

SMIX(λ): Enhancing Centralized Value Functions for Cooperative Multi-Agent Reinforcement Learning

Xinghu Yao*, Chao Wen*, Yuhui Wang and Xiaoyang Tan
 Nanjing University of Aeronautics and Astronautics, China
 {xinghuyao, chaowen, y.wang, x.tan}@nuaa.edu.cn

Abstract—Learning a stable and generalizable centralized value function (CVF) is a crucial but challenging task in multi-agent reinforcement learning (MARL), as it has to deal with the issue that the joint action space increases exponentially with the number of agents in such scenarios. This paper proposes an approach, named SMIX(λ), that uses an off-policy training to achieve this by avoiding the greedy assumption commonly made in CVF learning. As importance sampling for such off-policy training is both computationally costly and numerically unstable, we proposed to use the λ -return as a proxy to compute the TD error. With this new loss function objective, we adopt a modified QMIX network structure as the base to train our model. By further connecting it with the $Q(\lambda)$ approach from an unified expectation correction viewpoint, we show that the proposed SMIX(λ) is equivalent to $Q(\lambda)$ and hence shares its convergence properties, while without being suffered from the aforementioned curse of dimensionality problem inherent in MARL. Experiments on the StarCraft Multi-Agent Challenge (SMAC) benchmark demonstrate that our approach not only outperforms several state-of-the-art MARL methods by a large margin, but also can be used as a general tool to improve the overall performance of other CTDE-type algorithms by enhancing their CVFs.

Index Terms—Deep reinforcement learning (DRL), multi-agent reinforcement learning (MARL), multi-agent systems, StarCraft Multi-Agent Challenge (SMAC).

I. INTRODUCTION

RECENTLY, reinforcement learning (RL) has made great success in a variety of domains, from game playing [1], [2] to complex continuous control tasks [3]–[5]. However, many real-world problems are inherently multi-agent in nature, such as network packet routing [6], automatic control [7], [8], social dilemmas [9], consensus in multi-agent systems [10]–[12] and multi-player video games [13], which raises great challenges that are never encountered in single-agent settings.

In particular, the main challenges in multi-agent environments include the dimension of joint action space that grows exponentially with the number of agents [14], [15], unstable environments caused by the interaction of individual agents [16], [17], and multi-agent credit assignment in cooperative scenarios with global rewards [14], [15]. These challenges make it troublesome for both fully centralized methods which consider all agents as a single meta agent and fully decentralized methods which individually train each agent by treating other agents as part of the environment.

Recently the paradigm of centralized training with decentralized execution (CTDE) has become popular for multi-agent reinforcement learning [14], [15], [18], [19] due to its conceptual simplicity and practical effectiveness. Its key idea is to learn a centralized value function (CVF) shared by all the agents during training, while each agent acts in a decentralized manner during the execution phase. The CVF works as a proxy to the environment for each agent, through which individual value/advantage functions for each agent can be conveniently learned by incorporating appropriate credit assignment mechanism.

Unfortunately, the central role played by the centralized value function in the CTDE approach seems to receive inadequate attention in current practice - it is commonly treated in the same way as in single-agent settings [14], [15], [17], [20], leading to larger estimation error in multi-agent environments. Furthermore, to reduce the difficulty of decomposing the centralized value function to individual value functions, many algorithms impose extra structural assumptions onto the hypothesis space of the centralized value function during training. For example, *value decomposition networks* (VDN) [20], *monotonic value function factorization* (QMIX) [15], and *factorization with transformation* (QTRAN) [21] assume that the optimal joint action is equivalent to the collection of each agent’s optimal action.

On the other hand, performing an accurate estimation of centralized value function in multi-agent environments is inherently difficult due to the following reasons: 1) the “curse of dimensionality” [22] of the joint action space results in the sparsity of experiences; 2) the partial observability in multi-agent environments become even more severe than in single-agent settings; 3) the dynamics of multi-agent environments are complex and hard to model, partially due to the complicated interactions among agents. In practice, these factors usually contribute to an unreliable and unstable centralized value function with high bias and variance.

To tackle these difficulties, this work proposes a new sample efficient multi-agent reinforcement learning method, named SMIX(λ), under the CTDE framework. We summarize our major contributions in the following threefold.

Firstly, we propose a general optimization framework for CTDE (Centralized Training with Decentralized Execution) in the context of fully cooperative MARL (Multi-Agent Reinforcement Learning) and analyze its theoretical properties.

*Equal contribution.

Code is available at: <https://github.com/chaovven/SMIX>

While many previous CTDE-type MARL methods are built on the centralized greedy behavior (CGB) assumption, which states that a set of decentralized greedy policies (one for each agent) are collectively optimal for the centralized greedy policy. We emphasize more on the importance of improving generalization capability of the learnt policies by incorporating useful inductive bias (such as the non-negative constraint for multi-agent coordination). Relaxing CGB¹ is not only practical, but gives us the flexibility to choose a wider range of methods to train the centralized value function (CVF) as well, considering that many existing methods (especially those Q-learning based, e.g., [15]) have to rely on this assumption for tractable optimization in high-dimensional action space.

Secondly, we propose a novel variant of off-policy SARSA(λ) [23] algorithm for centralized value function estimation in the context of CTDE. The method is characterized by its capability to learn from multi-step lookahead data with improved sample efficiency, but does not suffer from the intensive computational cost originally involved in estimating the product of importance-sampling (IS) ratios in joint action space, hence being particularly suitable for multi-agent reinforcement learning. We further establish its convergence property by building its connection with the well-known $Q(\lambda)$ algorithm [24]. It is shown that our IS-free off-policy SARSA(λ) algorithm is general and is widely applicable to many other CTDE-type multi-agent algorithms such as counterfactual multi-agent (COMA) policy gradients [14], VDN [20] and QTRAN [21], to improve their performance.

Last but not least, with the flexibility provided by the proposed optimization framework, we propose a novel cooperative MARL method named SMIX(λ) to learn a set of decentralized policies, with each agent only knowing its own partial observation, action history and a joint reward shared by all agents. The method is based on the newly-developed IS-free off-policy SARSA(λ) algorithm, incorporated within the enhanced QMIX architecture [15] (hence the name SMIX(λ), meaning ‘‘SARSA(λ)-based MIXture networks’’). It is shown that the proposed SMIX(λ) approach is reliable, easy to implement and significantly outperforms several state-of-the-art CTDE methods including COMA, VDN, QMIX, and QTRAN, on the benchmark of the StarCraft Multi-Agent Challenge [25].

A preliminary version of this work appears in [26]. However, due to page limits, [26] fails to cover all important information about SMIX(λ). This expanded version aims to help readers gain a more comprehensive understanding of SMIX(λ). Specifically, we summarize our major contribution in the introduction section and the related work section (Section II) is added for introducing related works. The Section III is added for building a general CTDE optimization framework. Moreover, the detailed proofs of theorems and more derivation details are included in order to provide a more detailed description of the theoretical properties of QMIX, SMIX(λ) and $Q(\lambda)$ [24]. Besides, Figure 1 illustrates that the estimation method of the CVF in SMIX(λ) can be easily applied to other popular value-based and actor-critic-based CTDE methods. And more

experimental results have been added to show the effectiveness and generality of this estimation method. Last but not least, more implementation details are presented in Algorithm 1 and Section VI, which can improve the reproducibility of SMIX(λ).

In what follows, we first introduce the related work in Section II, then after a brief discussion about CTDE-type methods, a general CTDE optimization framework is established in Section III, the proposed SMIX(λ) method and its theoretical analysis are respectively described in Section IV and Section V. Main experimental results and ablation studies are given in Section VI and we conclude the paper in Section VII.

II. RELATED WORK

Deep reinforcement learning (DRL) has made significant progress in recent years with the powerful representation capabilities of deep neural networks [1], [2], [27]. However, challenges in multi-agent scenarios such as the unstable environments and curse of dimensionality make it hard to apply classic deep reinforcement learning methods to multi-agent environments [14], [28], [29].

The *centralized training with decentralized execution* (CTDE) paradigm provides a simple solution to the above issue by separating the agent learning and execution, under the greedy assumption that the optimal actions for individual agents lead to optimal joint action. It has gradually become the de facto standard in cooperative multi-agent scenarios due to its conceptual simplicity and practical effectiveness. Representative methods include *counterfactual multi-agent* (COMA) policy gradients [14], *value decomposition networks* (VDN) [20], *monotonic value function factorization* (QMIX) [15], and *factorization with transformation* (QTRAN) [21] – COMA is an on-policy actor-critic method that uses a carefully designed counterfactual baseline to perform credit assignment, while VDN, QMIX, and QTRAN are typical value-based CTDE methods by learning individual agents through learning a centralized value function first.

Our SMIX(λ) belongs to the CTDE framework as well, but we focus more on how to improve the sample efficiency and how to perform an accurate estimation of the centralized value function. Our key idea is to use off-policy training to achieve these goals while relaxing the greedy assumption in the learning stage. Although off-policy methods are known to improve the sample efficiency [1], [30], the popular importance sampling methods for off-policy training are problematic as these methods often involve calculating a product of a series of importance sampling ratios, which is not only computationally costly but has high variance [31] as well.

The estimation of the CVF plays a central role in the CTDE framework, as its bias and variance directly affect the performance of the whole system. Foerster et al. adopt a variant of TD(λ) [14] to balance the bias and variance in CVF estimation, but they use an on-policy training method which could be sample inefficient. Precup et al. propose an importance-sampling-based TD(λ) method in the single-agent setting and prove the convergence property with linear function approximation [32].

¹By ‘relaxing CGB’, we mean that in our framework the CGB assumption is not explicitly needed although it could be derived from some constraints such as the non-negative constraint.

Under the fully decentralized framework, [33] proposes two methods to stabilize the off-policy training process. For the first method, the authors proposed adding extra time tags onto every piece of information to be stored in the replay buffer. This allows to decay obsolete data. In the second method, each agent’s value function is conditioned on a fingerprint that disambiguates the age of the data sampled from the replay memory. Both methods help alleviate the non-stationary problem, but we adopt an alternative mechanism for this rather than manipulating the replay memory directly, as discussed in Section III-D.

It is worth mentioning that in practice, however, off-policy correction is not always needed in off-policy learning, especially when the behavior policy and target policy are close to each other. For example, Hernandez et al. find that it is possible to ignore off-policy correction over off-policy SARSA [23] and $Q(\sigma)$ [34] without seeing an adverse effect on the overall performance [35]. Fujimoto et al. show that the off-policy experiences generated during the interaction with the environment tend to be heavily correlated to the current policy, and their experimental results also reveal that the distribution of off-policy data during the training procedure is very close to that of the current policy [36]. Their analysis provides an intuitive explanation for why performance can be improved even without off-policy correction. Unfortunately, a notable gap remains between the empirical success and the underlying theoretical support. In Theorem 3 of Section V, we give a principled way to justify this ‘thumb of rule’, showing that a computationally efficient experience replay method such as ours in the context of MARL is not only feasible but theoretically sound as well.

There are other ways to deal with multi-agent problems. For example, agents can exchange information with each other through a communication channel [37], [38] or a shared network structure [39], [40]. The opponent modeling methods aim to infer other agents’ policies by interacting with the environment and observing other agents’ policy [41], [42]. These methods are designed to establish explicit or implicit connections among agents, and we aim to coordinate each agent through a centralized value function. Thus, these approaches are complementary to ours. Besides, multi-agent self-play [43] and task decomposition [44] have recently been shown to be useful in MARL.

Finally, there have been several attempts on the StarCraft Multi-Agent Challenge (SMAC) [25], including [14], [15], [33]. The results of [15] are the published state-of-the-art in value-based methods and [14] in actor-critic-based methods.

III. THE CTDE OPTIMIZATION FRAMEWORK

A. Problem Formulation

The cooperative multi-agent task we considered can be described as a variant of *decentralized partially observable Markov decision process* (Dec-POMDP) [45]. Specifically, this task can be defined as a tuple: $\mathcal{G} = \langle \mathcal{S}, \mathcal{A}, \mathcal{P}, r, \mathcal{Z}, \mathcal{O}, N, \gamma \rangle$, where $s \in \mathcal{S}$ denotes the true state of the environment, \mathcal{A} is the action set for each of N agents, and $\gamma \in [0, 1]$ is the discount factor. At each timestep, each agent $i \in$

$\{1, 2, \dots, N\}$ chooses an action $a^i \in \mathcal{A}$, forming a joint action $\mathbf{a} = \{a^1, a^2, \dots, a^N\} \in \mathcal{A}^N$. Then the environment gets into next state s' through a dynamic transition function $\mathcal{P}(s'|s, \mathbf{a}) : \mathcal{S} \times \mathcal{A}^N \times \mathcal{S} \mapsto [0, 1]$. All agents share the same reward function $r(s, \mathbf{a}) : \mathcal{S} \times \mathcal{A}^N \mapsto \mathbb{R}$. We consider a partial observable scenario² in which each agent draws partial observation $o \in \mathcal{O}$ from the observation function $\mathcal{Z}(s, i) : \mathcal{S} \times N \mapsto \mathcal{O}$. Each agent i also has an observation-action history $\tau^i \in \mathcal{T} \equiv (\mathcal{O} \times \mathcal{A})^*$, on which it conditions a stochastic policy. A stochastic policy is a mapping defined as $\pi(a|\tau) : \mathcal{T} \times \mathcal{A} \mapsto [0, 1]$.

In Dec-POMDP, the goal of a learning algorithm is to obtain a group of individual policies that can maximize the expected discounted returns $\mathbb{E}_{\mathbf{a} \in \pi, s \in \mathcal{S}} [\sum_{t=0}^{\infty} \gamma^t r(s, \mathbf{a})]$. To simplify notation, we denote joint quantities over agents in bold. We also omit the index i of each agent when there is no ambiguity in the following sections.

The *centralized training with decentralized execution* (CTDE) approach provides a solution to the Dec-POMDP problem by introducing a centralized structure to coordinate the decentralized policies. In the training phase of the CTDE paradigm, a centralized action-value function $Q([s, \boldsymbol{\tau}], \mathbf{a})$ (or simply expressed as $Q(\boldsymbol{\tau}, \mathbf{a})$) is learned from the local observation history of all agents (denoted as $\boldsymbol{\tau} = \{\tau^1, \tau^2, \dots, \tau^N\}$) and the global state (denoted as s), while during the execution phase, each agent’s policy π^i only relies on its own observation-action history τ^i .

B. A General CTDE Optimization Framework

In this paper, we consider a fully cooperative multi-agent scenario, in which the relationship between centralized value function and decentralized value functions satisfies particular coordinative constraints. Formally, we consider value-based CTDE as the following optimization problem.

$$\begin{aligned} & \underset{\boldsymbol{\pi}}{\text{maximize}} && \mathbb{E}_{s_0 \sim \rho_0(s_0), a^1 \sim \pi^1, \dots, a^N \sim \pi^N} [Q_{tot}^{\boldsymbol{\pi}}(s_0, \mathbf{a})] \\ & \text{subject to} && Q_{tot}(\boldsymbol{\tau}, \mathbf{a}) = f(Q^1(\tau^1, a^1), \dots, Q^N(\tau^N, a^N)), \end{aligned} \quad (1)$$

where f is a *state-dependent continuous* function, $\rho_0 : \mathcal{S} \rightarrow [0, 1]$ is the distribution of the initial state s_0 . The value function of each agent are combined to give the centralized value function Q_{tot} , through a *state-dependent continuous* function f implemented as a neural network.

Designing a proper optimization constraint in (1) is crucial because it directly affects the generalization ability and optimization cost of the algorithm. In this paper, we assume the relationship between centralized value function and decentralized functions satisfies the following assumption.

Assumption 1. In fully cooperative multi-agent scenarios, the coordination among agents can be constrained through $\frac{\partial Q_{tot}}{\partial Q^i} \geq 0, i \in \{1, \dots, N\}$.

This non-negative constraint (positive weight) first appears in VDN [20] and is generalized in QMIX [15] as a way to ensure tractable optimization in MARL. In this paper, however,

²In standard Dec-POMDP, the observation function $\mathcal{Z}(\mathbf{o}|\mathbf{a}, s')$ denotes the probability of the observing joint observation \mathbf{o} given that joint action \mathbf{a} was taken and led to state s' (cf., [45]).

we take it as a prior to capture the coordination information required to solve cooperative tasks. This effectively strikes a balance between coordinative expressivity and learning difficulty. As a result, each agent has a chance to influence the team reward positively instead of canceling out with each other in fully cooperative scenarios.

The following theorem guarantees that if assumption 1 is satisfied, then the optimal policy of each agent is conditioned on the optimal joint actions of other agents through the centralized observable critic Q_{tot}^π . This helps to address the non-stationary issue of MARL, as discussed later.

Theorem 1. *In optimization problem (1) with assumption 1 satisfied, if the joint action-value function Q_{tot}^π is already obtained, then we have:*

$$\operatorname{argmax}_{a^i} Q^i(\tau^i, a^i) = \operatorname{argmax}_{a^i} \max_{a^1, \dots, a^{i-1}, a^{i+1}, \dots, a^N} Q_{tot}^\pi(\tau, \mathbf{a}) \quad (2)$$

Proof: We denote $a^{1*} = \operatorname{argmax}_{a^1} Q^1(\tau^1, a^1)$. Since $\frac{\partial Q_{tot}^\pi}{\partial Q^1} \geq 0$, we have

$$\begin{aligned} & \max_{a^1, \dots, a^{i-1}, a^{i+1}, \dots, a^N} Q_{tot}^\pi(\tau, a^1, a^{i-1}, a^i, a^{i+1}, \dots, a^N) \\ &= \max_{a^2, \dots, a^{i-1}, a^{i+1}, \dots, a^N} Q_{tot}^\pi(\tau, a^{1*}, a^{i-1}, a^i, a^{i+1}, \dots, a^N). \end{aligned}$$

Similarly, we denote $a^{j*} = \operatorname{argmax}_{a^j} Q^j(\tau^j, a^j)$ for $j \in \{2, 3, i-1, i+1, N\}$. Then, due to $\frac{\partial Q_{tot}^\pi}{\partial Q^j} \geq 0$, we have

$$\begin{aligned} & \max_{a^2, \dots, a^{i-1}, a^{i+1}, \dots, a^N} Q_{tot}^\pi(\tau, a^{1*}, \dots, a^{i-1}, a^i, a^{i+1}, \dots, a^N) \\ &= \max_{a^3, \dots, a^{i-1}, a^{i+1}, \dots, a^N} Q_{tot}^\pi(\tau, a^{1*}, a^{2*}, a^{i-1}, a^i, a^{i+1}, \dots, a^N) \\ &= \dots \\ &= Q_{tot}^\pi(\tau, a^{1*}, a^{2*}, a^{(i-1)*}, a^i, a^{(i+1)*}, \dots, a^{N*}). \end{aligned}$$

Then, given $\frac{\partial Q_{tot}^\pi}{\partial Q^i} \geq 0$, we have

$$\begin{aligned} & \operatorname{argmax}_{a^i} Q_{tot}^\pi(\tau, a^{1*}, a^{2*}, a^{(i-1)*}, a^i, a^{(i+1)*}, \dots, a^{N*}) \\ &= \operatorname{argmax}_{a^i} Q^i(\tau^i, a^i) = a^{i*}. \end{aligned}$$

In an idealized setting where each agent observes the full state (each agent's partial observation τ^i is equivalent to the global state s), the optimal joint action-value function Q^* is a solution of a *multi-agent Markov decision process* (MMDP) which is itself equivalent to a standard MDP with \mathcal{A}^N as the action space [46]. In such cases, the joint action $\mathbf{a}_* = \operatorname{argmax}_{\mathbf{a}} Q^*(\tau, \mathbf{a})$ is a *Nash equilibrium* from the view of game theory [47]. Specifically, a Nash equilibrium is achieved among a group of agent if the following expression holds for all $i \in \{1, \dots, N\}$:

$$Q^*(\tau, \mathbf{a}_*) = Q^*(\tau, a_*^i, \mathbf{a}_*^{-i}) \geq Q^*(\tau, a^i, \mathbf{a}_*^{-i}), \forall a^i \in \mathcal{A}, \quad (3)$$

where each agent acts with the *best response* a_*^i to others.

We adopt a compact notation for the joint action of all agents except i as $\mathbf{a}_*^{-i} \triangleq [a_*^1, \dots, a_*^{i-1}, a_*^{i+1}, \dots, a_*^N]$. The following theorem shows that the *Nash equilibrium* \mathbf{a}_* can be obtained when each agent i acts greedy according to $Q^i(\tau^i, a^i)$ if the optimal centralized value function $Q^*(\tau, \mathbf{a})$ satisfies assumption 1.

Corollary 1. *If the optimal joint action-value function Q^* satisfies assumption 1, then the optimal joint policy $\mathbf{a}_* = \{a_*^1, a_*^2, \dots, a_*^N\}$, where $a_*^{i*} = \operatorname{argmax}_{a^i} Q^i(\tau^i, a^i)$ and $Q^i(\tau^i, a^i)$ is the decentralized value function obtained by solving (1).*

Proof: According to theorem 1, we have $a_*^{i*} = \operatorname{argmax}_{a^i} Q_{tot}^{\pi_*}(\tau, a^i, a_*^{-i*})$. Thus, we have:

$$Q_{tot}^{\pi_*}(\tau, a_*^{i*}, a_*^{-i*}) \geq Q_{tot}^{\pi_*}(\tau, a^i, a_*^{-i*}), \forall a^i \in \mathcal{A}.$$

The above expression holds for all $i \in \{1, 2, \dots, N\}$. ■

C. CTDE Under Centralized Greedy Behavior Assumption

To facilitate the freedom of each agent to make decision based on its local observation without consulting the centralized value function, the following centralized greedy behavior (CGB) assumption is usually adopted:

$$\operatorname{argmax}_{\mathbf{a}} Q_{tot}(\tau, \mathbf{a}) = \begin{pmatrix} \operatorname{argmax}_{a^1} Q^1(\tau^1, a^1) \\ \vdots \\ \operatorname{argmax}_{a^N} Q^N(\tau^N, a^N) \end{pmatrix}. \quad (4)$$

This assumption establishes a structural constraint between the centralized value function and the decentralized value functions, which can be thought of as a simplified credit assignment mechanism during the execution phase.

Many methods aim to impose a structure constraint between centralized value function and decentralized value functions to ensure the CGB assumption is satisfied. Specifically, VDN [20] uses the following additive combination:

$$Q_{tot}(\tau, \mathbf{a}; \boldsymbol{\theta}) = \sum_{i=1}^N \alpha_i Q^i(\tau^i, a^i; \theta^i), \alpha_i \geq 0, \quad (5)$$

where $\boldsymbol{\theta}$ is the collection of the parameter vector θ^i for each agent's action-value function. In VDN [20], all the combination coefficients $\alpha_i, i = 1, 2, \dots, N$ are set to 1. QMIX [15] extends this additive value factorization to a more general case by directly enforcing $\frac{\partial Q_{tot}}{\partial Q^i} \geq 0, i \in \{1, \dots, N\}$ via a non-negative mixing network f . With this, we can easily obtain the following theorem.

Theorem 2. *For QMIX, we have*

$$\begin{aligned} & \max_{\mathbf{a}} Q_{tot}(\tau, \mathbf{a}) = \\ & f\left(\tau, \operatorname{argmax}_{a^1} Q^1(\tau^1, a^1), \dots, \operatorname{argmax}_{a^N} Q^N(\tau^N, a^N)\right). \end{aligned} \quad (6)$$

Proof:

Due to $Q_{tot}(\tau, \mathbf{a}) = f(Q^1(\tau^1, a^1), \dots, Q^N(\tau^N, a^N))$, where f is a state-dependent continuous monotonic function, and $\mathbf{a} = (a_1, \dots, a_N)$. Thus, we have

$$\begin{aligned} & Q_{tot}(\tau, \operatorname{argmax}_{a^1} Q^1(\tau^1, a^1), \dots, \operatorname{argmax}_{a^N} Q^N(\tau^N, a^N)) \\ &= f\left(\max_{a^1} Q^1(\tau^1, a^1), \dots, \max_{a^N} Q^N(\tau^N, a^N)\right). \end{aligned}$$

Since $\frac{\partial f}{\partial Q^i} \geq 0$, given $(\bar{a}^2, \dots, \bar{a}^n)$, we have

$$f(Q^1(\tau^1, a^1), Q^2(\tau^2, \bar{a}^2), \dots, Q^N(\tau^N, \bar{a}^N)) \\ \leq f\left(\max_{a^1} Q^1(\tau^1, a^1), Q^2(\tau^2, \bar{a}^2), \dots, Q^N(\tau^N, \bar{a}^N)\right)$$

for any a^1 . Similarly, given $(\bar{a}^1, \dots, \bar{a}^{k-1}, \bar{a}^{k+1}, \dots, \bar{a}^N)$, we have

$$f(Q^1(\tau^1, \bar{a}^1), \dots, Q^k(\tau^k, a^k), \dots, Q^N(\tau^N, \bar{a}^N)) \\ \leq f\left(Q^1(\tau^1, \bar{a}^1), \dots, \max_{a^k} Q^k(\tau^k, a^k), \dots, Q^N(\tau^N, \bar{a}^N)\right)$$

for any a^k . Finally, for any (a^1, \dots, a^N) , we have

$$f(Q^1(\tau^1, a^1), \dots, Q^N(\tau^N, a^N)) \\ \leq f\left(\max_{a^1} Q^1(\tau^1, a^1), \dots, Q^N(\tau^N, a^N)\right) \\ \leq f\left(\max_{a^1} Q^1(\tau^1, a^1), \dots, \max_{a^N} Q^N(\tau^N, a^N)\right).$$

Therefore, we obtain

$$\max_{a^1, \dots, a^N} f(Q^1(\tau^1, a^1), \dots, Q^N(\tau^N, a^N)) \\ = f\left(\max_{a^1} Q^1(\tau^1, a^1), \dots, \max_{a^N} Q^N(\tau^N, a^N)\right),$$

which is the specific form of (6). \blacksquare

The above theorem shows that with the help of a non-negative mixing network f , performing Q-learning is tractable for $Q_{tot}(\tau, \mathbf{a})$ even in high dimensional joint action space. So does VDN. However, this non-negative constraint is only a sufficient condition of the CGB assumption, hence restricting the algorithm's representational complexity. To achieve the best generalization, QTRAN [21] allows to search the best hypothesis in a space equivalent to the one specified by the CGB condition. In other words, QTRAN works in a larger hypothesis space than both VDN and QMIX. However, optimizing in a larger hypothesis space requires more optimization efforts. Although a coordinate-decent-type method is proposed in QTRAN to address this issue, the method's scalability and the range of practical use can be limited. See more discussions on this in the experimental section.

D. Dealing with Non-Stationarity

The non-stationary problem is a key issue in the MARL setting [48]. An environment that contains multiple agents, seen from the angle of an individual agent, is constantly changing with the changes of any other agent. Through this lens, the CTDE (Centralized Training with Decentralized Execution) scheme provides a simple solution to this issue by introducing a fully observable critic (i.e., the centralized value function Q_{tot}). The fully observable critic accesses to the observation and actions of all agents, and consequently the environment becomes stationary even though the policy of other agents changes. The fully observable critic also helps to simplify the optimization problem and enables the algorithm to scale better to more complex scenarios.

Alternatively, the CTDE strategy can be thought of as a way to condition each agent's policy on other agents' joint policy

through the centralized observable critic (cf. Theorem 1). But this idea can also be implemented in a more straightforward way under fully decentralized settings without the need of the centralized observable critic, as in *Hyper Q-learning* [49], *Distributed Q-learning* [50] and *multi-agent fingerprints* [33]. Specifically, in *Hyper Q learning*, each agent learns a Q function over all possible opponent strategies. Similarly, the *Distributed Q-learning* [50] considers a decentralized cooperative multi-agent problem in fully observable environments and the joint action is available for all agents in the training time. The *Distributed Q-learning* algorithm updates the Q-values only when there is a guaranteed improvement. In *multi-agent fingerprints*, each agent's value function conditions on a fingerprint that disambiguates the age of the data sampled from the replay memory, which helps to stabilize the training process in fully decentralized settings.

IV. METHODS

In this section, we give the details of the proposed SMIX(λ) method, which is a SARSA(λ) [51] style off-policy method that aims at learning a flexible and generalizable centralized value function within the CTDE framework.

A. Motivation

One of the most popular methods to estimate Q_{tot}^π in (1) is based on the Q-learning, as in VDN, QMIX, and QTRAN. However, to make this computationally tractable in a joint action space, as discussed in III-C, the CGB assumption has to be made. This potentially restricts the range of possible learning algorithms for solving (1).

Alternatively, one can use an Expected-SARSA-based method. The Expected SARSA estimates its TD (temporal difference) target reinforcement signal with an expectation value of the next state-action pairs in an on-policy way [23]. In other words, it does not have to perform a greedy search over the joint action space and hence does not need the CGB assumption either. This allows us to decouple the learning algorithm and the CGB assumption. Through an iteration process over Q_{tot}^π , the optimal joint action-value function can be obtained [52]. Besides, it is well-known that the Q-learning algorithm can be viewed as a particular case of Expected SARSA in which the expectation over actions is replaced with a deterministically greedy one [23], [34].

In what follows, we will propose a new off-policy value function estimation method based on this idea and apply it to centralized value function estimation within the CTDE framework.

B. Importance-Sampling-Free Off-policy SARSA(λ)

Denoting the behavior policy as μ and the target policy as π , a general off-policy strategy to evaluate the Q value function (centralized value function Q_{tot} in multi-agent setting) for π using data τ generated by following μ can be expressed as follows [31],

$$Q(\tau, \mathbf{a}) \leftarrow Q(\tau, \mathbf{a}) + \mathbb{E}_\mu \left[\sum_{t \geq 0} \gamma^t \left(\prod_{i=1}^t \rho_i \right) \delta_t^\pi \right], \quad (7)$$

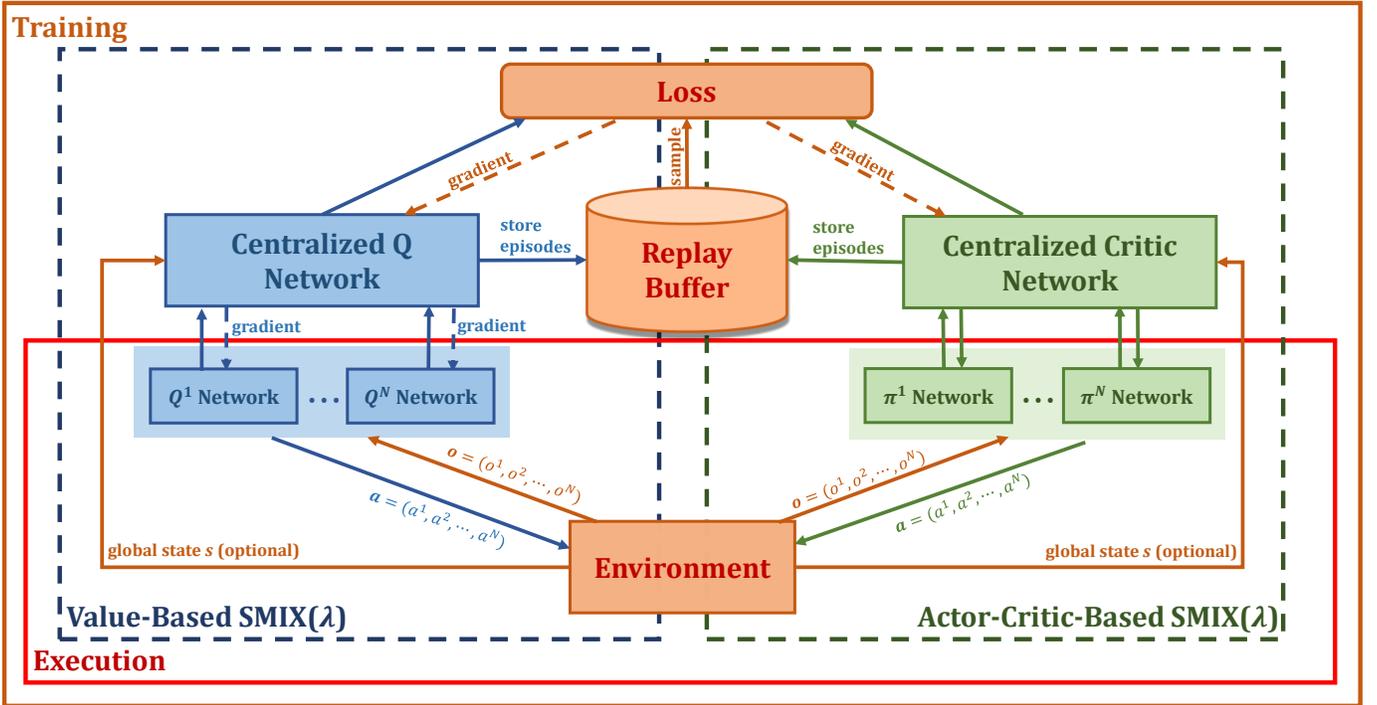


Fig. 1. Applying IS-free SARSA(λ) to the centralized value function estimation for value-based and actor-critic-based methods. (The left and right dashed boxes show the value-based SMIX(λ) and actor-critic based SMIX(λ) algorithms, and the two solid boxes represent the modules involved in the centralized training and decentralized execution respectively. Each agent’s Q network (or π network) only has access to its own observation, and the centralized Q network (or critic network) aggregates all agents’ observation information.)

where each ρ_i is a non-negative coefficient and satisfies $\prod_{i=1}^t \rho_i = 1$ when $t = 0$. The error term δ_t^π is generally written as the following expected TD-error,

$$\delta_t^\pi = \underbrace{r_{t+1} + \gamma \mathbb{E}_\pi Q(\tau_{t+1}, \cdot)}_{\text{1-step TD-target}} - Q(\tau_t, \mathbf{a}_t), \quad (8)$$

where $\mathbb{E}_\pi Q(\tau, \cdot) = \sum_{\mathbf{a}} \pi(\mathbf{a}|\tau) Q(\tau, \mathbf{a})$.³ In particular, for the importance sampling (IS) method, each ρ_i in (7) is defined as the relative probability of their trajectories occurring under the target policy π and behavior policy μ , also called importance sampling ratio, i.e., $\rho_i = \frac{\pi(\mathbf{a}_i|\tau_i)}{\mu(\mathbf{a}_i|\tau_i)}$.

Despite its theoretical soundness, the importance sampling (IS) method faces great challenges under the setting of multi-agent environments: 1) it suffers from large variance due to the product of the ratio [54], and 2) the “curse of dimensionality” issue of the joint action space makes it impractical to calculate the $\pi(\mathbf{a}_i|\tau_i)$ even for a single timestep i , when the number of agents is large. Previously, Liu et al. proposed a method that effectively addresses the first issue by avoiding calculating the product over the trajectories [54], but how to solve the second one remains open.

The above analysis highlights the need for exploring alternative approaches that can perform off-policy learning without

importance sampling in multi-agent settings. To achieve the above goal, the key idea of SMIX(λ) is to further simplify the coefficient ρ_i in (7), so as to reduce the variance of the importance sampling estimator and to potentially bypass the curse of dimensionality involved in calculating $\pi(\cdot|\tau)$.

Specifically, we relax each coefficient $\rho_i = 1.0$ in (7) use the λ -return [23] as the TD target estimator, which is defined as follows:

$$G_t^\lambda = (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} G_t^{(n)}, \quad (9)$$

where $G_t^{(n)} = r_{t+1} + \gamma r_{t+2} + \dots + \gamma^n \mathbb{E}_\pi Q(\tau_{t+n}, \mathbf{a}_{t+n}; \theta^-)$ is the n -step return and θ^- are the parameters of the target network.

Replacing 1-step TD-target in (8) with G_t^λ , and setting $\rho_i = 1.0$ for all i in (7), we have (the update step-size α is omitted for simplification),

$$Q(\tau, \mathbf{a}) \leftarrow Q(\tau, \mathbf{a}) + \mathbb{E}_\mu \left[\sum_{t \geq 0} \gamma^t (G_t^\lambda - Q(\tau_t, \mathbf{a}_t)) \right]. \quad (10)$$

In this paper, we call the method using (10) as the TD-target for off-policy value function estimation as “IS-free off-policy SARSA(λ)”. Next, we use this method to estimate the centralization value function in multi-agent reinforcement learning and analyze its theoretical properties of this method in section V.

³The policy evaluation strategy of many popular methods can be expressed as (7), including SARSA(λ) [51], off-policy importance sampling methods [32], off-policy $Q(\lambda)$ method [24], tree-backup method, TB(λ) [53] and Retrace(λ) [31]. These methods differ in the definition of the coefficient ρ_i and error term δ_t^π [24], [31].

C. The SMIX(λ) Algorithm

SMIX(λ) is trained end-to-end and the loss function for the centralized value function Q_{tot}^π has the following form:

$$\mathcal{L}_t(\theta) = \sum_{i=1}^{N_b} [(y_i^{tot} - Q_{tot}^\pi(\tau, \mathbf{a}; \theta))^2], \quad (11)$$

where $y_i^{tot} = G_t^\lambda$ is defined in (9) and is estimated through experience replaying, N_b is the batch size.

In implementation, SMIX(λ) use an experience replay [1] to store the most recent off-policy data. The experience replay usually stores a queue of experiences (tuples of [observation, action, reward, successor observation]). In SMIX(λ), however, each tuple in the experience replay stores one complete trajectory $(s_0, \tau_0, \mathbf{a}_0, r_1, \dots, s_{T-1}, \tau_{T-1}, \mathbf{a}_{T-1}, r_T, s_T)$ so as to evaluate the λ -return target.

The QMIX [15] structure is adopted as the basic deep network architecture for the proposed SMIX(λ). Each agent i has its own decentralized $Q^i(\tau^i, a^i)$ network composed of GRU [55] modules. Then all the individual Q^i values are passed into a mixing network to calculate the joint action-value Q_{tot}^π . The weight of the mixing network is generated by hypernetworks [56] using the global state s . All the neural networks are trained end-to-end and the centralized value function Q_{tot}^π is updated by minimizing the loss (11).

The general training procedure for SMIX(λ) is provided in Algorithm 1. Firstly, the replay buffer is filled with the trajectories and the oldest data is replaced when the buffer is full. Secondly, we sample a batch of episodes uniformly from the replay buffer to calculate the λ -return TD target. Then, the parameters of behavior network θ are updated by minimizing the loss function. Finally, we replace the target network's parameters θ^- with θ periodically and iterate the above process.

It is worth noting that our method of training a centralized value function is a general method and can be easily applied to other CTDE methods, include value-based methods (such as VDN [20]), actor-critic-based methods (e.g. COMA [14]), and even fully decentralized methods (e.g. independent Q-learning (IQL) [57]). Figure 1 gives the overall architecture of a generalized version of SMIX(λ).

D. Representational Complexity

The hypothesis space (or hypothesis set) \mathcal{H} is a space of all possible hypotheses for mapping inputs to outputs that can be searched [58], [59]. To learn a stable and generalizable CVF, choosing a suitable hypothesis space is essential, which is related not only to the characteristic of the problem domain but also to how the learned system is deployed. In particular, in multi-agent systems, all agents' joint action space increases exponentially with the increase of the number of agents, implying that the hypothesis space of CVF should be large enough to account for such complexity. However, to reduce optimization costs and enable decentralized execution, it is often necessary to impose some constraints on centrally valued functions. Figure 2(a) shows the relationship of hypothesis spaces under some different constraints.

The centralized value function Q_{tot} obtained by solving optimization problem (1) can represent any function that fits assumption 1. Our SMIX(λ) also uses the non-negative constraint during training. This makes QMIX [15] and SMIX(λ) share the same representational complexity of centralized value function. However, SMIX(λ) is more flexible than QMIX due to the decoupling of updating rule and CGB assumption during training. QTRAN [21] directly optimizes the joint action-value function, which gives this method a stronger representational complexity than QMIX and SMIX(λ). Figure 2(b) shows the relationship of representational complexity for several different algorithms.

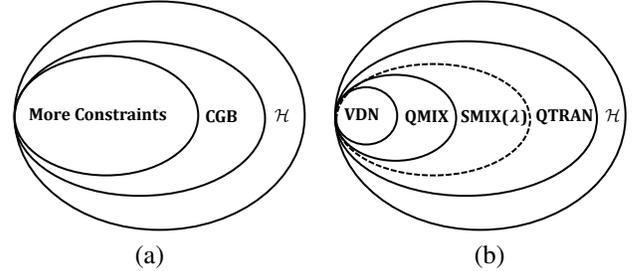


Fig. 2. (a) The size of hypothesis space corresponding to different constraints. (b) The relationship of representational complexity for several different algorithms.

We use an m -step cooperative matrix game [60] for two agents to illustrate the effects of representational complexity of QMIX, SMIX(λ) and QTRAN. Similar matrix games are usually used to study the algorithm's representational complexity [15], [21], [60]. In the m -step matrix game, zero rewards lead to termination, and the differentiating states are located at the terminal ends. Figure 3 illustrates the m -step matrix game for $m = 5$. The optimal policy is to take the top left joint action and finally take the bottom right action, giving an optimal total payoff of $m + 3$.

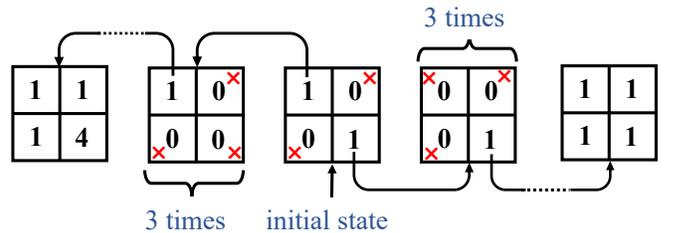


Fig. 3. m -step matrix game for $m = 5$ case.

Figure 4 gives the results on the m -step matrix game. One can see that although our SMIX(λ) learns the fast, QTRAN achieves the highest return. This is as expected as QTRAN has the most highest representational complexity among the compared ones, allowing it to achieve the lowest bias in this relatively simple scenario. However, with the increasing complexity of search space (as in SMAC [25]), a larger hypothesis space could turn out to be a disadvantage as the variance could dominate the generalization error, making the task of finding the best hypothesis quite challenging. In such cases, methods like ours would be a better choice, as illustrated in the experimental Section.

Algorithm 1 Training Procedure for SMIX(λ)

- 1: Initialize the behavior network with parameters θ , the target network with parameters θ^- , empty replay buffer \mathcal{D} to capacity $N_{\mathcal{D}}$, training batch size N_b
 - 2: **for** each training episode **do**
 - 3: **for** each episode **do**
 - 4: **for** $t = 1$ to $T - 1$ **do**
 - 5: Obtain the partial observation $\mathbf{o}_t = \{o_t^1, \dots, o_t^N\}$ for all agents and global state s_t
 - 6: Select action a_t^i according to ϵ -greedy policy w.r.t agent i 's decentralized value function Q^i for $i = 1, \dots, N$
 - 7: Execute joint action $\mathbf{a}_t = \{a^1, a^2, \dots, a^N\}$ in the environment
 - 8: Obtain the global reward r_{t+1} , the next partial observation o_{t+1}^i for each agent i and next global state s_{t+1}
 - 9: **end for**
 - 10: Store the episode in \mathcal{D} , replacing the oldest episode if $|\mathcal{D}| \geq N_{\mathcal{D}}$
 - 11: **end for**
 - 12: Sample a batch of N_b episodes $\sim \text{Uniform}(\mathcal{D})$
 - 13: Calculate λ -return targets y_i^{tot} according to (9) using θ^- for each timestep
 - 14: Update θ by minimizing $\sum_{t=1}^{T-1} \sum_{i=1}^{N_b} [(y_i^{tot} - Q_{tot}^{\pi}(\tau, \mathbf{a}; \theta))^2]$
 - 15: Replace target parameters $\theta^- \leftarrow \theta$ every C episodes
 - 16: **end for**
-

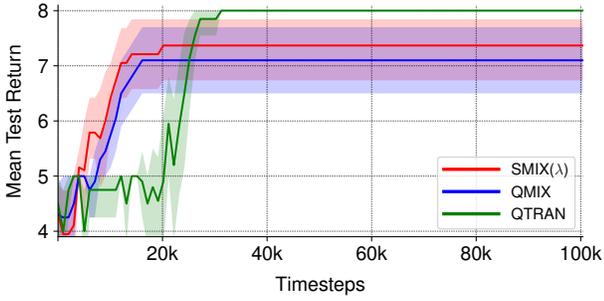


Fig. 4. Average return of QMIX, SMIX(λ) and QTRAN on 5-step matrix game for 100k training steps. (The mean and 95% confidence interval are shown across 20 independent runs.)

V. THERETICAL ANALYSIS

In this section, we give the convergence analysis of the proposed SMIX(λ) algorithm, by first building the connection between SMIX(λ) and a previous method named $Q(\lambda)$ [24], originally proposed for off-policy value function evaluation under single-agent settings.

Denoting G^{π} as the λ -return estimator (cf., 9) for the action-value of the target policy π , the goal of an off-policy method is to use the data from the behavior policy μ to correct G^{π} , in a way such that the following criterion is met,

$$\mathbb{E}_{\pi} [G^{\pi}] = \mathbb{E}_{\mu} [G^{\mu, \pi}], \quad (12)$$

where $G^{\mu, \pi}$ is the corrected return of off-policy data.

The most commonly used method for calculating the $G^{\mu, \pi}$ is the importance sampling (IS) method which multiply each reward with a weighted term to satisfy (12). Indeed, the motivation behind SMIX(λ) is to simplify the IS method so that it can be used in multi-agent settings. If we define the IS ratio at timestep t as: $\rho_t = \frac{\pi(\mathbf{a}_t | \tau_t)}{\mu(\mathbf{a}_t | \tau_t)}$, then the n -step return

using IS can be defined as:

$$\begin{aligned} G_t^{(n)} = & r_{t+1} + \gamma \rho_{t+1} r_{t+2} + \dots \\ & + \gamma^{n-1} \rho_{t+1} \dots \rho_{t+n-1} r_{t+n} \\ & + \gamma^n \rho_{t+1} \dots \rho_{t+n} \mathbb{E}_{\pi} Q^{\text{SMIX}}(\tau_{t+n}, \cdot). \end{aligned} \quad (13)$$

Thus, we have the following form of $G^{\mu, \pi}$:

$$G^{\mu, \pi} = (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} G_t^{(n)}. \quad (14)$$

Plugging (13) into (14), we have:

$$\begin{aligned} G^{\mu, \pi} \leftarrow & Q^{\text{SMIX}}(\tau_t, \mathbf{a}_t) + \sum_{k=t}^{\infty} \left(\prod_{i=t+1}^k \gamma \lambda \rho_i \right) \delta_k^{\pi}, \\ \delta_k^{\pi} = & (r_{k+1} + \gamma \rho_{k+1} Q^{\text{SMIX}}(\tau_{k+1}, \mathbf{a}_{k+1}) - Q^{\text{SMIX}}(\tau_k, \mathbf{a}_k)). \end{aligned} \quad (15)$$

In SMIX(λ), all the importance sampling factor are relaxed to 1.0, then corresponding to (7), the update rule of SMIX(λ) can be expressed as⁴,

$$\begin{aligned} Q^{\text{SMIX}}(\tau_t, \mathbf{a}_t) \leftarrow & Q^{\text{SMIX}}(\tau_t, \mathbf{a}_t) + \mathbb{E}_{\mu} \left[\sum_{k=t}^{\infty} \left(\prod_{i=t+1}^k \lambda \gamma \right) \delta_k^{\pi} \right], \\ \delta_k^{\pi} = & \left(r_{k+1} + \gamma \mathbb{E}_{\mu} Q^{\text{SMIX}}(\tau_{k+1}, \cdot) - Q^{\text{SMIX}}(\tau_k, \mathbf{a}_k) \right). \end{aligned} \quad (16)$$

In contrast with the multiplicative operation for off-policy learning, an additive-type operation is used in $Q(\lambda)$ [24]. In particular, an additive correction term $\Delta_r^{\mu, \pi}$, is added to each reward when calculating $G^{\mu, \pi}$ in (12) and get:

$$\begin{aligned} G_t^{(n)} = & (r_{t+1} + \Delta_{r_{t+1}}^{\mu, \pi}) + \dots + \gamma^{n-1} (r_{t+n} + \Delta_{r_{t+n}}^{\mu, \pi}) \\ & + \gamma^n \mathbb{E}_{\pi} Q^{Q(\lambda)}(\tau_{t+n}, \cdot). \end{aligned} \quad (17)$$

The major advantage of this additive off-policy correction is that there is no product of the ratio and no the joint policy

⁴We consider the expected form of $Q^{\text{SMIX}}(\tau_{k+1}, \mathbf{a}_{k+1})$ in (15) and the training data is sampled from a replay buffer.

$\pi(\mathbf{a}|\boldsymbol{\tau})$ involved, hence completely bypassing the limitations of the IS method⁵.

Specifically, the updating rule of $Q(\lambda)$ method is [24]:

$$Q^{Q(\lambda)}(\boldsymbol{\tau}_t, \mathbf{a}_t) \leftarrow Q^{Q(\lambda)}(\boldsymbol{\tau}_t, \mathbf{a}_t) + \mathbb{E}_{\boldsymbol{\mu}} \left[\sum_{k=t}^{\infty} \left(\prod_{i=t+1}^k \lambda \gamma \right) \hat{\delta}_k^{\pi} \right],$$

$$\hat{\delta}_k^{\pi} = (r_{k+1} + \Delta^{\boldsymbol{\mu}, \pi} r_{k+1}),$$

$$\Delta^{\boldsymbol{\mu}, \pi} r_{k+1} = \gamma \mathbb{E}_{\pi} Q^{Q(\lambda)}(\boldsymbol{\tau}_{k+1}, \cdot) - Q^{Q(\lambda)}(\boldsymbol{\tau}_k, \mathbf{a}_k). \quad (18)$$

By comparing (16) and (18), we see that our SMIX(λ) and off-policy $Q(\lambda)$ are essentially equivalent except that SMIX(λ) calculates $\mathbb{E}_{\boldsymbol{\mu}} Q^{\text{SMIX}}(\boldsymbol{\tau}_{k+1}, \cdot)$ in δ_k^{π} while $Q(\lambda)$ calculate $\mathbb{E}_{\pi} Q^{Q(\lambda)}(\boldsymbol{\tau}_{k+1}, \cdot)$ in $\hat{\delta}_k^{\pi}$.

The following theorem states that when π and $\boldsymbol{\mu}$ are sufficiently close, the difference between the output of SMIX(λ) and $Q(\lambda)$ is bounded. This implies that SMIX(λ) is consistent with the $Q(\lambda)$ algorithm.

Theorem 3. *Suppose we update the value function from $Q_n^{\text{SMIX}}(\boldsymbol{\tau}_t, \mathbf{a}_t) = Q_n^{Q(\lambda)}(\boldsymbol{\tau}_t, \mathbf{a}_t)$, where n represents the n -th update. Let $\epsilon = \max_{\boldsymbol{\tau}} \|\pi(\cdot|\boldsymbol{\tau}) - \boldsymbol{\mu}(\cdot|\boldsymbol{\tau})\|_1$, $M = \max_{\boldsymbol{\tau}, \mathbf{a}} |Q_n^{Q(\lambda)}(\boldsymbol{\tau}, \mathbf{a})|$. Then, the error between $Q_{n+1}^{\text{SMIX}}(\boldsymbol{\tau}_t, \mathbf{a}_t)$ and $Q_{n+1}^{Q(\lambda)}(\boldsymbol{\tau}_t, \mathbf{a}_t)$ can be bounded by the expression:*

$$|Q_{n+1}^{\text{SMIX}}(\boldsymbol{\tau}_t, \mathbf{a}_t) - Q_{n+1}^{Q(\lambda)}(\boldsymbol{\tau}_t, \mathbf{a}_t)| \leq \frac{\epsilon \gamma}{1 - \lambda \gamma} M. \quad (19)$$

Proof: First, we have,

$$\begin{aligned} \left| \delta_t^{\pi} - \hat{\delta}_t^{\pi} \right| &= \left| \gamma \mathbb{E}_{\boldsymbol{\mu}} Q_n^{\text{SMIX}}(\boldsymbol{\tau}_{t+1}, \cdot) - \gamma \mathbb{E}_{\pi} Q_n^{Q(\lambda)}(\boldsymbol{\tau}_{t+1}, \cdot) \right| \\ &= \gamma \left| \sum_{\mathbf{a}} \boldsymbol{\mu}(\mathbf{a}|\boldsymbol{\tau}_{t+1}) Q_n^{\text{SMIX}}(\boldsymbol{\tau}_{t+1}, \cdot) - \sum_{\mathbf{a}} \pi(\mathbf{a}|\boldsymbol{\tau}_{t+1}) Q_n^{Q(\lambda)}(\boldsymbol{\tau}_{t+1}, \cdot) \right| \\ &\leq \gamma \epsilon M. \end{aligned}$$

Thus,

$$\begin{aligned} &\left| Q_{n+1}^{\text{SMIX}}(\boldsymbol{\tau}_t, \mathbf{a}_t) - Q_{n+1}^{Q(\lambda)}(\boldsymbol{\tau}_t, \mathbf{a}_t) \right| \\ &= \left| \mathbb{E}_{\boldsymbol{\mu}} \left[\sum_{k=t}^{\infty} \left(\prod_{i=t+1}^k \lambda \gamma \right) \delta_k^{\pi} \right] - \mathbb{E}_{\boldsymbol{\mu}} \left[\sum_{k=t}^{\infty} \left(\prod_{i=t+1}^k \lambda \gamma \right) \hat{\delta}_k^{\pi} \right] \right| \\ &= \left| \mathbb{E}_{\boldsymbol{\mu}} \left[\sum_{k=t}^{\infty} \left(\prod_{i=t+1}^k \lambda \gamma \right) (\delta_k^{\pi} - \hat{\delta}_k^{\pi}) \right] \right| \\ &\leq \left| \mathbb{E}_{\boldsymbol{\mu}} \left[\sum_{k=t}^{\infty} \left(\prod_{i=t+1}^k \lambda \gamma \right) (\gamma \epsilon M) \right] \right| \\ &\leq \mathbb{E}_{\boldsymbol{\mu}} \left[\frac{1}{1 - \lambda \gamma} (\gamma \epsilon M) \right] = \frac{\epsilon \gamma}{1 - \lambda \gamma} M. \end{aligned}$$

Therefore, the expression (19) holds. \blacksquare

This theorem indicates that SMIX(λ) has the similar convergence property to $Q(\lambda)$ when the difference between behavior policy $\boldsymbol{\mu}$ and target policy π is bounded by ϵ , which is $\epsilon = \max_{\boldsymbol{\tau}} \|\pi(\cdot|\boldsymbol{\tau}) - \boldsymbol{\mu}(\cdot|\boldsymbol{\tau})\|_1$. In practice, this condition can

⁵But under the condition that the behavior policy $\boldsymbol{\mu}$ should be close to the target policy π , which under our experience replay setting should not be a problem (cf., [36]).

be easily implemented by periodically replacing the current behavior policy with the old target policy and limiting the size of the replay memory. The following theorem presents the convergence property of the $Q(\lambda)$ method [24].

Theorem 4. [24] *Consider the sequence of Q -functions computed according to (18) with fixed policy π and $\boldsymbol{\mu}$. Let $\epsilon = \max_{\boldsymbol{\tau}} \|\pi(\cdot|\boldsymbol{\tau}) - \boldsymbol{\mu}(\cdot|\boldsymbol{\tau})\|_1$. If $\lambda \epsilon < \frac{1-\gamma}{\gamma}$, then under the following conditions:*

- $\sum_{t \geq 0} \mathbb{P}\{\boldsymbol{\tau}_t, \mathbf{a}_t = \boldsymbol{\tau}, \mathbf{a}\} \geq D > 0$, where $\mathbb{P}(\boldsymbol{\tau}, \mathbf{a})$ represents the visit frequency,
- $\mathbb{E}_{\boldsymbol{\mu}_n} T_n^2 < \infty$, where T_n is the length of $\boldsymbol{\tau}_n$,
- $\sum_{n \geq 0} \alpha_n(\boldsymbol{\tau}, \mathbf{a}) = \infty, \sum_{n \geq 0} \alpha_n^2(\boldsymbol{\tau}, \mathbf{a}) < \infty$, where α_n is the step-size of the n -th iteration,

we have, almost surely:

$$\lim_{n \rightarrow \infty} Q_n^{Q(\lambda)}(\boldsymbol{\tau}, \mathbf{a}) = Q^{\pi}(\boldsymbol{\tau}, \mathbf{a}).$$

The above analysis shows that SMIX(λ) and $Q(\lambda)$ have similarities both formally and analytically. However, when applying them to the problem of multi-agent reinforcement learning, their computational complexity is fundamentally different. This is because to calculate the additive error correction term, $Q(\lambda)$ has to estimate the expectation over target policy π in (18), but this is unrealistic in the multi-agent setting since the dimension of the joint action space grows exponentially with the number of agents. By contrast, the SMIX(λ) relies on the experience replay technique to compute the expectation in (16), whose computational complexity grows only linearly with the number of training samples, regardless of the size of joint action space and the number of agents involved. Such scalability makes our method more appropriate for the task of multi-agent reinforcement learning, compared to the $Q(\lambda)$ and QTRAN (which suffers from the same problem as $Q(\lambda)$).

Finally, before ending this section, we summarize some of the key characteristics of QMIX and SMIX(λ) in Table I.

TABLE I
THE COMPARISON OF QMIX AND SMIX(λ).

Property	QMIX	SMIX(λ)
Uses experience replay	✓	✓
CVF estimation	Q-learning based	Expected-SARSA based
TD target	one-step return	λ -return
Algorithm's Flexibility	CGB assumption	No explicit CGB
Stable point of convergence	Q^*	Q^*

VI. EXPERIMENTS

In this section, we first describe the environmental setup and the implementation details of our method. Then we give the experimental results and ablation study. The code of SMIX(λ) is available at: <https://github.com/chaovven/SMIX>.

A. Environmental Setup

We evaluate our SMIX(λ) in the StarCraft Multi-Agent Challenge (SMAC) [25] environment. The SMAC is chosen as our testbed mainly because of the following two reasons: (1)

TABLE II
THE SCENARIOS CONSIDERED IN OUR EXPERIMENTS.

Name	Ally Units	Enemy Units	Type
3m	3 Marines	3 Marines	homogeneous & symmetric
8m	8 Marines	8 Marines	homogeneous & symmetric
2s3z	2 Stalkers & 3 Zealots	2 Stalkers & 3 Zealots	heterogeneous & symmetric
3s5z	3 Stalkers & 5 Zealots	3 Stalkers & 5 Zealots	heterogeneous & symmetric
2m_vs_1z	2 Marines	1 Zealot	asymmetric
2s_vs_1sc	2 Stalkers	1 Spine Crawler	asymmetric
3s_vs_3z	3 Stalkers	3 Zealots	asymmetric
1c3s5z	1 Colossi, 3 Stalkers & 5 Zealots	1 Colossi, 3 Stalkers & 5 Zealots	heterogeneous & symmetric
MMM	1 Medivac, 2 Marauders & 7 Marines	1 Medivac, 2 Marauders & 7 Marines	heterogeneous & symmetric

SMAC provides a set of rich cooperative scenarios that challenge algorithms to handle significant partial observability and credit assignment problem [14]. These problems bring a great challenge for centralized value function estimation. (2) SMAC also provides an open-source Python-based implementation of several key algorithms, which allows for fair comparisons between different methods.

SMAC is based on the popular real-time strategy (RTS) game StarCraft II⁶. Each unit can be seen as an individual agent which has a complex set of micro-actions. Different from the full StarCraft II game, SMAC focuses on fully cooperative, decentralized *micromanagement* multi-agent problems. *Micromanagement* means the task of controlling individual or grouped units to fight enemy units. While high-level strategies such as economy and resource management, known as *macromanagement*, are not considered in SMAC.

SMAC provides several challenging micro scenarios that aim to evaluate different aspects of cooperative behaviors of a group of agents. In each scenario, two groups of agents are placed on the map with random initial positions within groups at the beginning of each episode. Each agent can only receive local observations within its *sight range*, which introduces significant partial observability. Extra global state information is available during centralized training. The units of the first group (allied units) are controlled by decentralized agents, while the units of the other group (enemy units) are controlled by built-in heuristic game AI bot with difficulty ranging from *very easy* to *cheat insane*. In our experiments, we set the difficulty of the game AI bot to *very difficult* for all experiments. Available actions for each agent contain: `move[direction]`, `attack[enemy id]`, `stop`, and `noop`. Agents receive a joint reward equal to the total damage dealt on the enemy units. We use the default setting for the reward. Refer to [25] for more details.

The following 3 types of scenarios are considered in our experiments: (1) *homogeneous and symmetric units*, (2) *heterogeneous and symmetric units*, (3) *asymmetric units*. The list of scenarios considered in our experiment is presented in Table II. Figure 5 shows the screenshots of two SMAC scenarios used in our experiments.



(a) 3 Stalkers vs 5 Zealots

(b) 8 Marines vs 8 Marines

Fig. 5. Screenshots of two SMAC scenarios.

We use *test win rate* as the evaluation metric, which is proposed in [25] and is the default evaluation metric in the SMAC environment. The test win rate is evaluated in the following procedure: the training process is interrupted after every 20,000 timesteps, then 24 independent test episodes are run with each agent performing greedy action selection in a decentralized way. Test win rate refers to the percentage of episodes where the agents defeat all enemy units within the time limit.

B. Implementation Details

The agent network architecture of SMIX(λ) consists of a 64-dimensional GRU [55]. One 64-dimensional fully connected layer with ReLU activation function before GRU is applied for processing the input. The layer after GRU is a fully connected layer of 64 units, which outputs the decentralized action-values $Q^i(\tau, \cdot)$ of agent i . All agent networks share parameters for reducing the number of parameters to be learned. Thus the agent's one-hot index i is concatenated onto each agent's observations. The agent's previous action is also concatenated to the input.

Based on the basic network architecture of QMIX [15], SMIX(λ) performs the centralized value function estimation with λ -return ($\lambda = 0.8$) calculated from a batch of 32 episodes. The batch is sampled uniformly from a replay buffer that stores the most recent 1500 episodes. We run 4 episodes simultaneously. Then we perform training on those fully unrolled episodes. The target network is updated after every 200 training episodes. λ is set to 0.8.

The ϵ -greedy method is used in the training procedure for exploration. ϵ is annealed linearly from 1.0 to 0.05 across the

⁶StarCraft II is a trademark of Blizzard EntertainmentTM.

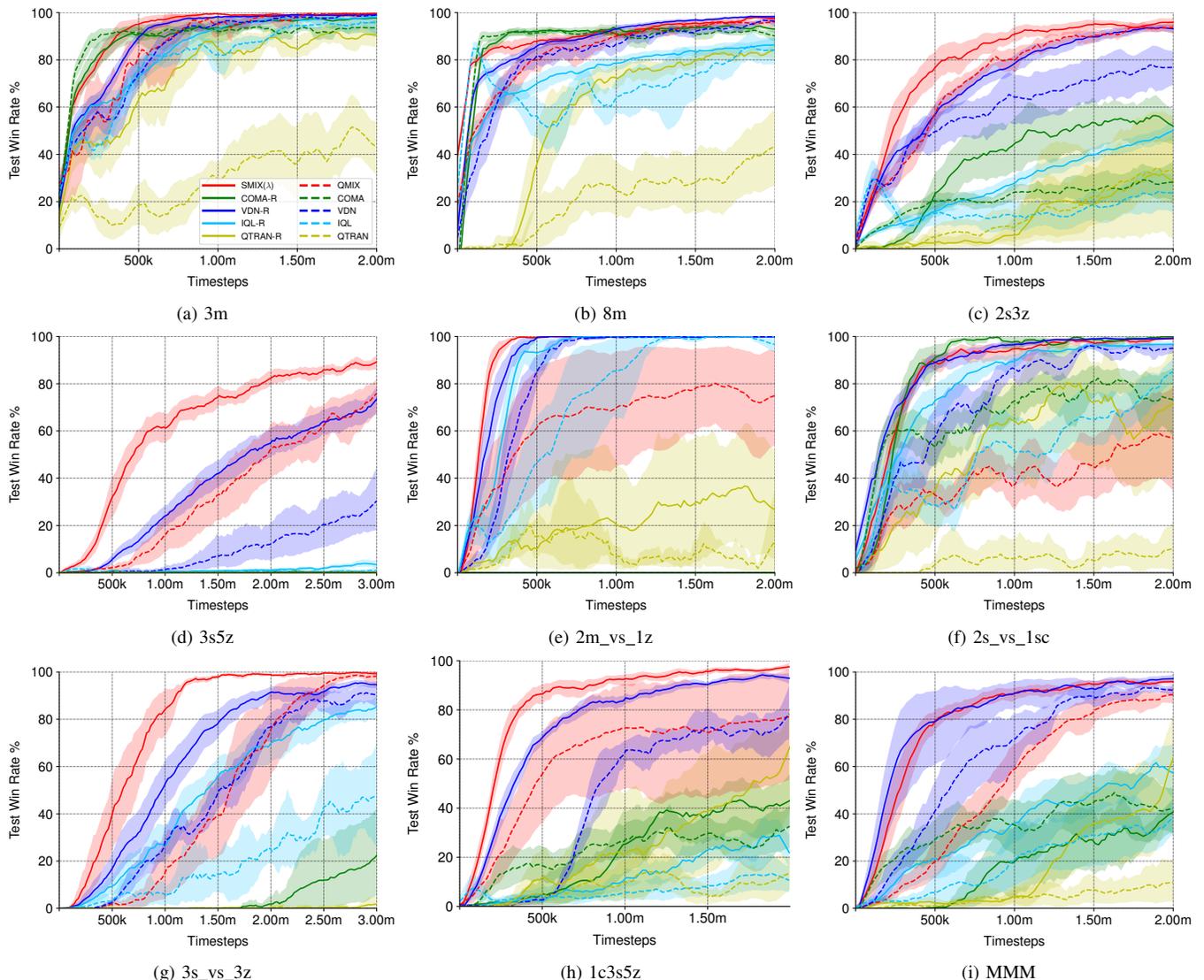


Fig. 6. Test win rates for SMIX(λ), revised methods (COMA-R, VDN-R, IQL-R, QTRAN-R) and comparison methods (QMIX, COMA, VDN, IQL, QTRAN) in nine different scenarios. Revised methods are created by replacing original methods’ CVF estimation procedures with the one adopted by SMIX(λ). The performance of our method and revised methods and plotted with solid line, with their counterparts shown with dashed lines of the same color. The mean and 95% confidence interval are shown across 10 independent runs. The legend in (a) applies across all plots.

first 50k timesteps for all experiments. The discount factor γ is set to 0.99, and the RMSprop optimizer is used with learning rate $lr = 0.0005$ and $\alpha = 0.99$ without weight decay or momentum during training.

C. Comparative Evaluation

We compare our SMIX(λ) with state-of-the-art algorithms QMIX [15] and COMA [14], which currently perform the best on the SMAC benchmark. VDN [20] and IQL [57] are chosen as baselines for comparisons. QTRAN [21] is also implemented in SMAC for comprehensive comparison.

The results of all methods in the training process are plotted in Figure 6 and we also provide quantitative comparisons of our methods and their counterparts after training for 1 million steps in Table III. Overall, SMIX(λ) significantly outperforms all the comparison methods in heterogeneous or asymmetric scenarios (i.e., scenarios except 3m and 8m), while performing

comparably to them in homogeneous and symmetric scenarios (i.e., 3m and 8m) both in terms of the learning speed and final performance.

In homogeneous and symmetric scenarios such as 3m and 8m, COMA is only slightly faster than SMIX(λ) but underperforms SMIX(λ) in terms of the final performance. In asymmetric (e.g., 3s_vs_3z, 2s_vs_1sc) or heterogeneous (e.g., 2s3z, 3s5z, 1c3s5z) maps, COMA fails to solve these scenarios effectively. It’s worth noting that QTRAN does not work very well in these complex environments, which is also consistent with the experimental results of other authors [60]. One possible reason is that QTRAN needs to address a very large optimization problem in the high-dimensional joint action space. This highlights the importance of improving the performance in value function estimation.

Due to the poor performance of COMA and QTRAN, QMIX can be seen as the state-of-the-art on this benchmark.

However, the learning speed of $\text{SMIX}(\lambda)$ is almost twice as fast as QMIX. In 3s5z, $\text{SMIX}(\lambda)$ (solid red line) achieves a nearly 90% win rate, while the best comparison method QMIX (dotted red line) achieves about only 70% test win rate. In 2s_vs_1sc, $\text{SMIX}(\lambda)$ also requires less than half the number of samples of QMIX and other comparison methods to reach the asymptotic performance. The largest performance gap can be seen in 3s_vs_3z map (Figure 6g). QMIX needs to be trained for nearly 3 million timesteps to achieve a 100% test win rate, while half of the timesteps are sufficient for $\text{SMIX}(\lambda)$ to achieve the same win rate. In MMM, we can find an interesting result that the VDN can achieve better performance than QMIX (see Figure 6i). This indicates that a simpler network structure can also have enough representative capacity and the reason for VDN’s superior performance is that a simpler network architecture only needs a relatively small number of samples for training. Furthermore, by incorporating the proposed centralized training method, VDN’s performance can be further improved, which implies that the bottleneck of VDN and QMIX may be the bias and variance of the CVF estimation. On the whole, the superior performance of $\text{SMIX}(\lambda)$ using λ -return with off-policy episodes presents a clear benefit over the one-step estimation of QMIX.

D. Generalizing $\text{SMIX}(\lambda)$ to Other MARL Algorithms

$\text{SMIX}(\lambda)$ focuses on centralized value function evaluation with λ -return calculated from off-policy episodes. This method could ideally be generalized to other MARL algorithms incorporating critic estimation, including critic-only and actor-critic algorithms.

To demonstrate the benefits of our approach, we generalize the CVF estimation procedure of $\text{SMIX}(\lambda)$ to the following algorithms: COMA, VDN, IQL and QTRAN. We achieve these by replacing their original value function estimation procedure with ours (see Section IV), yielding four revised algorithms called COMA-R, VDN-R, IQL-R and QTRAN-R respectively. Figure 6 gives the comparisons between our methods and their counterparts. Overall, most of the extended methods (solid line) perform on par or significantly better than their counterparts (in the same color but dashed line) in most scenarios both in terms of the final win rate and learning speed.

VDN-R considerably improves the performance of VDN. Especially under challenging scenarios such as 3s5z, VDN-R achieves about 75% final win rate, which is more than twice as that of VDN (nearly 30%). Such improvement may be contributed to λ -return and the independence of the unrealistic centralized greedy assumption during the learning phase.

Furthermore, we find that VDN-R performs even better than QMIX in most scenarios. Recall that VDN uses a linear combination of decentralized Q-values (and so does our VDN-R) and QMIX extends VDN by combining decentralized Q-values in a non-linear way. Thus, QMIX can represent a richer class of CVF than VDN. However, our results indicate that the performance bottleneck of VDN may not be the limited representational capacity, but how to effectively balance the bias and variance in the estimation of CVF.

Similarly, by comparing QTRAN-R and QTRAN in Figure 6, one can see significant performance improvement by

replacing QTRAN’s CVF with ours, indicating that the CVF estimation plays an essential role in the QTRAN system and our proposed method is beneficial.

Also, COMA-R improves COMA’s performance in most scenarios, which can be considered as a success of utilizing the off-policy data, as COMA also adopts λ -return but uses only the on-policy data.

Another observation is that our method also works for IQL, which is a fully decentralized MARL algorithm. This suggests that our method is not limited to *centralized* value function estimation but also applicable to *decentralized* cases.

It is worth mentioning that the extended methods may not make improvements if the original methods do not work, e.g., QTRAN, COMA, IQL, and their counterparts do not work in 3s5z scenario (see Figure 6d and Table III). The reason may be that the main limitations of QTRAN, COMA and IQL on 3s5z do not lie in the inaccurate value function estimation, but rather in other problems, e.g., scaling not well to a large number of agents and multi-agent credit assignment problem.

E. Sample Efficiency Analysis

To study the sample efficiency of $\text{SMIX}(\lambda)$, we fix the amount of samples used in training procedure to compare the sample efficiency of different algorithms. The quantitative comparisons of different methods after training for 1 million steps (1 million interactions with the environment) are presented in Table III.

Overall, $\text{SMIX}(\lambda)$ achieves the best performance consistently across all scenarios among compared algorithms, after being trained from the same number of interaction with the environment. This can be partially explained by the fact that our algorithm is off-policy by design. In other words, it has the capability to learn from the previous experience/policies, hence improving the sample efficiency.

F. Ablation Study

We perform the ablation experiments to investigate the necessity of balancing the bias and variance and the influence of utilizing the off-policy data.

λ -Return vs. n-Step Returns. To investigate the necessity of balancing the bias and variance in multi-agent problems, we adjust the parameter λ , where larger λ corresponds to smaller bias and larger variance whereas smaller λ indicates the opposite. Especially, $\lambda = 0$ is equivalent to *one-step return* (corresponding to the largest bias and the smallest variance); $\lambda = 1$ is equivalent to *Monte-Carlo (MC) return* (∞ -step, corresponding to the smallest bias and the largest variance). We also evaluate a variant named $\text{SMIX}(n)$, which uses n-step return in place of λ -return as the TD target, i.e., $y_t^{\text{tot}} = \sum_{i=1}^n \gamma^{i-1} r_{t+i} + \gamma^n Q(\tau_{t+n}, \mathbf{a}_{t+n}; \theta^-)$.

As Figure 7a and 7d show, $\text{SMIX}(\lambda)$ with $\lambda = 0.8$ consistently achieves the best performance in selected scenarios. The method with $\lambda = 1$ (MC return, blue line) performs the worst in 3s5z, while $\lambda = 0$ (one-step return, green line) performs the worst in 2s_vs_1sc. These results reveal that the large variance of MC return or large bias of one-step return may degrade the performance. Similar results could also be

TABLE III
MEAN, STANDARD DEVIATION, AND MEDIAN OF TEST WIN RATE PERCENTAGES AFTER TRAINING FOR 1 MILLION TIMESTEPS IN NINE DIFFERENT SCENARIOS.

Algorithms		SMIX(λ)	QMIX	COMA-R	COMA	VDN-R	VDN	IQL-R	IQL	QTRAN-R	QTRAN
3m	mean \pm std	99 (± 0)	95 (± 3)	93 (± 8)	92 (± 2)	98 (± 0)	95 (± 2)	91 (± 4)	83 (± 9)	84 (± 10)	29 (± 11)
	median	99	95	97	93	98	95	94	86	86	29
8m	mean \pm std	91 (± 3)	90 (± 3)	92 (± 2)	90 (± 2)	94 (± 3)	86 (± 5)	80 (± 5)	59 (± 15)	72 (± 6)	24 (± 17)
	median	90	89	93	91	93	87	79	58	72	21
2s3z	mean \pm std	90 (± 4)	81 (± 7)	44 (± 18)	24 (± 6)	78 (± 14)	64 (± 16)	32 (± 8)	14 (± 10)	5 (± 8)	9 (± 10)
	median	91	81	47	24	79	71	31	13	2	5
3s5z	mean \pm std	61 (± 11)	16 (± 12)	0 (± 0)	0 (± 0)	29 (± 12)	1 (± 2)	0 (± 0)	0 (± 0)	0 (± 0)	0 (± 0)
	median	62	11	0	0	26	0	0	0	0	0
2m_vs_1z	mean \pm std	99 (± 0)	69 (± 38)	0 (± 0)	0 (± 0)	99 (± 0)	99 (± 0)	99 (± 0)	85 (± 27)	20 (± 29)	7 (± 4)
	median	100	99	0	0	99	99	100	99	4	6
2s_vs_1sc	mean \pm std	94 (± 5)	39 (± 19)	97 (± 4)	77 (± 11)	96 (± 2)	86 (± 8)	92 (± 6)	51 (± 22)	63 (± 24)	6 (± 11)
	median	96	45	100	78	97	88	94	54	70	0
3s_vs_3z	mean \pm std	84 (± 14)	15 (± 20)	0 (± 0)	0 (± 0)	67 (± 25)	27 (± 9)	35 (± 21)	5 (± 4)	0 (± 0)	0 (± 0)
	median	88	9	0	0	83	27	31	6	0	0
1c3s5z	mean \pm std	92 (± 3)	72 (± 32)	24 (± 19)	27 (± 13)	84 (± 3)	61 (± 5)	11 (± 5)	5 (± 6)	21 (± 29)	9 (± 10)
	median	93	88	18	28	85	61	11	3	9	8
MMM	mean \pm std	91 (± 4)	59 (± 15)	20 (± 17)	34 (± 18)	91 (± 9)	72 (± 17)	35 (± 14)	20 (± 19)	5 (± 4)	4 (± 5)
	median	91	61	22	39	94	78	34	15	2	3

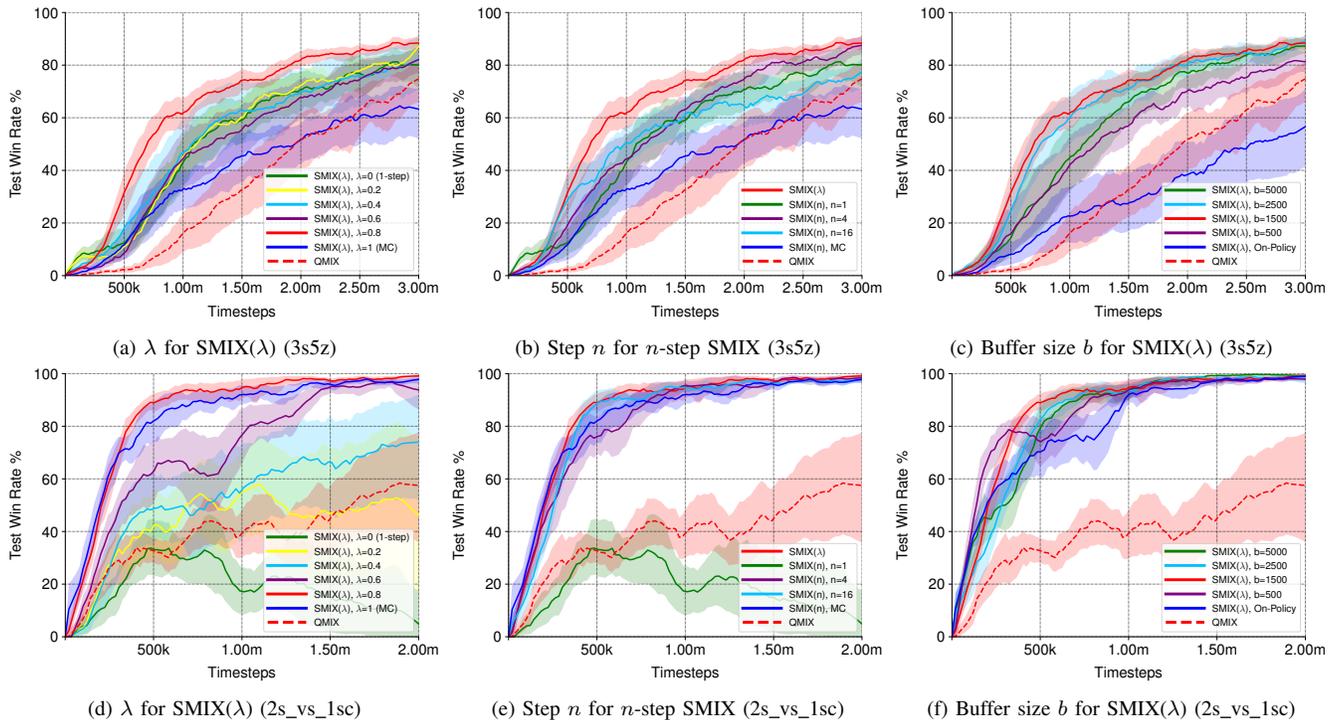


Fig. 7. Sensitivity of SMIX(λ) to selected hyperparameters in two different scenarios. The mean and 95% confidence interval is shown across 10 independent runs. The performance of the baseline (QMIX) is shown as a dashed red line. (a) and (d) show the sensitivity of SMIX(λ) to the value of λ ; (b) and (e) show the results using n -step TD with different backup steps; (c) and (f) show the comparison between SMIX(λ) and its on-policy version.

seen in SMIX(n) (Figure 7b and 7e), where SMIX(n) with $n = 4$ performs the best in 3s5z, while the one with $n = 16$ performs the best in 2s_vs_1sc. It is not easy to find the same n for SMIX(n) as SMIX(λ) which sets $\lambda = 0.8$ and performs consistently well across different maps. These results are consistent with the results in literature [61] that multi-step returns are beneficial for experience replay based algorithms. In summary, it is necessary to balance the trade-off between bias and variance in multi-agent problems, and λ -return could

serve as a convenient method to achieve such a trade-off.

Incorporating Off-Policy Data vs. Pure On-Policy Data.

To investigate the influence of utilizing the off-policy data, we perform experiments to compare SMIX(λ) against its on-policy version by scaling the size of the replay buffer. The on-policy version of SMIX(λ) corresponds to SMIX(λ) with buffer size $b = 4$ (the most recent 4 episodes in the replay buffer are all on-policy data), while the off-policy SMIX(λ) are the ones with buffer size $b > 4$, where the percentage of

off-policy data increases with the size of the replay buffer.

As shown in Figure 7c and 7f, all the variants of SMIX(λ) incorporating off-policy data ($b > 4$) perform better than the on-policy version ($b = 4$) in selected scenarios. Notably, the performance of SMIX(λ) with $b = 1500$ is almost twice that of the on-policy version both in terms of the final win rate and learning speed in 3s5z. Note that 3s5z (8 units) map is more complex than 2s_vs_1sc (2 units) in terms of the number of agents, and consequently, the joint action space of the former is much larger. However, more off-policy data does not always lead to better performance, as the method with $b = 5000$ (green line) performs worse than the one with $b = 1500$ (solid red line) in both scenarios. Actually, the buffer size is corresponding to the ϵ in Theorem 3 which measures the mismatch between the target policy π and the behavior policy μ . A smaller buffer size makes SMIX(λ) less sample efficient but a larger buffer size results in a looser error bound which biases the CVF estimation. This may explain why the performance degrades once the buffer size exceeds a threshold value. And our experimental results suggest that a moderate buffer size of 1500 could be a good candidate.

Effects of Non-negative Constraints. To illustrate the rationality of assumption 1, we compare the performance of SMIX(λ) and QMIX with their ablations without this non-negative constraint.

Figure 8 gives comparative results of SMIX(λ), and QMIX with their ablations without “positive weights”. We can see that all these algorithms suffer from drastic performance fluctuation when the additive constraint is removed. This means that non-negative constraint does play a positive role in multi-agent learning. This figure also shows that even without the requirement of non-negative weights, our SMIX(λ) is still able to achieve competitive performance with the non-negative version of QMIX.

TABLE IV
THE SCALABILITY OF SMIX(λ) AND QMIX AFTER TRAINING FOR 1 MILLION TIMESTEPS.

Algorithms		SMIX(λ)	QMIX
3m	mean \pm std	99 (± 0)	95 (± 3)
	median	99	95
8m	mean \pm std	91 (± 3)	90 (± 3)
	median	90	89
25m	mean \pm std	75 (± 26)	30 (± 17)
	median	93	24

Scalability. The results in Table IV show the scalability of SMIX(λ) and QMIX after training for 1 million steps. Overall, the performance of both methods decreases along with the increasing number of agents. However, our SMIX(λ) still outperforms QMIX, especially in hard scenarios 25m. Specifically, with 3 agents (the 3m map), SMIX(λ) achieves the best performance among the compared methods with a 99% win rate. By increasing the number of agents to 8 (the 8m map), the performance of all the methods decreases due to

the higher degree of challenging of the task, while our method still performs best among the compared ones⁷. Finally, when the number of agents been increased to 25 (the 25m map), the performance of QMIX decreases dramatically, which is not the case for SMIX(λ). These results show that the centralized value function estimation method used in SMIX(λ) method has better scalability and performs more robust in challenging tasks than QMIX. Finally, it is worth mentioning that for our experiments with up to 25 agents, the joint action space would be as large as $|\mathcal{A}|^{25}$, which imposes a great challenge to any MARL method.

VII. CONCLUSIONS & FUTURE WORK

One of the central challenges in multi-agent reinforcement learning with CTDE settings is to estimate the centralized value function. However, the sparse experiences and unstable nature of the multi-agent environments make this become a challenging task. To address this issue, we present the SMIX(λ) approach, by enhancing the quality of centralized value function in three aspects: (1) removing the greedy assumption to help to learn a more flexible functional structure, (2) using off-policy learning to alleviate the problem of sparse experiences and to improve exploration, and (3) using λ -return to balance the bias and variance of the algorithm. Our analysis indicates that SMIX(λ) has nice convergence guarantee through off-policy learning without importance sampling, which gives potential advantages in multi-agent settings. It is shown that the proposed SMIX(λ) approach significantly outperforms several state-of-the-art CTDE methods on the benchmark of the StarCraft Multi-Agent Challenge, and is beneficial to other CTDE methods as well by replacing their centralized value function estimator with ours.

Our future work will focus on incorporating the communication and opponent modeling methods into SMIX(λ) to further tackle the non-stationarity issue during the execution. We also aim to make SMIX(λ) perform more efficiently in dealing with a large numbers of agents.

REFERENCES

- [1] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski *et al.*, “Human-level control through deep reinforcement learning,” *Nature*, vol. 518, no. 7540, p. 529, 2015.
- [2] D. Silver, J. Schrittwieser, K. Simonyan, I. Antonoglou, A. Huang, A. Guez, T. Hubert, L. Baker, M. Lai, A. Bolton *et al.*, “Mastering the game of go without human knowledge,” *Nature*, vol. 550, no. 7676, p. 354, 2017.
- [3] D. Zhao and Y. Zhu, “Meca near-optimal online reinforcement learning algorithm for continuous deterministic systems,” *IEEE transactions on neural networks and learning systems*, vol. 26, no. 2, pp. 346–356, 2014.
- [4] B. Kiumarsi, K. G. Vamvoudakis, H. Modares, and F. L. Lewis, “Optimal and autonomous control using reinforcement learning: A survey,” *IEEE transactions on neural networks and learning systems*, vol. 29, no. 6, pp. 2042–2062, 2017.
- [5] Z. Yang, K. Merrick, L. Jin, and H. A. Abbass, “Hierarchical deep reinforcement learning for continuous action control,” *IEEE transactions on neural networks and learning systems*, vol. 29, no. 11, pp. 5174–5184, 2018.
- [6] D. Ye, M. Zhang, and Y. Yang, “A multi-agent framework for packet routing in wireless sensor networks,” *sensors*, vol. 15, no. 5, pp. 10026–10047, 2015.

⁷See Table III for more results on 3m and 8m.

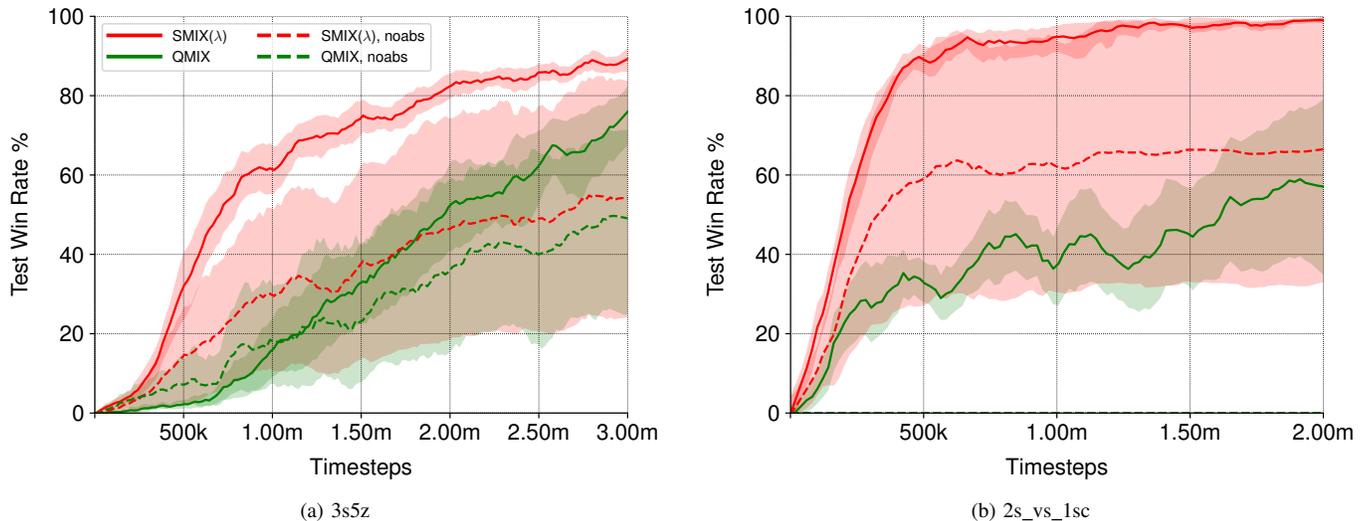


Fig. 8. The effects of non-negative constraints of SMIX(λ) and QMIX in two different scenarios. Performances of those without non-negative constraints (denoted by “noabs”) are shown as dashed lines. Methods plotted with solid lines are enforced with non-negative constraint. The mean and 95% confidence interval is shown across 6 independent runs.

- [7] E. Van der Pol and F. A. Oliehoek, “Coordinated deep reinforcement learners for traffic light control,” *Proceedings of the Learning, Inference and Control of Multi-Agent Systems*, 2016.
- [8] W. Liu and J. Huang, “Cooperative adaptive output regulation for lower triangular nonlinear multi-agent systems subject to jointly connected switching networks,” *IEEE transactions on neural networks and learning systems*, vol. 31, no. 5, pp. 1724–1734, 2019.
- [9] C. Yu, M. Zhang, F. Ren, and G. Tan, “Emotional multiagent reinforcement learning in spatial social dilemmas,” *IEEE transactions on neural networks and learning systems*, vol. 26, no. 12, pp. 3083–3096, 2015.
- [10] G. Jing, Y. Zheng, and L. Wang, “Consensus of multiagent systems with distance-dependent communication networks,” *IEEE transactions on neural networks and learning systems*, vol. 28, no. 11, pp. 2712–2726, 2016.
- [11] X. Bu, Z. Hou, and H. Zhang, “Data-driven multiagent systems consensus tracking using model free adaptive control,” *IEEE transactions on neural networks and learning systems*, vol. 29, no. 5, pp. 1514–1524, 2017.
- [12] Y. Zheng, J. Ma, and L. Wang, “Consensus of hybrid multi-agent systems,” *IEEE transactions on neural networks and learning systems*, vol. 29, no. 4, pp. 1359–1365, 2017.
- [13] M. Jaderberg, W. M. Czarnecki, I. Dunning, L. Marris, G. Lever, A. G. Castaneda, C. Beattie, N. C. Rabinowitz, A. S. Morcos, A. Ruderman *et al.*, “Human-level performance in 3d multiplayer games with population-based reinforcement learning,” *Science*, vol. 364, no. 6443, pp. 859–865, 2019.
- [14] J. N. Foerster, G. Farquhar, T. Afouras, N. Nardelli, and S. Whiteson, “Counterfactual multi-agent policy gradients,” in *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [15] T. Rashid, M. Samvelyan, C. S. Witt, G. Farquhar, J. Foerster, and S. Whiteson, “Qmix: Monotonic value function factorisation for deep multi-agent reinforcement learning,” in *International Conference on Machine Learning*, 2018, pp. 4292–4301.
- [16] G. J. Laurent, L. Matignon, L. Fort-Piat *et al.*, “The world of independent learners is not markovian,” *International Journal of Knowledge-based and Intelligent Engineering Systems*, vol. 15, no. 1, pp. 55–64, 2011.
- [17] R. Lowe, Y. Wu, A. Tamar, J. Harb, O. P. Abbeel, and I. Mordatch, “Multi-agent actor-critic for mixed cooperative-competitive environments,” in *Advances in Neural Information Processing Systems*, 2017, pp. 6379–6390.
- [18] F. A. Oliehoek, M. T. Spaan, and N. Vlassis, “Optimal and approximate q-value functions for decentralized pomdps,” *Journal of Artificial Intelligence Research*, vol. 32, pp. 289–353, 2008.
- [19] L. Kraemer and B. Banerjee, “Multi-agent reinforcement learning as a rehearsal for decentralized planning,” *Neurocomputing*, vol. 190, pp. 82–94, 2016.
- [20] P. Sunehag, G. Lever, A. Gruslys, W. M. Czarnecki, V. Zambaldi, M. Jaderberg, M. Lanctot, N. Sonnerat, J. Z. Leibo, K. Tuyls *et al.*, “Value-decomposition networks for cooperative multi-agent learning based on team reward,” in *International Conference on Autonomous Agents and MultiAgent Systems*, 2018, pp. 2085–2087.
- [21] K. Son, D. Kim, W. J. Kang, D. E. Hostallero, and Y. Yi, “Qtran: Learning to factorize with transformation for cooperative multi-agent reinforcement learning,” in *International Conference on Machine Learning*, 2019, pp. 5887–5896.
- [22] Bellman, *Dynamic Programming*, ser. Rand Corporation research study. Princeton University Press, 1957.
- [23] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*. MIT press, 2018.
- [24] A. Harutyunyan, M. G. Bellemare, T. Stepleton, and R. Munos, “Q (λ) with off-policy corrections,” in *Proceedings of the International Conference on Algorithmic Learning Theory*. Springer, 2016, pp. 305–320.
- [25] M. Samvelyan, T. Rashid, C. Schroeder de Witt, G. Farquhar, N. Nardelli, T. G. Rudner, C.-M. Hung, P. H. Torr, J. Foerster, and S. Whiteson, “The starcraft multi-agent challenge,” in *International Conference on Autonomous Agents and MultiAgent Systems*, 2019, pp. 2186–2188.
- [26] C. Wen, X. Yao, Y. Wang, and X. Tan, “Smix(λ): Enhancing centralized value functions for cooperative multi-agent reinforcement learning,” in *Proceedings of the Thirty-Fourth AAAI Conference on Artificial Intelligence*, 2020.
- [27] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, “Proximal policy optimization algorithms,” *arXiv preprint arXiv:1707.06347*, 2017.
- [28] T. T. Nguyen, N. D. Nguyen, and S. Nahavandi, “Deep reinforcement learning for multiagent systems: A review of challenges, solutions, and applications,” *IEEE transactions on cybernetics*, 2020.
- [29] P. Hernandez-Leal, B. Kartal, and M. E. Taylor, “A survey and critique of multiagent deep reinforcement learning,” *Autonomous Agents and Multi-Agent Systems*, vol. 33, no. 6, pp. 750–797, 2019.
- [30] T. Schaul, J. Quan, I. Antonoglou, and D. Silver, “Prioritized experience replay,” *arXiv preprint arXiv:1511.05952*, 2015.
- [31] R. Munos, T. Stepleton, A. Harutyunyan, and M. Bellemare, “Safe and efficient off-policy reinforcement learning,” in *Advances in Neural Information Processing Systems*, 2016, pp. 1054–1062.
- [32] D. Precup, R. S. Sutton, and S. Dasgupta, “Off-policy temporal-difference learning with function approximation,” in *International Conference on Machine Learning*, 2001, pp. 417–424.
- [33] J. Foerster, N. Nardelli, G. Farquhar, T. Afouras, P. H. Torr, P. Kohli, and S. Whiteson, “Stabilising experience replay for deep multi-agent reinforcement learning,” in *International Conference on Machine Learning*, vol. 70. JMLR.org, 2017, pp. 1146–1155.
- [34] K. De Asis, J. F. Hernandez-Garcia, G. Z. Holland, and R. S. Sutton, “Multi-step reinforcement learning: A unifying algorithm,” in *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.

- [35] J. F. Hernandez-Garcia and R. S. Sutton, "Understanding multi-step deep reinforcement learning: A systematic study of the dqn target," *arXiv preprint arXiv:1901.07510*, 2019.
- [36] S. Fujimoto, D. Meger, and D. Precup, "Off-policy deep reinforcement learning without exploration," in *International Conference on Machine Learning*, 2019, pp. 2052–2062.
- [37] J. N. Foerster, Y. M. Assael, N. de Freitas, and S. Whiteson, "Learning to communicate to solve riddles with deep distributed recurrent q-networks," *arXiv preprint arXiv:1602.02672*, 2016.
- [38] A. Singh, T. Jain, and S. Sukhbaatar, "Learning when to communicate at scale in multiagent cooperative and competitive tasks," in *International Conference on Learning Representations*, 2019.
- [39] J. Jiang and Z. Lu, "Learning attentional communication for multi-agent cooperation," in *Advances in Neural Information Processing Systems*, 2018, pp. 7254–7264.
- [40] S. Iqbal and F. Sha, "Actor-attention-critic for multi-agent reinforcement learning," in *International Conference on Machine Learning*, 2019, pp. 2961–2970.
- [41] J. Foerster, R. Y. Chen, M. Al-Shedivat, S. Whiteson, P. Abbeel, and I. Mordatch, "Learning with opponent-learning awareness," in *Proceedings of the 17th International Conference on Autonomous Agents and MultiAgent Systems*, 2018, pp. 122–130.
- [42] A. Letcher, J. Foerster, D. Balduzzi, T. Rocktäschel, and S. Whiteson, "Stable opponent shaping in differentiable games," in *International Conference on Machine Learning*, 2019.
- [43] S. Sukhbaatar, Z. Lin, I. Kostrikov, G. Synnaeve, A. Szlam, and R. Fergus, "Intrinsic motivation and automatic curricula via asymmetric self-play," in *International Conference on Learning Representations*, 2018.
- [44] C. Sun, W. Liu, and L. Dong, "Reinforcement learning with task decomposition for cooperative multiagent systems," *IEEE transactions on neural networks and learning systems*, 2020.
- [45] F. A. Oliehoek and Amato, *A concise introduction to decentralized POMDPs*. Springer, 2016, vol. 1.
- [46] C. Boutilier, "Planning, learning and coordination in multiagent decision processes," in *Proceedings of the 6th conference on Theoretical aspects of rationality and knowledge*, 1996, pp. 195–210.
- [47] J. Hu and M. P. Wellman, "Nash q-learning for general-sum stochastic games," *Journal of machine learning research*, vol. 4, no. Nov, pp. 1039–1069, 2003.
- [48] G. Papoudakis, F. Christianos, A. Rahman, and S. V. Albrecht, "Dealing with non-stationarity in multi-agent deep reinforcement learning," *arXiv preprint arXiv:1906.04737*, 2019.
- [49] G. Tesauro, "Extending q-learning to general adaptive multi-agent systems," in *Advances in neural information processing systems*, 2004, pp. 871–878.
- [50] M. Lauer and M. Riedmiller, "An algorithm for distributed reinforcement learning in cooperative multi-agent systems," in *In Proceedings of the Seventeenth International Conference on Machine Learning*. Citeseer, 2000.
- [51] R. Sutton, A. R. Mahmood, D. Precup, and H. Hasselt, "A new q (λ) with interim forward view and monte carlo equivalence," in *International Conference on Machine Learning*, 2014, pp. 568–576.
- [52] H. Van Seijen, H. Van Hasselt, S. Whiteson, and M. Wiering, "A theoretical and empirical analysis of expected sarsa," in *2009 IEEE Symposium on Adaptive Dynamic Programming and Reinforcement Learning*. IEEE, 2009, pp. 177–184.
- [53] D. Precup, R. S. Sutton, and S. Singh, "Eligibility traces for off-policy policy evaluation," in *International Conference on Machine Learning*, 2000.
- [54] Q. Liu, L. Li, Z. Tang, and D. Zhou, "Breaking the curse of horizon: Infinite-horizon off-policy estimation," in *Advances in Neural Information Processing Systems*, 2018, pp. 5356–5366.
- [55] J. Chung, C. Gulcehre, K. Cho, and Y. Bengio, "Empirical evaluation of gated recurrent neural networks on sequence modeling," in *Advances in Neural Information Processing Systems*, 2014.
- [56] D. Ha, A. Dai, and Q. V. Le, "Hypernetworks," in *International Conference on Learning Representations*, 2017.
- [57] M. Tan, "Multi-agent reinforcement learning: Independent vs. cooperative agents," in *International Conference on Machine Learning*, 1993, pp. 330–337.
- [58] S. Shalev-Shwartz and S. Ben-David, *Understanding machine learning: From theory to algorithms*. Cambridge university press, 2014.
- [59] S. J. Russell and P. Norvig, *Artificial intelligence: a modern approach*. Malaysia; Pearson Education Limited., 2016.
- [60] A. Mahajan, T. Rashid, M. Samvelyan, and S. Whiteson, "Maven: Multi-agent variational exploration," in *Advances in Neural Information Processing Systems*, 2019, pp. 7613–7624.
- [61] W. Fedus, P. Ramachandran, R. Agarwal, Y. Bengio, H. Larochelle, M. Rowland, and W. Dabney, "Revisiting fundamentals of experience replay," in *International Conference on Machine Learning*, 2020.