse Espeholt, Hubert Soyer, Remi Munos, Karen Simonyan, Vlad Mnih, Tom Ward, Yotam  
576 Doron, Vlad Firoiu, Tim Harley, Iain Dunning, et al. Impala: Scalable distributed deep-rl with im-  
577 portance weighted actor-learner architectures. In *International conference on machine learning*,  
578 pp. 1407–1416. PMLR, 2018.
- 579  
580 Geoff Gordon and Ryan Tibshirani. Karush-kuhn-tucker conditions. *Optimization*, 10(725/36):725,  
581 2012.
- 582  
583 Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy  
584 maximum entropy deep reinforcement learning with a stochastic actor. In *International confer-*  
585 *ence on machine learning*, pp. 1861–1870. PMLR, 2018.
- 586  
587 Pablo Hernandez-Leal, Michael Kaisers, Tim Baarslag, and Enrique Munoz De Cote. A sur-  
588 vey of learning in multiagent environments: Dealing with non-stationarity. *arXiv preprint*  
589 *arXiv:1707.09183*, 2017.
- 590  
591 Nan Jiang and Lihong Li. Doubly robust off-policy value evaluation for reinforcement learning. In  
592 *International conference on machine learning*, pp. 652–661. PMLR, 2016.
- 593  
594 Sham Kakade and John Langford. Approximately optimal approximate reinforcement learning.  
595 In *Proceedings of the Nineteenth International Conference on Machine Learning*, pp. 267–274,  
596 2002.
- 597  
598 Jakub Grudzien Kuba, Ruiqing Chen, Muning Wen, Ying Wen, Fanglei Sun, Jun Wang, and Yaodong  
599 Yang. Trust region policy optimisation in multi-agent reinforcement learning. *arXiv preprint*  
600 *arXiv:2109.11251*, 2021.

- 594 Karol Kurach, Anton Raichuk, Piotr Stanczyk, Michal Zajac, Olivier Bachem, Lasse Espeholt, Car-  
595 los Riquelme, Damien Vincent, Marcin Michalski, Olivier Bousquet, et al. Google research  
596 football: A novel reinforcement learning environment. In *Proceedings of the AAAI Conference*  
597 *on Artificial Intelligence*, volume 34, pp. 4501–4510, 2020.
- 598  
599 Hepeng Li and Haibo He. Multiagent trust region policy optimization. *IEEE Transactions on Neural*  
600 *Networks and Learning Systems*, 2023.
- 601  
602 Hepeng Li, Nicholas Clavette, and Haibo He. An analytical update rule for general policy optimiza-  
603 tion. In *International Conference on Machine Learning*, pp. 12696–12716. PMLR, 2022.
- 604  
605 Qiang Liu, Lihong Li, Ziyang Tang, and Dengyong Zhou. Breaking the curse of horizon: Infinite-  
606 horizon off-policy estimation. *Advances in neural information processing systems*, 31, 2018.
- 607  
608 A Rupam Mahmood, Dmytro Korenkevych, Gautham Vasan, William Ma, and James Bergstra.  
609 Benchmarking reinforcement learning algorithms on real-world robots. In *Conference on robot*  
610 *learning*, pp. 561–591. PMLR, 2018.
- 611  
612 Georgios Papoudakis, Filippos Christianos, Lukas Schäfer, and Stefano V Albrecht. Comparative  
613 evaluation of cooperative multi-agent deep reinforcement learning algorithms. *arXiv preprint*  
614 *arXiv:2006.07869*, 2020.
- 615  
616 Tabish Rashid, Mikayel Samvelyan, Christian Schroeder, Gregory Farquhar, Jakob Foerster, and  
617 Shimon Whiteson. Qmix: Monotonic value function factorisation for deep multi-agent reinforce-  
618 ment learning. In *International Conference on Machine Learning*, pp. 4295–4304. PMLR, 2018.
- 619  
620 Mikayel Samvelyan, Tabish Rashid, Christian Schroeder De Witt, Gregory Farquhar, Nantas  
621 Nardelli, Tim GJ Rudner, Chia-Man Hung, Philip HS Torr, Jakob Foerster, and Shimon Whiteson.  
622 The starcraft multi-agent challenge. *arXiv preprint arXiv:1902.04043*, 2019.
- 623  
624 John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region  
625 policy optimization. In *International conference on machine learning*, pp. 1889–1897. PMLR,  
626 2015.
- 627  
628 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy  
629 optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 630  
631 David Silver, Guy Lever, Nicolas Heess, Thomas Degris, Daan Wierstra, and Martin Riedmiller.  
632 Deterministic policy gradient algorithms. In *International conference on machine learning*, pp.  
633 387–395. Pmlr, 2014.
- 634  
635 Mingfei Sun, Sam Devlin, Jacob Austin Beck, Katja Hofmann, and Shimon Whiteson. Monotonic  
636 improvement guarantees under non-stationarity for decentralized ppo. 2022.
- 637  
638 Peter Sunehag, Guy Lever, Audrunas Gruslys, Wojciech Marian Czarnecki, Vinicius Zambaldi, Max  
639 Jaderberg, Marc Lanctot, Nicolas Sonnerat, Joel Z Leibo, Karl Tuyls, et al. Value-decomposition  
640 networks for cooperative multi-agent learning. *arXiv preprint arXiv:1706.05296*, 2017.
- 641  
642 Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- 643  
644 Richard S Sutton, David McAllester, Satinder Singh, and Yishay Mansour. Policy gradient meth-  
645 ods for reinforcement learning with function approximation. *Advances in neural information*  
646 *processing systems*, 12, 1999.
- 647  
648 Ziyang Tang, Yihao Feng, Lihong Li, Dengyong Zhou, and Qiang Liu. Doubly robust bias reduction  
649 in infinite horizon off-policy estimation. In *International Conference on Learning Representa-*  
650 *tions*, 2020.
- 651  
652 Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control.  
653 In *2012 IEEE/RSJ international conference on intelligent robots and systems*, pp. 5026–5033.  
654 IEEE, 2012.
- 655  
656 Xihuai Wang, Zheng Tian, Ziyu Wan, Ying Wen, Jun Wang, and Weinan Zhang. Order matters:  
657 Agent-by-agent policy optimization. *arXiv preprint arXiv:2302.06205*, 2023.

648 Yuhui Wang, Hao He, Xiaoyang Tan, and Yaozhong Gan. Trust region-guided proximal policy  
649 optimization. *Advances in Neural Information Processing Systems*, 32, 2019.

650 Ying Wen, Hui Chen, Yaodong Yang, Minne Li, Zheng Tian, Xu Chen, and Jun Wang. A game-  
651 theoretic approach to multi-agent trust region optimization. In *International Conference on Dis-*  
652 *tributed Artificial Intelligence*, pp. 74–87. Springer, 2022.

653 Zifan Wu, Chao Yu, Deheng Ye, Junge Zhang, Hankz Hankui Zhuo, et al. Coordinated proximal  
654 policy optimization. *Advances in Neural Information Processing Systems*, 34:26437–26448,  
655 2021.

656 Chao Yu, Akash Velu, Eugene Vinitzky, Jiaxuan Gao, Yu Wang, Alexandre Bayen, and Yi Wu. The  
657 surprising effectiveness of ppo in cooperative multi-agent games. *Advances in Neural Information*  
658 *Processing Systems*, 35:24611–24624, 2022.

## 664 A PROOFS

### 665 A.1 USEFUL LEMMAS

666 **Lemma 1** (Multi-agent Policy Performance Difference Lemma). *Given any joint policies  $\bar{\pi}$  and  $\pi$ ,*  
667 *the difference between the performance of two joint policies can be expressed as:*

$$671 \mathcal{J}(\bar{\pi}) - \mathcal{J}(\pi) = \frac{1}{1-\gamma} \mathbb{E}_{(s,\mathbf{a}) \sim (d^{\bar{\pi}}, \bar{\pi})} [A^{\pi}(s, \mathbf{a})]$$

672 where  $d^{\pi} = (1-\gamma) \sum_{t=0}^{\infty} \gamma^t \Pr(s_t = s | \pi)$  is the normalized discounted state visitation distribution.

673 *Proof.* A corollary of the Policy Performance Difference Lemma, see Lemma 1.16 (Agarwal et al.,  
674 2019).  $\square$

675 **Definition 1.** *A coupling of two probability distributions  $\mu$  and  $\nu$  is a pair of random variables*  
676  *$(X, Y)$  such that the marginal distribution of  $X$  is  $\mu$  and the marginal distribution of  $Y$  is  $\nu$ . A*  
677 *coupling  $(X, Y)$  satisfies the following constraints:  $\Pr(X = x) = \mu(x)$  and  $\Pr(Y = y) = \nu(y)$ .*

678 **Proposition 2.** *For any coupling  $(X, Y)$  that  $D_{TV}(\mu, \nu) \leq \Pr(X \neq Y)$ .*

679 **Proposition 3.** *There exists a coupling  $(X, Y)$  that  $D_{TV}(\mu, \nu) \leq \Pr(X \neq Y)$ .*

680 **Corollary 1.** *For all  $s$ , there exists a coupling  $(\pi(\cdot|s), \bar{\pi}(\cdot|s))$ , that  $\Pr(\mathbf{a} = \bar{\mathbf{a}}) \geq 1 - D_{TV}^{max}(\pi, \bar{\pi})$ ,*  
681 *for  $\mathbf{a} \sim \pi(\cdot|s), \bar{\mathbf{a}} \sim \bar{\pi}(\cdot|s)$*

682 **Corollary 2.** *For all  $s$ ,  $D_{TV}^{max}(\pi(\cdot|s), \bar{\pi}(\cdot|s)) \leq \sum_{i=1}^n D_{TV}(\pi^i(\cdot|s), \bar{\pi}^i(\cdot|s))$ .*

683 **Definition 2.** *If  $(\pi, \bar{\pi})$  is an  $\alpha$ -coupled policy pair, then  $(\mathbf{a}, \bar{\mathbf{a}}|s)$  satisfies  $\Pr(\mathbf{a} \neq \bar{\mathbf{a}}|s) \leq \alpha$  for all*  
684  *$s$ , and  $\mathbf{a} \sim \pi(\cdot|s), \bar{\mathbf{a}} \sim \bar{\pi}(\cdot|s)$ .*

685 From Corollaries 1 and 2, we know that given any joint policy pair  $\pi$  and  $\bar{\pi}$ , select  $\alpha =$   
686  $D_{TV}^{max}(\pi(\cdot|s), \bar{\pi}(\cdot|s))$ , then  $(\pi, \bar{\pi})$  is an  $\alpha$ -coupled policy pair that for all  $s$ ,  $\Pr(\mathbf{a} \neq \bar{\mathbf{a}}|s) \leq$   
687  $D_{TV}^{max}(\pi(\cdot|s), \bar{\pi}(\cdot|s)) \leq \sum_{i=1}^n \alpha^i$ , where  $\alpha^i = D_{TV}^{max}(\pi^i, \bar{\pi}^i)$ .

688 **Lemma 2.** *Given any joint policies  $\pi_1, \pi_2$ , and  $\pi_3$ , if  $\pi_1, \pi_2$  is a coupled policy pair, the following*  
689 *inequality holds:*

$$690 \left| \mathbb{E}_{\mathbf{a}_1 \sim \pi_1} [A^{\pi_3}(s, \mathbf{a}_1)] - \mathbb{E}_{\mathbf{a}_2 \sim \pi_2} [A^{\pi_3}(s, \mathbf{a}_2)] \right| \leq 2\epsilon^{\pi_3} \cdot D_{TV}^{max}(\pi_1, \pi_2) \leq 2\epsilon^{\pi_3} \sum_{i=1}^n \alpha_{\pi_1, \pi_2}^i,$$

691 where  $\alpha_{\pi_1, \pi_2}^i = D_{TV}^{max}(\pi_1^i, \pi_2^i)$  and  $\epsilon^{\pi_3} = \max_{s, \mathbf{a}} |A^{\pi_3}(s, \mathbf{a})|$ .

702 *Proof.*

$$\begin{aligned}
703 & \left| \mathbb{E}_{\mathbf{a}_1 \sim \pi_1} [A^{\pi_3}(s, \mathbf{a}_1)] - \mathbb{E}_{\mathbf{a}_2 \sim \pi_2} [A^{\pi_3}(s, \mathbf{a}_2)] \right| \\
704 & = \left| Pr(\mathbf{a}_1 \neq \mathbf{a}_2 | s) \mathbb{E}_{(\mathbf{a}_1, \mathbf{a}_2) \sim (\pi_1, \pi_2)} [A^{\pi_3}(s, \mathbf{a}_1) - A^{\pi_3}(s, \mathbf{a}_2)] \right| \\
705 & \leq \sum_{i=1}^n \alpha_{\pi_1, \pi_2}^i \mathbb{E}_{(\mathbf{a}_1, \mathbf{a}_2) \sim (\pi_1, \pi_2)} [|A^{\pi_3}(s, \mathbf{a}_1) - A^{\pi_3}(s, \mathbf{a}_2)|] \\
706 & \leq \sum_{i=1}^n \alpha_{\pi_1, \pi_2}^i \cdot 2 \max_{s, \mathbf{a}} |A^{\pi_3}(s, \mathbf{a})|
\end{aligned}$$

712  $\square$

713 **Lemma 3.** *Given any joint policies  $\pi_1, \pi_2$ , if  $\pi_1, \pi_2$  is a coupled policy pair, the following inequality holds:*

$$\left| \mathbb{E}_{\mathbf{a}_1 \sim \pi_1} [A^{\pi_2}(s, \mathbf{a}_1)] \right| \leq 2\epsilon^{\pi_2} \cdot D_{TV}^{max}(\pi_1, \pi_2) \leq 2\epsilon^{\pi_2} \sum_{i=1}^n \alpha_{\pi_1, \pi_2}^i,$$

719 *Proof.* By Lemma 2, the inequality clearly holds.  $\square$

720 **Lemma 4.** *Given any joint policies  $\pi^1, \pi^2$  and  $\pi^3$ , if  $(\pi^1, \pi^2)$  and  $(\pi^2, \pi^3)$  are coupled policy pairs, the following inequality holds:*

$$\begin{aligned}
723 & \left| \mathbb{E}_{(s_t, \mathbf{a}_t) \sim (Pr^{\pi^2}, \pi^2)} [A^{\pi^1}] - \mathbb{E}_{(s_t, \bar{\mathbf{a}}_t) \sim (Pr^{\pi^3}, \pi^2)} [A^{\pi^1}] \right| \\
724 & \leq 4\epsilon^{\pi^1} D_{TV}^{max}(\pi^1, \pi^2) (1 - (1 - D_{TV}^{max}(\pi^2, \pi^3))^t)
\end{aligned}$$

725 where  $\epsilon^{\pi^1} = \max_{s, \mathbf{a}} |A^{\pi^1}(s, \mathbf{a})|$  and we denote  $A(s, \mathbf{a})$  as  $A$  for brevity.

726 *Proof.* Let  $n_t$  represent the times  $\mathbf{a} \neq \bar{\mathbf{a}}$  ( $\pi^2$  disagrees with  $\pi^3$ ) before timestamp  $t$ .

$$\begin{aligned}
728 & \left| \mathbb{E}_{(s_t, \mathbf{a}_t) \sim (Pr^{\pi^2}, \pi^2)} [A^{\pi^1}] - \mathbb{E}_{(s_t, \bar{\mathbf{a}}_t) \sim (Pr^{\pi^3}, \pi^2)} [A^{\pi^1}] \right| \\
729 & = Pr(n_t > 0) \cdot \left| \mathbb{E}_{(s_t, \mathbf{a}_t) \sim (Pr^{\pi^2}, \pi^2) | n_t > 0} [A^{\pi^1}] - \mathbb{E}_{(s_t, \bar{\mathbf{a}}_t) \sim (Pr^{\pi^3}, \pi^2) | n_t > 0} [A^{\pi^1}] \right| \\
730 & \stackrel{(1)}{=} (1 - Pr(n_t = 0)) \cdot E \\
731 & \leq (1 - \prod_{h=1}^t Pr(\mathbf{a}_h = \bar{\mathbf{a}}_h) | \mathbf{a}_h \sim \pi^2(\cdot | s_h), \bar{\mathbf{a}}_h \sim \pi^3(\cdot | s_h)) \cdot E \\
732 & \stackrel{(2)}{\leq} (1 - \prod_{h=1}^t (1 - D_{TV}^{max}(\pi^2, \pi^3))) \cdot E \\
733 & = (1 - (1 - D_{TV}^{max}(\pi^2, \pi^3))^t) \cdot E \\
734 & \stackrel{(3)}{\leq} (1 - (1 - D_{TV}^{max}(\pi^2 | \pi^3))^t) \cdot 4 \cdot D_{TV}^{max}(\pi^1, \pi^2) \cdot \epsilon^{\pi^1} \\
735 & = 4\epsilon^{\pi^1} D_{TV}^{max}(\pi^1, \pi^2) (1 - (1 - D_{TV}^{max}(\pi^2, \pi^3))^t)
\end{aligned}$$

736 In (1), we denote  $E = \left| \mathbb{E}_{(s, \mathbf{a}) \sim (d^{\pi^2}, \pi^2) | n_t > 0} [A^{\pi^1}] - \mathbb{E}_{(s, \bar{\mathbf{a}}) \sim (d^{\pi^3}, \pi^2) | n_t > 0} [A^{\pi^1}] \right|$ . (2) follows the definition of  $\alpha$ -coupled policy pair. (3) follows the Lemma 3.  $\square$

740 **Lemma 5.** *Given any joint policies  $\pi_1, \pi_2, \pi_3$ , and  $\pi_4$ , if  $\pi_1, \pi_2$  is a coupled policy pair, the following inequality holds:*

$$\left| \mathbb{E}_{\mathbf{a}_1 \sim \pi_1} [A^{\pi_3}(s, \mathbf{a}_1)] - \mathbb{E}_{\mathbf{a}_2 \sim \pi_2} [A^{\pi_4}(s, \mathbf{a}_2)] \right| \leq \epsilon^{\pi_3, \pi_4} \cdot D_{TV}^{max}(\pi_1, \pi_2) \leq \epsilon^{\pi_3, \pi_4} \sum_{i=1}^n \alpha_{\pi_1, \pi_2}^i,$$

741 where  $\alpha_{\pi_1, \pi_2}^i = D_{TV}^{max}(\pi_1^i, \pi_2^i)$  and  $\epsilon^{\pi_3, \pi_4} = \max_{s, \mathbf{a}} |A^{\pi_3}(s, \mathbf{a}) - A^{\pi_4}(s, \mathbf{a})|$ .

756 *Proof.*

$$\begin{aligned}
757 & \left| \mathbb{E}_{\mathbf{a}_1 \sim \pi_1} [A^{\pi_3}(s, \mathbf{a}_1)] - \mathbb{E}_{\mathbf{a}_2 \sim \pi_2} [A^{\pi_4}(s, \mathbf{a}_2)] \right| \\
758 & = \left| Pr(\mathbf{a}_1 \neq \mathbf{a}_2 | s) \mathbb{E}_{(\mathbf{a}_1, \mathbf{a}_2) \sim (\pi_1, \pi_2)} [A^{\pi_3}(s, \mathbf{a}_1) - A^{\pi_4}(s, \mathbf{a}_2)] \right| \\
759 & \leq \sum_{i=1}^n \alpha_{\pi_1, \pi_2}^i \mathbb{E}_{(\mathbf{a}_1, \mathbf{a}_2) \sim (\pi_1, \pi_2)} [|A^{\pi_3}(s, \mathbf{a}_1) - A^{\pi_4}(s, \mathbf{a}_2)|] \\
760 & \leq \sum_{i=1}^n \alpha_{\pi_1, \pi_2}^i \cdot \max_{s, \mathbf{a}} |A^{\pi_3}(s, \mathbf{a}) - A^{\pi_4}(s, \mathbf{a})|
\end{aligned}$$

□

761  
762  
763  
764  
765  
766  
767  
**Lemma 6.** *Given any joint policies  $\pi_1, \pi_2, \pi_3, \pi_4, \pi_5$  and  $\pi_6$ , if  $\pi_1, \pi_2$  is a coupled policy pair, the following inequality holds:*

$$770 \left| \mathbb{E}_{(s, \mathbf{a}_1) \sim (d^{\pi_5} - d^{\pi_6}, \pi_1)} [A^{\pi_3}(s, \mathbf{a}_1)] - \mathbb{E}_{(s, \mathbf{a}_2) \sim (d^{\pi_5} - d^{\pi_6}, \pi_2)} [A^{\pi_4}(s, \mathbf{a}_2)] \right| \leq \epsilon^{\pi_3, \pi_4} \cdot \delta^{\pi_5, \pi_6} \sum_{i=1}^n \alpha_{\pi_1, \pi_2}^i$$

771  
772  
773 where  $\delta^{\pi_5, \pi_6} = \sum_s |d^{\pi_5}(s) - d^{\pi_6}(s)|$ .

774  
775 *Proof.*

$$\begin{aligned}
776 & \left| \mathbb{E}_{(s, \mathbf{a}_1) \sim (d^{\pi_5} - d^{\pi_6}, \pi_1)} [A^{\pi_3}(s, \mathbf{a}_1)] - \mathbb{E}_{(s, \mathbf{a}_2) \sim (d^{\pi_5} - d^{\pi_6}, \pi_2)} [A^{\pi_4}(s, \mathbf{a}_2)] \right| \\
777 & = \left| \sum_s [d^{\pi_5}(s) - d^{\pi_6}(s)] \mathbb{E}_{\mathbf{a}_1 \sim \pi_1} [A^{\pi_3}(s, \mathbf{a}_1)] - \sum_s [d^{\pi_5}(s) - d^{\pi_6}(s)] \mathbb{E}_{\mathbf{a}_2 \sim \pi_2} [A^{\pi_4}(s, \mathbf{a}_2)] \right| \\
778 & = \left| \sum_s [d^{\pi_5}(s) - d^{\pi_6}(s)] [\mathbb{E}_{\mathbf{a}_1 \sim \pi_1} [A^{\pi_3}(s, \mathbf{a}_1)] - \mathbb{E}_{\mathbf{a}_2 \sim \pi_2} [A^{\pi_4}(s, \mathbf{a}_2)]] \right| \\
779 & \leq \sum_s |d^{\pi_5}(s) - d^{\pi_6}(s)| \left| \mathbb{E}_{\mathbf{a}_1 \sim \pi_1} [A^{\pi_3}(s, \mathbf{a}_1)] - \mathbb{E}_{\mathbf{a}_2 \sim \pi_2} [A^{\pi_4}(s, \mathbf{a}_2)] \right| \\
780 & \stackrel{(1)}{\leq} \epsilon^{\pi_3, \pi_4} \sum_{i=1}^n \alpha_{\pi_1, \pi_2}^i \sum_s |d^{\pi_5}(s) - d^{\pi_6}(s)| \\
781 & = \epsilon^{\pi_3, \pi_4} \cdot \delta^{\pi_5, \pi_6} \sum_{i=1}^n \alpha_{\pi_1, \pi_2}^i
\end{aligned}$$

782  
783  
784  
785  
786  
787  
788  
789  
790 where (1) follows Lemma 5. □

## 791 A.2 PROOFS OF VANILLA SURROGATE OBJECTIVE

792  
793 Recall that  $\mathcal{L}_{\hat{\pi}^{i-1}}^{\text{VAN}}(\hat{\pi}^i) = \mathcal{J}(\hat{\pi}^{i-1}) + \frac{1}{1-\gamma} \mathbb{E}_{(s, \mathbf{a}) \sim (d^\pi, \hat{\pi}^i)} [A^\pi(s, \mathbf{a})]$ .

794  
795 *Proof.*

$$\begin{aligned}
796 & \left| \mathcal{J}(\hat{\pi}^i) - \mathcal{J}(\hat{\pi}^{i-1}) - \frac{1}{1-\gamma} \mathbb{E}_{(s, \mathbf{a}) \sim (d^\pi, \hat{\pi}^i)} [A^\pi(s, \mathbf{a})] \right| \\
797 & \leq \frac{1}{1-\gamma} \left| \mathbb{E}_{(s, \mathbf{a}) \sim (d^{\hat{\pi}^i}, \hat{\pi}^i)} [A^{\pi^{i-1}}(s, \mathbf{a})] - \mathbb{E}_{(s, \mathbf{a}) \sim (d^\pi, \hat{\pi}^i)} [A^\pi(s, \mathbf{a})] \right| \\
798 & \leq \frac{1}{1-\gamma} \left| \mathbb{E}_{(s, \mathbf{a}) \sim (d^{\hat{\pi}^i}, \hat{\pi}^i)} [A^{\pi^{i-1}}(s, \mathbf{a})] - \mathbb{E}_{(s, \mathbf{a}) \sim (d^\pi, \hat{\pi}^i)} [A^{\pi^{i-1}}(s, \mathbf{a})] \right| \\
799 & + \frac{1}{1-\gamma} \left| \mathbb{E}_{(s, \mathbf{a}) \sim (d^\pi, \hat{\pi}^i)} [A^{\pi^{i-1}}(s, \mathbf{a})] - \mathbb{E}_{(s, \mathbf{a}) \sim (d^\pi, \hat{\pi}^i)} [A^\pi(s, \mathbf{a})] \right| \\
800 & \stackrel{(1)}{\leq} 4\alpha^i \epsilon^i \sum_{t=0}^{\infty} \gamma^t (1 - (1 - \sum_{j=1}^{i-1} \alpha^j)^t) + \frac{1}{1-\gamma} \Delta^i \\
801 & \leq 4\alpha^i \epsilon^i \left( \frac{1}{1-\gamma} - \frac{1}{1-\gamma(1-\alpha^i - \sum_{j=1}^{i-1} \alpha^j)} \right) + \frac{1}{1-\gamma} \Delta^i
\end{aligned}$$

where (1) follows Lemma 4,  $\epsilon^i = \epsilon^{\hat{\pi}^{i-1}} = \max_{s,\mathbf{a}} |A^{\hat{\pi}^{i-1}}|$ ,  $\Delta^i = \max_{s,\mathbf{a}} |A^{\hat{\pi}^{i-1}} - A^\pi|$ ,  $\alpha^i = \alpha_{\hat{\pi}^{i-1}, \hat{\pi}^i}^i = D_{TV}^{\max}(\pi^i, \bar{\pi}^i)$ , where  $D_{TV}(p, q)$  is the total variation distance between distributions  $p$  and  $q$  and we define  $D_{TV}^{\max}(\pi, \bar{\pi}) = \max_s D_{TV}(\pi(\cdot|s), \bar{\pi}(\cdot|s))$ .  $\square$

### A.3 PROOFS OF MONOTONIC POLICY IMPROVEMENT OF OUR OPSPO

**Theorem 1 (Corrected Single-Agent Monotonic Bound).** *For agent  $i$ , let  $\epsilon^i = \max_{s,\mathbf{a}} |A^{\hat{\pi}^{i-1}}(s, \mathbf{a})|$ ,  $\xi^i = \max_{s,\mathbf{a}} |A^{\pi, \hat{\pi}^{i-1}}(s, \mathbf{a}) - A^{\hat{\pi}^{i-1}}(s, \mathbf{a})|$ ,  $\delta^i = \sum_s |d^{\pi, \hat{\pi}^{i-1}}(s) - d^{\hat{\pi}^{i-1}}(s)|$ ,  $\alpha^i = D_{TV}^{\max}(\pi^i, \bar{\pi}^i)$ , where  $D_{TV}(p, q)$  is the total variation distance between distributions  $p$  and  $q$  and we define  $D_{TV}^{\max}(\pi, \bar{\pi}) = \max_s D_{TV}(\pi(\cdot|s), \bar{\pi}(\cdot|s))$ , then we have:*

$$\begin{aligned} |\mathcal{J}(\hat{\pi}^i) - \mathcal{L}_{\hat{\pi}^{i-1}}^{\text{Our}}(\hat{\pi}^i)| &\leq 4\alpha^i \epsilon^i \left( \frac{1}{1-\gamma} - \frac{1}{1-\gamma(1-\alpha^i)} \right) + \frac{1}{1-\gamma} \alpha^i \delta^i \xi^i + \frac{1}{1-\gamma} \xi^i \\ &\leq \frac{4\gamma \epsilon^i}{(1-\gamma)^2} (\alpha^i)^2 + \frac{1}{1-\gamma} \alpha^i \delta^i \xi^i + \frac{1}{1-\gamma} \xi^i. \end{aligned} \quad (21)$$

*Proof.* Recall that  $\mathcal{L}_{\hat{\pi}^{i-1}}^{\text{Our}}(\hat{\pi}^i) = \mathcal{J}(\hat{\pi}^{i-1}) + \frac{1}{1-\gamma} \mathbb{E}_{(s,\mathbf{a}) \sim (d^{\pi, \hat{\pi}^{i-1}}, \hat{\pi}^i)} [A^{\pi, \hat{\pi}^{i-1}}(s, \mathbf{a})]$ .

$$\begin{aligned} &\left| \mathcal{J}(\hat{\pi}^i) - \mathcal{J}(\hat{\pi}^{i-1}) - \frac{1}{1-\gamma} \mathbb{E}_{(s,\mathbf{a}) \sim (d^{\pi, \hat{\pi}^{i-1}}, \hat{\pi}^i)} [A^{\pi, \hat{\pi}^{i-1}}] \right| \\ &\leq \frac{1}{1-\gamma} \left| \mathbb{E}_{(s,\mathbf{a}) \sim (d^{\hat{\pi}^i}, \hat{\pi}^i)} [A^{\hat{\pi}^{i-1}}] - \mathbb{E}_{(s,\mathbf{a}) \sim (d^{\hat{\pi}^{i-1}}, \hat{\pi}^i)} [A^{\hat{\pi}^{i-1}}] \right| \\ &+ \frac{1}{1-\gamma} \left| \mathbb{E}_{(s,\mathbf{a}) \sim (d^{\hat{\pi}^{i-1}}, \hat{\pi}^i)} [A^{\hat{\pi}^{i-1}}] - \mathbb{E}_{(s,\mathbf{a}) \sim (d^{\pi, \hat{\pi}^{i-1}}, \hat{\pi}^i)} [A^{\pi, \hat{\pi}^{i-1}}] \right| \\ &\leq \frac{1}{1-\gamma} \left| \mathbb{E}_{(s,\mathbf{a}) \sim (d^{\hat{\pi}^i}, \hat{\pi}^i)} [A^{\hat{\pi}^{i-1}}] - \mathbb{E}_{(s,\mathbf{a}) \sim (d^{\hat{\pi}^{i-1}}, \hat{\pi}^i)} [A^{\hat{\pi}^{i-1}}] \right| \\ &+ \frac{1}{1-\gamma} \left| \mathbb{E}_{(s,\mathbf{a}) \sim (d^{\hat{\pi}^{i-1}}, \hat{\pi}^i)} [A^{\hat{\pi}^{i-1}}] - \mathbb{E}_{(s,\mathbf{a}) \sim (d^{\hat{\pi}^{i-1}}, \hat{\pi}^i)} [A^{\pi, \hat{\pi}^{i-1}}] - \mathbb{E}_{(s,\mathbf{a}) \sim (d^{\pi, \hat{\pi}^{i-1}}, d^{\hat{\pi}^{i-1}}, \hat{\pi}^i)} [A^{\pi, \hat{\pi}^{i-1}}] \right| \\ &\leq \frac{1}{1-\gamma} \left| \mathbb{E}_{(s,\mathbf{a}) \sim (d^{\hat{\pi}^i}, \hat{\pi}^i)} [A^{\hat{\pi}^{i-1}}] - \mathbb{E}_{(s,\mathbf{a}) \sim (d^{\hat{\pi}^{i-1}}, \hat{\pi}^i)} [A^{\hat{\pi}^{i-1}}] \right| \\ &+ \frac{1}{1-\gamma} \left| \mathbb{E}_{(s,\mathbf{a}) \sim (d^{\hat{\pi}^{i-1}}, \hat{\pi}^i)} [A^{\hat{\pi}^{i-1}}] - \mathbb{E}_{(s,\mathbf{a}) \sim (d^{\hat{\pi}^{i-1}}, \hat{\pi}^i)} [A^{\pi, \hat{\pi}^{i-1}}] \right| \\ &+ \frac{1}{1-\gamma} \left| - \mathbb{E}_{(s,\mathbf{a}) \sim (d^{\pi, \hat{\pi}^{i-1}}, d^{\hat{\pi}^{i-1}}, \hat{\pi}^i)} [A^{\pi, \hat{\pi}^{i-1}}] \right| \\ &\stackrel{(1)}{\leq} \frac{1}{1-\gamma} \left| \mathbb{E}_{(s,\mathbf{a}) \sim (d^{\hat{\pi}^i}, \hat{\pi}^i)} [A^{\hat{\pi}^{i-1}}] - \mathbb{E}_{(s,\mathbf{a}) \sim (d^{\hat{\pi}^{i-1}}, \hat{\pi}^i)} [A^{\hat{\pi}^{i-1}}] \right| \\ &+ \frac{1}{1-\gamma} \left| \mathbb{E}_{(s,\mathbf{a}) \sim (d^{\hat{\pi}^{i-1}}, \hat{\pi}^i)} [A^{\hat{\pi}^{i-1}}] - \mathbb{E}_{(s,\mathbf{a}) \sim (d^{\hat{\pi}^{i-1}}, \hat{\pi}^i)} [A^{\pi, \hat{\pi}^{i-1}}] \right| \\ &+ \frac{1}{1-\gamma} \left| \mathbb{E}_{(s,\mathbf{a}) \sim (d^{\pi, \hat{\pi}^{i-1}}, d^{\hat{\pi}^{i-1}}, \hat{\pi}^i)} [A^{\hat{\pi}^{i-1}}] - \mathbb{E}_{(s,\mathbf{a}) \sim (d^{\pi, \hat{\pi}^{i-1}}, d^{\hat{\pi}^{i-1}}, \hat{\pi}^i)} [A^{\pi, \hat{\pi}^{i-1}}] \right| \\ &\stackrel{(2)}{\leq} 4\epsilon^i \alpha^i \left( \frac{1}{1-\gamma} - \frac{1}{1-\gamma(1-\alpha^i)} \right) + \frac{1}{1-\gamma} \xi^i + \frac{1}{1-\gamma} \alpha^i \delta^i \xi^i \end{aligned}$$

(1) uses  $\mathbb{E}_{(s,\mathbf{a}) \sim (d^{\pi, \hat{\pi}^{i-1}}, d^{\hat{\pi}^{i-1}}, \hat{\pi}^i)} [A^{\hat{\pi}^{i-1}}] = 0$ . (2) uses Lemma 4 and Lemma 6.  $\square$

**Theorem 2 (Corrected Joint Monotonic Bound).** *For each agent  $i \in \mathcal{N}$ , let  $\epsilon^i = \max_{s,\mathbf{a}} |A^{\hat{\pi}^{i-1}}(s, \mathbf{a})|$ ,  $\epsilon = \max_i \epsilon^i$ ,  $\xi^i = \max_{s,\mathbf{a}} |A^{\pi, \hat{\pi}^{i-1}}(s, \mathbf{a}) - A^{\hat{\pi}^{i-1}}(s, \mathbf{a})|$ ,  $\delta^i =$*

$\sum_s |d^{\pi, \hat{\pi}^{i-1}}(s) - d^{\hat{\pi}^{i-1}}(s)|$ ,  $\alpha^i = D_{TV}^{max}(\pi^i, \bar{\pi}^i)$ , then we have:

$$\begin{aligned} |\mathcal{J}(\bar{\pi}) - \mathcal{G}_{\pi}^{Our}(\bar{\pi})| &\leq 4\epsilon \sum_{i=1}^n \alpha^i \left( \frac{1}{1-\gamma} - \frac{1}{1-\gamma(1-\alpha^i)} \right) + \frac{\sum_i^n \alpha^i \delta^i \xi^i}{1-\gamma} + \frac{\sum_{i=1}^n \xi^i}{1-\gamma} \\ &\leq \frac{4\gamma\epsilon}{(1-\gamma)^2} \sum_i^n (\alpha^i)^2 + \frac{\sum_i^n \alpha^i \delta^i \xi^i}{1-\gamma} + \frac{\sum_{i=1}^n \xi^i}{1-\gamma}. \end{aligned} \quad (22)$$

*Proof.* Recall that  $\mathcal{G}_{\pi}^{Our}(\bar{\pi}) = \mathcal{J}(\pi) + \frac{1}{1-\gamma} \sum_{i=1}^n \mathbb{E}_{(s, \mathbf{a}) \sim (d^{\pi, \hat{\pi}^{i-1}}, \hat{\pi}^i)} [A^{\pi, \hat{\pi}^{i-1}}(s, \mathbf{a})]$ .

$$\begin{aligned} &\left| \mathcal{J}(\bar{\pi}) - \mathcal{J}(\pi) - \frac{1}{1-\gamma} \sum_{i=1}^n \mathbb{E}_{(s, \mathbf{a}) \sim (d^{\pi, \hat{\pi}^{i-1}}, \hat{\pi}^i)} [A^{\pi, \hat{\pi}^{i-1}}] \right| \\ &\leq \left| \mathcal{J}(\hat{\pi}^n) - \mathcal{J}(\hat{\pi}^{n-1}) + \dots + \mathcal{J}(\hat{\pi}^1) - \mathcal{J}(\hat{\pi}^0) - \frac{1}{1-\gamma} \sum_{i=1}^n \mathbb{E}_{(s, \mathbf{a}) \sim (d^{\pi, \hat{\pi}^{i-1}}, \hat{\pi}^i)} [A^{\pi, \hat{\pi}^{i-1}}] \right| \\ &\leq \sum_{i=1}^n \left| \mathcal{J}(\hat{\pi}^i) - \mathcal{J}(\hat{\pi}^{i-1}) - \frac{1}{1-\gamma} \mathbb{E}_{(s, \mathbf{a}) \sim (d^{\pi, \hat{\pi}^{i-1}}, \hat{\pi}^i)} [A^{\pi, \hat{\pi}^{i-1}}] \right| \\ &\leq 4\epsilon \sum_{i=1}^n \alpha^i \left( \frac{1}{1-\gamma} - \frac{1}{1-\gamma(1-\alpha^i)} \right) + \frac{\sum_{i=1}^n \xi^i}{1-\gamma} + \frac{\sum_i^n \alpha^i \delta^i \xi^i}{1-\gamma} \end{aligned}$$

□

#### A.4 MONOTONIC POLICY IMPROVEMENT OF MAPPO, CoPPO, HAPPO AND A2PO

We use the formats of the monotonic bounds of MAPPO, CoPPO, HAPPO and A2PO given in (Wang et al., 2023).

#### A.5 PROOFS OF STABILITY OF CORRECTED CLIP RANGE

Recall that  $\tilde{\mathcal{L}}_{\hat{\pi}^{i-1}}^{\text{Native}}(\hat{\pi}^i)$ :

$$\mathbb{E}_{(s, \mathbf{a}) \sim (d^{\pi, \hat{\pi}^i})} \left[ \min \left( r^i \mathbf{r}^{i-1} \omega^{i-1}(s) A^{\pi, \hat{\pi}^{i-1}}, \text{clip} \left( r^i \mathbf{r}^{i-1}, l, h \right) \omega^{i-1}(s) A^{\pi, \hat{\pi}^{i-1}} \right) \right], \quad (23)$$

where  $r^i = \bar{\pi}^i(a^i|s)/\pi^i(a^i|s)$ ,  $\mathbf{r}^{i-1} = \hat{\pi}^{i-1}(\mathbf{a}|s)/\pi(\mathbf{a}|s)$ ,  $l = 1 - \epsilon$ ,  $h = 1 + \epsilon$ , and  $A^{\pi, \hat{\pi}^{i-1}}$  is short for  $A^{\pi, \hat{\pi}^{i-1}}(s, \mathbf{a})$ .

If we use corrected clip range  $(l\mathbf{r}^{i-1}, h\mathbf{r}^{i-1})$ , then we have  $\mathcal{L}_{\hat{\pi}^{i-1}}^{\text{Our}}(\hat{\pi}^i)$ :

$$\mathbb{E}_{(s, \mathbf{a}) \sim (d^{\pi, \hat{\pi}^i})} \left[ \min \left( r^i \mathbf{r}^{i-1} \omega^{i-1}(s) A^{\pi, \hat{\pi}^{i-1}}, \text{clip} \left( r^i \mathbf{r}^{i-1}, l\mathbf{r}^{i-1}, h\mathbf{r}^{i-1} \right) \omega^{i-1}(s) A^{\pi, \hat{\pi}^{i-1}} \right) \right], \quad (24)$$

we denote  $\Pi_{\text{opt}}^i$  as the optimal policy set maximizing  $\mathcal{L}_{\hat{\pi}^{i-1}}^{\text{Our}}(\hat{\pi}^i)$  (Eq. 24).

**Lemma 7.**  $\Pi_{\text{opt}}^i = \{\hat{\pi}^i \mid \text{for all state and action pair } (s, \mathbf{a}) \text{ that } A^{\pi, \hat{\pi}^{i-1}} < 0, \hat{\pi}^i(\mathbf{a}|s) \leq \pi(\mathbf{a}|s)l\mathbf{r}^{i-1}; \text{ for all state and action pair } (s, \mathbf{a}) \text{ that } A^{\pi, \hat{\pi}^{i-1}} > 0, \hat{\pi}^i(\mathbf{a}|s) \geq \min(\pi(\mathbf{a}|s)h\mathbf{r}^{i-1}, 1)\}$ .

*Proof.* Firstly, we prove that a policy  $\hat{\pi}_{\text{opt}}^i$  meeting the conditions in  $\Pi_{\text{opt}}^i$  is the optimal solution maximizing the objective in  $\mathcal{L}_{\hat{\pi}^{i-1}}^{\text{Our}}(\hat{\pi}^i)$ .

Given any  $(s, \mathbf{a})$ , if  $A^{\pi, \hat{\pi}^{i-1}}(s, \mathbf{a}) < 0$ ,  $\mathcal{L}_{\hat{\pi}^{i-1}}^{\text{Our}}(\hat{\pi}^i, s, \mathbf{a})$  could be written as:

$$\mathcal{L}_{\hat{\pi}^{i-1}}^{\text{Our}}(\hat{\pi}^i, s, \mathbf{a}) = \begin{cases} lr^{i-1}\omega^{i-1}(s)A^{\pi, \hat{\pi}^{i-1}}(s, \mathbf{a}), & r^i r^{i-1} \leq lr^{i-1} \\ r^i r^{i-1}\omega^{i-1}(s)A^{\pi, \hat{\pi}^{i-1}}(s, \mathbf{a}), & r^i r^{i-1} > lr^{i-1} \end{cases} \quad (25)$$

$\mathcal{L}_{\hat{\pi}^{i-1}}^{\text{Our}}(\hat{\pi}_{\text{opt}}^i, s, \mathbf{a})$  falls in the first case, because  $\hat{\pi}_{\text{opt}}^i$  meeting the condition in  $\Pi_{\text{opt}}^i$  satisfies  $\frac{\hat{\pi}_{\text{opt}}^i(\mathbf{a}|s)}{\pi(\mathbf{a}|s)} \leq lr^{i-1}$  when  $A^{\pi, \hat{\pi}^{i-1}}(s, \mathbf{a}) < 0$ .

Thus, if  $A^{\pi, \hat{\pi}^{i-1}}(s, \mathbf{a}) < 0$ , then  $\mathcal{L}_{\hat{\pi}^{i-1}}^{\text{Our}}(\hat{\pi}^i, s, \mathbf{a}) \leq \mathcal{L}_{\hat{\pi}^{i-1}}^{\text{Our}}(\hat{\pi}_{\text{opt}}^i, s, \mathbf{a})$  for any  $\hat{\pi}^i$ .

Given any  $(s, \mathbf{a})$ , if  $A^{\pi, \hat{\pi}^{i-1}}(s, \mathbf{a}) > 0$ ,  $\mathcal{L}_{\hat{\pi}^{i-1}}^{\text{Our}}(\hat{\pi}^i, s, \mathbf{a})$  could be written as:

$$\mathcal{L}_{\hat{\pi}^{i-1}}^{\text{Our}}(\hat{\pi}^i, s, \mathbf{a}) = \begin{cases} hr^{i-1}\omega^{i-1}(s)A^{\pi, \hat{\pi}^{i-1}}(s, \mathbf{a}), & r^i r^{i-1} \geq hr^{i-1} \\ r^i r^{i-1}\omega^{i-1}(s)A^{\pi, \hat{\pi}^{i-1}}(s, \mathbf{a}), & r^i r^{i-1} < hr^{i-1} \end{cases} \quad (26)$$

$\mathcal{L}_{\hat{\pi}^{i-1}}^{\text{Our}}(\hat{\pi}_{\text{opt}}^i, s, \mathbf{a})$  also falls in the first case, because  $\hat{\pi}_{\text{opt}}^i$  meeting the condition in  $\Pi_{\text{opt}}^i$  satisfies  $\frac{\hat{\pi}_{\text{opt}}^i(\mathbf{a}|s)}{\pi(\mathbf{a}|s)} \geq hr^{i-1}$  when  $A^{\pi, \hat{\pi}^{i-1}}(s, \mathbf{a}) > 0$ .

Thus, if  $A^{\pi, \hat{\pi}^{i-1}}(s, \mathbf{a}) > 0$ , then  $\mathcal{L}_{\hat{\pi}^{i-1}}^{\text{Our}}(\hat{\pi}^i, s, \mathbf{a}) \leq \mathcal{L}_{\hat{\pi}^{i-1}}^{\text{Our}}(\hat{\pi}_{\text{opt}}^i, s, \mathbf{a})$  for any  $\hat{\pi}^i$ .

Based on such fact, we have proven that a policy  $\hat{\pi}_{\text{opt}}^i$  meeting the conditions in  $\Pi_{\text{opt}}^i$  is the optimal solution.

Secondly, we prove that a policy  $\hat{\pi}_0^i$  not meeting conditions in  $\Pi_{\text{opt}}^i$  is not the optimal solution of maximizing the objective in Eq. 24. In order to prove this, we construct a policy  $\hat{\pi}_{\text{opt}}^i$  satisfying conditions in  $\Pi_{\text{opt}}^i$ . Then,  $\mathcal{L}_{\hat{\pi}^{i-1}}^{\text{Our}}(\hat{\pi}_0^i, s, \mathbf{a}) \leq \mathcal{L}_{\hat{\pi}^{i-1}}^{\text{Our}}(\hat{\pi}_{\text{opt}}^i, s, \mathbf{a})$  for any state and action pair  $(s, \mathbf{a})$ . Based on such fact, we have proven that a policy not meeting the conditions in  $\Pi_{\text{opt}}^i$  is not the optimal solution of maximizing the objective in Eq. 24.

Finally, combining the above results, we prove that  $\Pi_{\text{opt}}^i$  described in Lemma 7 contains all the optimal solutions of maximizing Eq. 24.

□

**Theorem 3 (Stability of Corrected Clip Range).** *Let  $\Pi_{\text{opt}}^i$  as the optimal policy set maximizing  $\mathcal{L}_{\hat{\pi}^{i-1}}^{\text{Our}}(\hat{\pi}^i)$ ,  $\hat{\pi}_{*,\text{off}}^i \in \Pi_{\text{opt}}^i$  denotes the optimal joint policy, which achieves the minimum KL divergence over all optimal joint policies, i.e.,  $D_{\text{KL}}(\hat{\pi}^{i-1}(\cdot|s_t), \hat{\pi}_{*,\text{off}}^i(\cdot|s_t)) \leq D_{\text{KL}}(\hat{\pi}^{i-1}(\cdot|s_t), \hat{\pi}_{\text{opt}}^i(\cdot|s_t))$  for  $\hat{\pi}_{\text{opt}}^i \in \Pi_{\text{opt}}^i$  at any timestep  $t$ , and let  $\hat{\pi}_{*,\text{on}}^i$  have the similar definition for PPO with data collected by  $\hat{\pi}^{i-1}$  and clip range  $(l, h)$ , we have  $\max_t D_{\text{KL}}(\hat{\pi}^{i-1}(\cdot|s_t), \hat{\pi}_{*,\text{off}}^i(\cdot|s_t)) = \max_t D_{\text{KL}}(\hat{\pi}^{i-1}(\cdot|s_t), \hat{\pi}_{*,\text{on}}^i(\cdot|s_t))$  for all timestep  $t$ .*

*Proof.* we denote  $D_{\text{KL}}(\hat{\pi}^{i-1}(\cdot|s_t), \hat{\pi}_{*,\text{off}}^i(\cdot|s_t))$  as  $D_{\text{KL}}^{s_t}(\hat{\pi}^{i-1}, \hat{\pi}_{*,\text{off}}^i)$  and  $D_{\text{KL}}(\hat{\pi}^{i-1}(\cdot|s_t), \hat{\pi}_{*,\text{on}}^i(\cdot|s_t))$  as  $D_{\text{KL}}^{s_t}(\hat{\pi}^{i-1}, \hat{\pi}_{*,\text{on}}^i)$ . In the proof, we need to prove that  $D_{\text{KL}}^{s_t}(\hat{\pi}^{i-1}, \hat{\pi}_{*,\text{off}}^i) = D_{\text{KL}}^{s_t}(\hat{\pi}^{i-1}, \hat{\pi}_{*,\text{on}}^i)$  for any timestep  $t$ . Specifically, we prove this in two cases:  $A^{\pi, \hat{\pi}^{i-1}}(s_t, \mathbf{a}_t) \leq 0$  and  $A^{\pi, \hat{\pi}^{i-1}}(s_t, \mathbf{a}_t) > 0$ .

If  $A^{\pi, \hat{\pi}^{i-1}}(s_t, \mathbf{a}_t) \leq 0$ , the optimal policy  $\hat{\pi}_{*,\text{off}}^i$  can be derived by solving the following constraint optimization problem according to Lemma 7:

$$\begin{aligned} \min_{\hat{\pi}^i} \quad & \sum_{\mathbf{a}} \hat{\pi}^{i-1}(\mathbf{a}|s_t) \log \frac{\hat{\pi}^{i-1}(\mathbf{a}|s_t)}{\hat{\pi}^i(\mathbf{a}|s_t)} \\ \text{s.t.} \quad & \hat{\pi}^i(\mathbf{a}_t|s_t) \leq lr^{i-1}\pi(\mathbf{a}_t|s_t), \\ & \sum_{\mathbf{a}} \hat{\pi}^i(\mathbf{a}|s_t) = 1, \\ & \hat{\pi}^i(\mathbf{a}|s_t) > 0, \end{aligned} \quad (27)$$

where  $\mathbf{a}_t$  denotes the action at timestep  $t$ . By using the Karush-Kuhn-Tucker conditions (Gordon & Tibshirani, 2012), we get:

$$\hat{\pi}_{*,\text{off}}^i(\mathbf{a}|s_t) = \begin{cases} \frac{\hat{\pi}^{i-1}(\mathbf{a}|s_t)(1-\pi(\mathbf{a}_t|s_t)lr^{i-1})}{1-\hat{\pi}^{i-1}(\mathbf{a}|s_t)}, & \mathbf{a} \neq \mathbf{a}_t \\ \pi(\mathbf{a}_t|s_t)lr^{i-1}, & \mathbf{a} = \mathbf{a}_t \end{cases} \quad (28)$$

The corresponding KL divergence is:

$$D_{\text{KL}}^{s_t}(\hat{\pi}^{i-1}, \hat{\pi}_{*,\text{off}}^i) = (1 - \hat{\pi}^{i-1}(\mathbf{a}|s_t)) \log \frac{1 - \hat{\pi}^{i-1}(\mathbf{a}|s_t)}{1 - \hat{\pi}^{i-1}(\mathbf{a}|s_t) \cdot l} - \hat{\pi}^{i-1}(\mathbf{a}|s_t) \log(l) \quad (29)$$

For  $D_{\text{KL}}^{s_t}(\hat{\pi}^{i-1}, \hat{\pi}_{*,\text{on}}^i)$ , we can directly applying Eq. (26) of appendix in (Wang et al., 2019) in our setting. Then we can get  $D_{\text{KL}}^{s_t}(\hat{\pi}^{i-1}, \hat{\pi}_{*,\text{off}}^i)$  equals  $D_{\text{KL}}^{s_t}(\hat{\pi}^{i-1}, \hat{\pi}_{*,\text{on}}^i)$ , when  $A^{\pi, \hat{\pi}^{i-1}}(s_t, \mathbf{a}_t) \leq 0$ .

If  $A^{\pi, \hat{\pi}^{i-1}}(s_t, \mathbf{a}_t) > 0$ , the optimal policy  $\hat{\pi}_{*,\text{off}}^i$  can be derived by solving the following constraint optimization problem according to Lemma 7:

$$\begin{aligned} \min_{\hat{\pi}^i} \quad & \sum_{\mathbf{a}} \hat{\pi}^{i-1}(\mathbf{a}|s_t) \log \frac{\hat{\pi}^{i-1}(\mathbf{a}|s_t)}{\hat{\pi}^i(\mathbf{a}|s_t)} \\ \text{s.t.} \quad & \hat{\pi}^i(\mathbf{a}_t|s_t) \geq \min(hr^{i-1}\pi(\mathbf{a}_t|s_t), 1), \\ & \sum_{\mathbf{a}} \hat{\pi}^i(\mathbf{a}|s_t) = 1, \\ & \hat{\pi}^i(\mathbf{a}|s_t) > 0, \end{aligned} \quad (30)$$

By using the KKT conditions, we get:

$$\hat{\pi}_{*,\text{off}}^i(\mathbf{a}|s_t) = \begin{cases} \frac{\hat{\pi}^{i-1}(\mathbf{a}|s_t)(1-\min(hr^{i-1}\pi(\mathbf{a}_t|s_t), 1))}{1-\hat{\pi}^{i-1}(\mathbf{a}|s_t)}, & \mathbf{a} \neq \mathbf{a}_t \\ \min(hr^{i-1}\pi(\mathbf{a}_t|s_t), 1), & \mathbf{a} = \mathbf{a}_t \end{cases} \quad (31)$$

When  $A^{\pi, \hat{\pi}^{i-1}}(s_t, \mathbf{a}_t) > 0$  and  $hr^{i-1}\pi(\mathbf{a}_t|s_t) \leq 1$ , the KL divergence is:

$$D_{\text{KL}}^{s_t}(\hat{\pi}^{i-1}, \hat{\pi}_{*,\text{off}}^i) = (1 - \hat{\pi}^{i-1}(\mathbf{a}|s_t)) \log \frac{1 - \hat{\pi}^{i-1}(\mathbf{a}|s_t)}{1 - \hat{\pi}^{i-1}(\mathbf{a}|s_t) \cdot h} - \hat{\pi}^{i-1}(\mathbf{a}|s_t) \log(h). \quad (32)$$

For  $D_{\text{KL}}^{s_t}(\hat{\pi}^{i-1}, \hat{\pi}_{*,\text{on}}^i)$ , we can directly applying Eq. (28) of appendix in (Wang et al., 2019) in our setting. Then we can get  $D_{\text{KL}}^{s_t}(\hat{\pi}^{i-1}, \hat{\pi}_{*,\text{off}}^i)$  equals  $D_{\text{KL}}^{s_t}(\hat{\pi}^{i-1}, \hat{\pi}_{*,\text{on}}^i)$ , when  $A^{\pi, \hat{\pi}^{i-1}}(s_t, \mathbf{a}_t) > 0$  and  $hr^{i-1}\pi(\mathbf{a}_t|s_t) \leq 1$ . when  $A^{\pi, \hat{\pi}^{i-1}}(s_t, \mathbf{a}_t) > 0$  and  $hr^{i-1}\pi(\mathbf{a}_t|s_t) > 1$ , we have  $D_{\text{KL}}^{s_t}(\hat{\pi}^{i-1}, \hat{\pi}_{*,\text{off}}^i) = +\infty = D_{\text{KL}}^{s_t}(\hat{\pi}^{i-1}, \hat{\pi}_{*,\text{on}}^i)$ .

Combining above results on two cases ( $A^{\pi, \hat{\pi}^{i-1}}(s_t, \mathbf{a}_t) \leq 0$  and  $A^{\pi, \hat{\pi}^{i-1}}(s_t, \mathbf{a}_t) > 0$ ), we have proven  $D_{\text{KL}}^{s_t}(\hat{\pi}^{i-1}, \hat{\pi}_{*,\text{off}}^i) = D_{\text{KL}}^{s_t}(\hat{\pi}^{i-1}, \hat{\pi}_{*,\text{on}}^i)$  for any timestep  $t$ . Based on such fact, we can conclude that  $\max_t D_{\text{KL}}^{s_t}(\hat{\pi}^{i-1}, \hat{\pi}_{*,\text{off}}^i) = \max_t D_{\text{KL}}^{s_t}(\hat{\pi}^{i-1}, \hat{\pi}_{*,\text{on}}^i)$ .  $\square$

## B EXPERIMENTAL DETAILS

### B.1 PSEUDO CODE

The pseudo code for our OPSPO is given in Alg. 1. The pseudo code for learning state density ratio  $\omega^i(s)$  is given in Alg. 2.

### B.2 IMPLEMENTATION OF STATE DISTRIBUTION CORRECTION

**Algorithm 1:** Off-Policyness-aware Sequential Policy Optimization (OPSPO)

---

```

1026 Algorithm 1: Off-Policyness-aware Sequential Policy Optimization (OPSPO)
1027
1028 1 Initial the joint policy  $\pi_0 = \{\pi_0^1, \dots, \pi_0^n\}$ , and the global value function  $V$ .
1029 2 for iteration  $m = 1, 2, \dots$  do
1030   3 Collect data using  $\pi_{m-1} = \{\pi_{m-1}^1, \dots, \pi_{m-1}^n\}$ .
1031   4 for Order  $i = 1, \dots, n$  do
1032     5 Joint policy  $\hat{\pi}^i = \{\pi_0^1, \dots, \pi_m^i, \pi_{m-1}^{i+1}, \dots, \pi_{m-1}^n\}$ .
1033     6 Compute state density ratio  $\omega^{i-1}$  via Alg. 2.
1034     7 Compute the advantage estimation as  $A^{\pi, \hat{\pi}^{i-1}}(s, \mathbf{a})$  via Eq. 15.
1035     8 Compute the value target  $V_{\text{Target}}^{\pi, \hat{\pi}^{i-1}}(s_t)$  via Eq. 14.
1036     9 Compute the clip range  $(l^i r^{i-1}, h^i r^{i-1})$ .
1037   10 for  $P$  epochs do
1038     11  $\pi_m^i = \arg \max_{\pi_m^i} \tilde{\mathcal{L}}_{\hat{\pi}^{i-1}}^{\text{Our}}(\hat{\pi}^i)$  as in Eq. 20.
1039     12  $V = \arg \min_V \mathbb{E}_{s \sim d^\pi} \|V_{\text{Target}}^{\pi, \hat{\pi}^{i-1}}(s) - V(s)\|^2$ .
1040
1041
1042
1043

```

---

**Algorithm 2:** Optimization of state density ratio  $\omega^i(s)$ 


---

```

1044 Algorithm 2: Optimization of state density ratio  $\omega^i(s)$ 
1045
1046 1 Input: Transition data  $\mathcal{D}$  from  $\omega_n$  behavior joint policies,  $\pi_{k-\omega_n+1}, \dots, \pi_k$ ; a target policy
1047    $\hat{\pi}_k^i$ . Discount factor  $\gamma \in (0, 1)$ , starting state  $\mathcal{D}_0$  from initial distribution,  $T = \omega_e, K = f_e$ ,
1048    $\pi_{\text{mix}} = \frac{1}{\omega_n} \sum_{j=k-\omega_n+1}^k \pi_j(\mathbf{a}|s)$ .
1049 2 Initial the density ratio  $\omega(s) = \omega_{\theta^i}(s)$  to be a neural network parameterized by  $\theta^i$ ,
1050    $f(s) = f_{\psi^i}(s)$  to be a neural network parameterized by  $\psi^i$ .
1051    $\omega_{\Delta}^{\text{mix}}(s, \mathbf{a}, s') = \left( \omega(s') - \omega(s) \frac{\hat{\pi}_k^i(\mathbf{a}|s)}{\pi_{\text{mix}}(\mathbf{a}|s)} \right)$ .
1052 3 for iteration  $1, 2, \dots, T$  do
1053   4 Randomly choose a batch  $\mathcal{M}$  uniformly from the transition data  $\mathcal{D}$  and a batch  $\mathcal{M}_0$ 
1054     uniformly from start states  $\mathcal{D}_0$ .
1055   5 for iteration  $= 1, 2, \dots, K$  do
1056     6 Update the parameter  $\psi^i$  by  $\psi^i \leftarrow \psi^i + \epsilon_{\psi^i} \nabla_{\psi^i} \hat{L}(\omega_{\theta^i}, f_{\psi^i})$ , where
1057       
$$\hat{L}(\omega_{\theta^i}, f_{\psi^i}) = \gamma \frac{1}{|\mathcal{M}|} \sum_{(s, \mathbf{a}, s') \in \mathcal{M}} \omega_{\Delta}^{\text{mix}}(s, \mathbf{a}, s') f(s') - (1-\gamma) \frac{1}{|\mathcal{M}_0|} \sum_{s \in \mathcal{M}_0} (1-\omega(s)) f(s)$$

1058     7 Update the parameter  $\theta^i$  by  $\theta^i \leftarrow \theta^i - \epsilon_{\theta^i} \nabla_{\theta^i} \hat{L}(\omega_{\theta^i}, f_{\psi^i})$ .
1059
1060
1061
1062 8 Output: the density ratio  $\omega^i = \omega_{\theta^i}$ .
1063
1064
1065

```

---

For the training of state density ratio  $\omega^i(s)$ , we adapt the algorithm 2 in (Tang et al., 2020) to perform minimax optimization to train a neural network parameterized  $\omega^i(s; \theta^i)$  and a neural network parameterized test function  $f^i(s; \psi^i)$ . Moreover, to alleviate the partial coverage issue and better predict  $\omega^i(s)$ , we use a multi-behavior policies version (Chen et al., 2019) of BCH. Compared with the original BCH, this variant allows us to use data collected by previous  $\omega_n$  policies, i.e.,  $\pi_{k-\omega_n+1}, \dots, \pi_k$ . The corresponding min-max problem formation is:

$$\min_{\omega} \max_{f \in \mathcal{F}} \gamma \mathbb{E}_{(s, \mathbf{a}, s') \sim d_{\text{mix}}^{\pi}} \left[ \omega_{\Delta}^{\text{mix}}(s, \mathbf{a}, s') f(s') \right]^2 + (1-\gamma) \mathbb{E}_{s \sim d_{\text{mix}, 0}^{\pi}} \left[ (1-\omega(s)) f(s) \right]^2. \quad (33)$$

where  $\pi_{\text{mix}} = \frac{1}{\omega_n} \sum_{j=k-\omega_n+1}^k \pi_j(\mathbf{a}|s)$ ,  $\pi_k$  is the latest behavior policy,  $\omega_n$  is the number of behavior policies,  $d_{\text{mix}}^{\pi}$  is state distribution under  $\pi_{\text{mix}}$ ,  $\omega_{\Delta}^{\text{mix}}(s, \mathbf{a}, s') = \left( \omega(s') - \omega(s) \frac{\pi^i(\mathbf{a}|s)}{\pi_{\text{mix}}(\mathbf{a}|s)} \right)$ ,  $d_0^{\pi}$  is an initial state distribution under  $\pi_{\text{mix}}$ . A detail description can be found in Alg. 2.

Following the suggestion of previous work (Wang et al., 2023), we adopt a parameter sharing setting in SMAC. This makes the sequential updating corrupted, making it very difficult to learn the exact state density ratio  $\omega^i(s)$  by solving Eq. 17. To this end, we use step-wise weighted importance

sampling to approximate  $\omega^{i-1}(s_t)$ . Given  $m$  observed trajectories  $\tau_{j=1}^m$ , for the  $j$ -th trajectory  $\tau_j$  we define  $\omega_j^{i-1}(s_t) = \frac{1}{Z_t} \prod_{k=0}^t \frac{\hat{\pi}^{i-1}(\mathbf{a}_k | s_k)}{\pi(\mathbf{a}_k | s_k)}$ ,  $(\mathbf{a}_k, s_k) \sim \tau_j$ , where  $Z_t = \sum_{j=1}^m \omega_j^{i-1}(s_t)$ .

### B.3 HYPER-PARAMETERS

We list the hyper-parameters used for each task of SMAC in Tab 2. Other parameters use the default settings in A2PO (Wang et al., 2023).

Table 2: Hyper-parameters in SMAC.

Tasks	ppo epoch	$\gamma$	$\epsilon_1$	$\epsilon_2$
8m vs 9m	15	0.95	0.2	0.05
5m vs 6m	10	0.93	0.1	0.05
3s5z	10	0.95	0.2	0.05
10m vs 11m	10	0.95	0.2	0.05
MMM2	10	0.95	0.2	0.05
3s5z vs 3s6z	8	0.90	0.2	0.1
27m vs 30m	8	0.95	0.2	0.05
6h vs 8z	8	0.95	0.2	0.1
corridor	8	0.95	0.2	0.1

For MA-MuJoCo, the output from the last layer is processed by a Tanh layer and the action distribution is modeled as a Gaussian distribution initialized with mean as 0 and log std as -0.5. The probability output of different actions are averaged when computing the policy ratio. We list the hyper-parameters used for each task of MA-MuJoCo in Tab 3. The parameters not mentioned are consistent with A2PO.

Table 3: Hyper-parameters in MA-MuJoCo.

Tasks	ppo epoch	$\gamma$	$\epsilon_1$	$\epsilon_2$	$\omega_e$	$f_e$	$\omega_n$
Ant-v2 8x1	8	0.93	0.2	0.1	10	5	20
Walker2d-v2 6x1	8	0.93	0.2	0.1	10	5	20
Hopper-v2 3x1	8	0.95	0.1	0.05	10	5	20
HalfCheetah-v2 6x1	8	0.93	0.2	0.1	10	5	20
Humanoid-v2 9 8	8	0.90	0.2	0.05	10	5	20
HumanoidStandup-v2 9 8	8	0.93	0.2	0.05	10	5	20

For GRF, We list the hyper-parameters used in the 5-vs-5 scenario in Tab. 4. The parameters not mentioned are consistent with A2PO.

Table 4: Hyper-parameters in GRF.

Hyperparameters	Values
ppo epoch	15
$\gamma$	0.95
$\epsilon_1$	0.2
$\epsilon_2$	0.1
$\omega_e$	5
$\omega_f$	5
$\omega_n$	10