

# SCHOLASTIC-ACTOR-CRITIC FOR MULTI AGENT RE-INFORCEMENT LEARNING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

1 The Actor-Critic framework of multi-agent reinforcement learning(MARL) is gathering  
 2 more attention nowadays. Centralized training with decentralized execution  
 3 allows the policies to use extra information to ease the training while enhancing  
 4 overall performance. In such a framework, the quality of critic profoundly impacts  
 5 the final average rewards. Thus we present a method, called Scholastic-Actor-  
 6 Critic(SMAC), that involves a more powerful critic to maintain efficiency in ample  
 7 knowledge acquisition. The headmaster critic is designed to group agents with  
 8 proper size and proper timing, while other critics update simultaneously at the  
 9 decision time. The learning rule includes additional terms account for the impact  
 10 of other agents within a group. Our method receives higher payouts compared to  
 11 other state-of-the-art methods and is robust against the explosion of dimension  
 12 during training. We apply our method to the Coin Game, the Cooperative Treasure  
 13 Collection(CTC) (Lerer & Peysakhovich, 2017) and a dynamic battle game,  
 14 MAgent(Zheng et al., 2018). Experiment results are all satisfying.

## 15 1 INTRODUCTION

16 MARL(Multi-Agent Reinforcement Learning) is gathering more attention in deep learning researches.  
 17 Artificial agents thus perform better to interact both with other agents and humans in complex partially  
 18 competitive or sequential dilemma occasions. MARL is a big topic with fully cooperative settings,  
 19 competitive settings and mixed settings. It is still challenging to make decisions with inadequate  
 20 information in applications, such as playing games, advertising and self-driving cars.

21 The ability to maintain cooperation and competition in a variety of complicated situations is essential  
 22 in MARL. Early works focus on improving policy or value constructing methods (Foerster et al.,  
 23 2018b) (Silver et al., 2016) (Sukhbaatar et al., 2017)(Gupta et al., 2017), promoting more effectively  
 24 opponent modeling methods (He et al., 2016)(Foerster et al., 2018a)(Metz et al., 2016)(Tesauro, 2004)  
 25 and enhancing communication between opponents (Foerster et al., 2017) (Lerer & Peysakhovich,  
 26 2017) (Das et al., 2017) (Foerster et al., 2016) (Mordatch & Abbeel, 2018) (Sukhbaatar et al., 2016)  
 27 (Lauer & Riedmiller, 2000) (Matignon et al., 2007) (Omidshafiei et al., 2017).

28 In cooperative-and-competitive settings, Iterated Prisoners' Dilemma is a traditional problem, in  
 29 which selfish actions usually lead to an overall bad result. At this time, cooperation maximizes social  
 30 welfare, which leads to an average best outcome. In this setting, the measurement is the total of  
 31 rewards of all agents, while randomly initialized agents usually pursue independent gradient descent  
 32 on the specific value function. Lerer & Peysakhovich (2017) and Leibo et al. (2017) point out that  
 33 reciprocity among agents results in a higher average reward. Peng et al. (2017) and Evans & Gao  
 34 (2016) find that even in strongly adversarial settings, reciprocity shows its nontrivial value.

35 In traditional Q methods, each agent's policy changes over time, resulting in a non-stationary  
 36 environment. In a non-stationary environment, agents are not able to make good use of naive  
 37 experience replay. Recent years Lowe et al. (2017) propose the actor-critic framework(also called  
 38 MADDPG), which combines offline and online learning, which enhances the ability for multi-agent  
 39 learning. Then, (Yang et al., 2018)(MF-MARL), Iqbal & Sha (2018)(MAAC) and Jiang & Lu (2018)  
 40 explore policy and communication optimizations within the Actor-Critic framework.

41 We here propose the Scholastic-Multi-Actor-Critic method(SMAC), which aims to improve the ability  
 42 of the critic. We want to train a more powerful critic, the headmaster critic that enables actors to

43 communicate more efficiently during training. The SMAC learns to control when and how an agent  
 44 receives information from others. That is, the access of observations of an agent depends on the  
 45 critic. This optional additional term when applied to a group of agents, leads to extra reciprocity and  
 46 cooperation. The policy gradient is consistent with prior works presented by Sandholm and Crites  
 47 Sandholm & Crites (1996) and Foerster et al. (2018a).

48 Our approach enables high dimensional settings. We deploy experiments on the Coin Game4.1.1, the  
 49 Cooperative Treasure Collection4.1.2 and the MAgent(Zheng et al., 2018). Our algorithm leads to  
 50 the overall highest average return on these games. All agents using our method achieve the stable  
 51 equilibrium with less training resources.

## 52 2 RELATED WORKS

53 As mentioned above, interactions between agents can either be cooperative, competitive or usually  
 54 both. Model-free reinforcement learning algorithms in this domain could be concluded to value-based  
 55 methods, policy-based methods and actor-critic methods.

56 MADDPG (Lowe et al., 2017) combines offline and online learning that enhances the ability of  
 57 multi-agent learning. It allows the policies to use extra information to ease the training. The critic is  
 58 enlarged with extra information about the policies of other agents, while each actor only has access  
 59 to local information. Local actors are used at the execution phase after training.

60 COMA(Counterfactual Multi-Agent Policy Gradients) raised by Foerster et al. (2018b) is aimed  
 61 to solve multi-agent credit assignment in cooperative settings. Before, each agent trains with his  
 62 own critic so that the information sharing between them is insufficient, resulting in poor cooperation  
 63 between agents. Therefore, the centralized critic firstly introduced in COMA to give a preliminary  
 64 solution to this problem.

65 MF-MARL, the Mean Field Multi-Agent Reinforcement method developed by Yang et al. (2018) try  
 66 to model opponents by the use of Mean Field Theory under Q-learning and Actor-Critic methods. It  
 67 uses numerical techniques that greatly reduce the cost of modeling opponents.

68 Somewhat like COMA(Foerster et al., 2018b), MAAC (Iqbal & Sha, 2018)(Multi-Actor-Attention-  
 69 Critic) considers to make full use of information and takes the attention mechanism within the  
 70 centralized critic network. The experiment result shows that as the scale is growing, this method  
 71 demonstrates its great effect. However, the requirement of computing is too high. On the other hand,  
 72 ATOC (Jiang & Lu, 2018)(Learning Attentional Communication) decides to find a good communica-  
 73 tion group for the initiator agents by attention methods, too. Nevertheless, the determination of the  
 74 initiator is very vague, and as the decisive role, if the initial selection is not appropriate, the entire  
 75 model will collapse.

## 76 3 METHODS

### 77 3.1 BACKGROUND

#### 78 3.1.1 STOCHASTIC GAME AND DEEP Q-NETWORKS

79 A multi-agent stochastic game  $G$  is formulated by a tuple  $G = \langle S, A, P, O, R, n, \gamma \rangle$ .  $S$  denotes the  
 80 state space, the configurations for all agents. Each agent takes  $a_i \in A$  at every time step, forming  
 81 joint actions  $\mathbf{a} \in \mathbf{A} \equiv A^n$ . To choose actions, each agent uses a policy  $\pi_{\theta_i} : \mathcal{O}_i \times \mathcal{A}_i$ , which  
 82 produces the next state according to the state transition function.  $P(s'|s, \mathbf{a}) : S \times \mathbf{A} \times S \rightarrow [0, 1]$   
 83 denotes transition probabilities of states, and  $o_i \in \mathbf{O}$  denotes observations. The reward function  
 84  $r^i(s, \mathbf{a}) : S \times \mathbf{A} \rightarrow \mathbb{R}$  specify rewards and  $\gamma \in [0, 1)$  is the discount factor, and for each agent,  
 85  $R_t^i = \sum_{l=0}^{\infty} \gamma^l r_{t+l}^i$ . Policy gradient methods update an agent’s policy, parameterised by  $\theta^i$ .

86 Provided and initial state  $s$ , the value function of agent  $i$  under the joint policy  $\pi$  could be formulated  
 87 as:

$$v_{\pi}^j(s) = v^j(s; \pi) = \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{\pi, p} \left[ r_t^j | s_0 = s, \pi \right] \quad (1)$$

88 We define the Q-function within the framework of N-agent games based on the Bellman equation in  
89 (1) such that the Q-function  $Q_\pi^i$  for agent  $i$  under policy  $\pi$  could be recursively formulated as

$$Q^\pi(s, a) = \mathbb{E}_{s'} [r(s, a) + \gamma \mathbb{E}_{a' \sim \pi} [Q^\pi(s', a')]] \quad (2)$$

90 , and deep Q-networks learn the action-value function  $Q^*$  by minimizing the loss in (3):

$$\mathcal{L}(\theta) = \mathbb{E}_{s, a, r, s'} [(Q^*(s, a | \theta) - y)^2], \quad (3)$$

91 and

$$y = r + \gamma \max_{a'} \bar{Q}^*(s', a') \quad (4)$$

92 where  $\bar{Q}^*$  is the target Q function and its parameters update periodically with the most recent  
93  $\theta$ , which stabilize the learning. Besides, the experience replay buffer  $D = (s, a, r, s')$  also used  
94 to stabilization. However, because agents are independently updating their policies as learning  
95 progresses, the environment appears non-stationary from the view of any one agent, violating Markov  
96 assumptions required for convergence of Q-learning. Foerster et al. (2017)'s approach point out,  
97 another difficulty is that the experience replay buffer cannot be used in such a setting since in general.

### 98 3.1.2 POLICY GRADIENTS

99 Policy gradient techniques (Sutton et al., 2000) aims to estimate the gradient of an agent's expected  
100 returns with respect to the parameters of its policy. This gradient estimate takes the following form as  
101 (5):

$$\nabla_\theta J(\pi_\theta) = \mathbb{E}_{a \sim \pi_\theta} \left[ \nabla_\theta \log(\pi_\theta(a_t | s_t)) \sum_{t'=t}^{\infty} \gamma^{t'-t} r_{t'}(s_{t'}, a_{t'}) \right] \quad (5)$$

### 102 3.1.3 ACTOR-CRITIC METHODS

103 The term  $\sum_{t'=t}^{\infty} \gamma^{t'-t} r_{t'}(s_{t'}, a_{t'})$  in the policy gradient estimator leads to high variance, as returns  
104 can vary drastically between training episodes. The Actor-critic method (Konda & Tsitsiklis, 2000)  
105 aims to ameliorate this issue by using a function to approximate the expected returns. Moreover, it  
106 replaces the original return term in the policy gradient estimator with this function. Given a state and  
107 action, an agent under actor-critic methods learns a function to estimate expected discounted returns  
108 as:  $Q_\psi(s_t, a_t) = \mathbb{E} \left[ \sum_{t'=t}^{\infty} \gamma^{t'-t} r_{t'}(s_{t'}, a_{t'}) \right]$ , it updates by minimizing the regression loss of:

$$\mathcal{L}_Q(\psi) = \mathbb{E}_{s, a, r, s'} [(Q_\psi(s, a) - y)^2] \quad (6)$$

109 where

$$y = r(s, a) + \gamma \mathbb{E}_{a' \sim \pi(s')} [Q_{\bar{\psi}}(s', a')] \quad (7)$$

110 in which  $Q_{\bar{\psi}}$  is the target Q-value function. A recent approach (Haarnoja et al., 2018) applies a  
111 soft value function by modifying the policy gradient to incorporate an entropy term to encourage  
112 exploration and avoid converging to non-optimal deterministic policies. It could be formulated as:

$$\nabla_\theta J(\pi_\theta) = \mathbb{E}_{a \sim \pi_\theta} [\nabla_\theta \log(\pi_\theta(a|s)) (\alpha \log(\pi_\theta(a|s)) - Q_\psi(s, a) + b(s))] \quad (8)$$

113 where  $b(s)$  is a state-dependent baseline. The loss function for temporal difference learning is also  
114 revised with a new target, that is:

$$y = r(s, a) + \gamma \mathbb{E}_{a' \sim \pi(s')} [Q_{\bar{\psi}}(s', a') - \alpha \log(\pi_\theta(a'|s'))] \quad (9)$$

### 115 3.2 SCHOLASTIC-ACTOR-CRITIC

116 Our method obeys the same paradigm of training critics centrally and executing learned policies  
117 distributedly. That is proposed to overcome the challenge of non-stationary environments. The main  
118 idea behind our approach is group discussion, which encourages agents to emulate those better than  
119 themselves with high efficiency. We design a more powerful critic, the headmaster critic, to learn  
120 how to group agents and determine when to communicate, that has the same effect of the attention  
121 mechanism. The additional critic has a global perspective of all agents and focuses on agents with  
122 highest and lowest rewards. Accounting for the impacts from opponents, observations and actions  
123 incorporate information into the estimation of each agent's value function in the same group.

## 124 3.2.1 ASSIGNMENT OF GROUPS

125 Expand the setting of MAAC(Iqbal & Sha, 2018), we introduce a headmaster critic to assign  
 126 communication groups. The critic randomly selects  $n$  collections with random size  $s$  and changes  
 127 every  $k$  epochs. After selecting  $n$  collections, we take the average contributions from each group (super  
 128 agent,  $sa$ ), and apply the following loss function:

$$\mathcal{L}_Q(\psi) = \sum_{i=1}^N \mathbb{E}_{(o,sa,r,o') \sim D} \left[ \left( Q_i^\psi(o, sa) - y_i \right)^2 \right] \quad (10)$$

129 where

$$y_i = r_i + \gamma \mathbb{E}_{sa' \sim \pi_{\bar{\theta}}(o')} \left[ Q_i^{\bar{\psi}}(o', sa') \right] \quad (11)$$

130 The action-value  $Q_i^\psi(o, sa)$  function estimates outcomes in group  $i$  from 1 to  $n$ , which receives  
 131 observations and actions of agents. To avoid the degradation, we set threshold for  $n$  as  $n/2$  and  $s$  as  
 132  $s > 1$ .

## 133 3.2.2 CRITICS IN GROUPS

134 Critics within the same group updated together to minimize a joint regression loss function:

$$\mathcal{L}_Q(\psi) = \sum_{i=1}^N \mathbb{E}_{(o,a,r,o') \sim D} \left[ \left( Q_i^\psi(o, a) - y_i \right)^2 \right] \quad (12)$$

135 Note that  $Q_i^\psi(o, a)$ , the action-value estimate for agent  $i$ , receives observations and actions for partial  
 136 agents. Where,

$$y_i = r_i + \gamma \mathbb{E}_{a' \sim \pi_{\bar{\theta}}(o')} \left[ Q_i^{\bar{\psi}}(o', a') - \alpha \log(\pi_{\bar{\theta}_i}(a'_i | o'_i)) + \Gamma \right] \quad (13)$$

137

$$\Gamma = \omega \log(\pi_{\bar{\theta}_i}(a'_i | o'_{others})) + \sigma \log(\pi_{\bar{\theta}_i}(a'_i | o'_{others})) \quad (14)$$

138 in which  $\psi$  and  $\theta$  are the parameters of the target critics and target policies, respectively.

## 139 3.2.3 AGENTS IN GROUPS

140 To calculate the Q-value function  $Q_i^\psi(o, a)$  for agent  $i$ , the critic receives the observations  $o =$   
 141  $(o_1, \dots, o_N)$  and actions  $a = (a_1, \dots, a_N)$  for all agents in a group. Then other agents' contributions  
 142 could be formulated as 15. where  $g_i$  is a two-layer MLP(multi-layer perceptron) embedding function  
 143 and  $f_i$  is a softmax function. It could be formulated as:

$$Q_i^\psi(o, a) = f_i(g_i(o_i, a_i)) \quad (15)$$

144 As shown in Foerster et al. (2018b), an advantage function using a baseline that only marginalizes out  
 145 the actions of the given agent from  $Q$ . It helps in credit assigning. In other words, by comparing the  
 146 value of specific actions to an average action, an agent could learn whether the action he made would  
 147 cause an increase in expected return. Thus the individual policies are updated with the following  
 148 gradient:

$$\nabla_{\theta_i} J(\pi_{\theta}) = \mathbb{E}_{a \sim \pi_{\theta}} \left[ \nabla_{\theta_i} \log(\pi_{\theta_i}(a_i | o_i)) \left( \alpha \log(\pi_{\theta_i}(a_i | o_i)) - Q_i^\psi(o, a) + b(o, a_{others}) \right) \right] \quad (16)$$

$$A_i(o, a) = Q_i^\psi(o, a) - b(o, a_{others}) \quad (17)$$

$$b(o, a_{others}) = \mathbb{E}_{a_i \sim \pi_i(o_i)} \left[ Q_i^\psi(o, (a_i, a_{others})) \right]$$

149  $b(o, a)$  is the multi-agent baseline that used to calculate the advantage function.

150 We implement a more general and flexible form of a multi-agent baseline. We do not apply a global  
151 reward, but naturally decompose an agent’s encoding observations and the average of encodings of  
152 other agents.

$$\mathbb{E}_{a_i \sim \pi_i(o_i)} \left[ Q_i^\psi(o, (a_i, a_{others})) \right] = \sum_{a'_i \in A_i} \pi(a'_i | o_i) Q_i(o, (a'_i, a_{others})) \quad (18)$$

153 As shown above, we output the value for every action and add an observation-encoder as  $E_i = g_i(o_i)$ .  
154 For each agent, using these encodings in place of the  $E_i = g_i(o_i, a_i)$  described above, and modify  $f_i$   
155 such that it outputs a value for each possible action. We can estimate the expectation by sampling  
156 actions from our policy and averaging their Q-values. So we do not need to add any parameters in the  
157 case of continuous policies.

## 158 4 EXPERIMENTS

### 159 4.1 SETUP

160 We operate our algorithms in various settings, including the Coin Game 4.1.1, Cooperative Treasure  
161 Collection(CTC) (Lerer & Peysakhovich, 2017) 4.1.2 and MAgent(Zheng et al., 2018) (a cooperative-  
162 competitive battle game in the Open-source MAgent system) that tests capabilities of our approach  
163 and baselines. The three games we raised, from simple to complex, are all facing iterated prisoners  
164 dilemmas(Luce & Raiffa, 1958). For each setting, we study the scalability of different methods as the  
165 number of agents grows and evaluate their ability to attend to information relevant to rewards.

#### 166 4.1.1 COIN GAME

167 The Coin Game is a higher dimensional alternative of IPD (iterated prisoners dilemma), which  
168 is convenient to make comparisons to previous works. As shown in 1, two agents with red and  
169 blue colors are tasked to collect coins which are either red or blue on the grids. A new coin with  
170 random color appears randomly after the last one is picked up. Agents move to a coin’s position and  
171 both receive a point after picking it up while the agent with a different color loses 2 points. When  
172 they only pick up coins with their own color, the total return is maximized. While players usually  
173 pick up different ones. Therefore the maximum achievable collective return is approximately 50 in  
174 expectation if neither agent chooses to defect and both agents collect all coins of their own color. In  
175 this game we define niceness as  $n(s_t, a_t)$  to be part of the measurement. If an agent takes action  $a_t^i$ ,  
176 picks up a coin which penalizes the other players,  $n(s_t, a_t) = -1$ . We use recent defections as the  
177 measure of niceness  $N(T) = \sum_{i=1}^t \lambda^{t-i} n(s_i, a_i)$  at time  $T$ .

#### 178 4.1.2 COOPERATIVE TREASURE COLLECTION

179 Cooperative Treasure Collection(CTC), as shown in 1, is a variant of Coin Game in which agents  
180 play roles as hunter or bank. ”Hunter”s are tasked to collect the treasure of any color and deposit  
181 them into the corresponding colored bank. The ”Bank”s are tasked to gather as much treasure as  
182 possible from the ”Hunter”s simply. Agents could see each others’ positions and concern their own.  
183 ”Hunter”s receive a global reward for the successful collection of treasure, and all agents receive a  
184 global reward of the depositing amount. ”Hunter”s will additionally penalized for colliding with each  
185 other. As such, the task contains a mixture of shared and individual rewards. It requires different  
186 ”modes of attention” which depends on the agent’s state and other agents’ potential actions that affects  
187 its rewards.

#### 188 4.1.3 MAGENT

189 The mixed cooperative-competitive battle game, MAgent(Zheng et al., 2018), is a more complex  
190 multi-player environment. Agents are divided into armies, and required to take a series of actions  
191 while exact discounted reward cannot be assessed. Each army consists of homogeneous agents, and  
192 the goal of them is to get more rewards by collaborating with teammates to defeat all opponents.

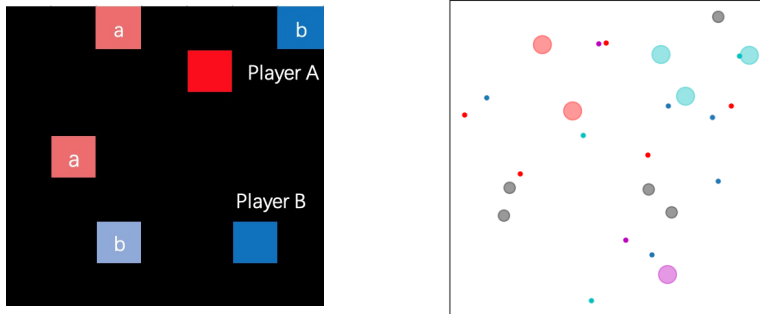


Figure 1: The Coin Game and the Cooperative Treasure Collection Game

Table 1: The average rewards compared to other methods with growing of the scale in the convergent training stages.

Game	Agents	MADDPG+SAC	MARL	MAAC	ATOC	Ours(SMAC)
CTC	8	-3.9	3.4	-4.7	3.1	2.8
	16	17.6	11.7	0.8	1.5	3.4
	32	32.1	14.8	10.1	13.0	13.2
	64	41.2	18.9	23.3	24.2	24.5
	128	77.3	29.5	64.1	65.8	78.1
MAgent	8	-	3.4	4.9	-2.7	0.8
	16	-	14.7	27.9	26.5	27.0
	32	-	32.5	29.5	28.6	30.7
	64	-	34.8	35.4	39.1	41.5
	128	-	35.6	56.1	40.6	57.7

*\*Note that the number of agents for each group in MAgent is half of the total. And all values are normalized into 0 to 100.*

193 Agents can take actions to either move to or attack others on nearby grids. Ideally, the agents are able  
 194 to learn skills such as chasing to hunt, escaping from enemies or working with teammates.

## 195 4.2 BASELINES

196 We have compared our method to recently proposed state-of-art methods in the multi-agent learning  
 197 field: (1)DDPG(Lillicrap et al., 2015), (2)MADDPG(Lowe et al., 2017), (3)MF-MARL(Yang et al.,  
 198 2018), (4)MAAC(Iqbal & Sha, 2018) (5)ATOC(Jiang & Lu, 2018).

199 As mentioned in MAAC(Iqbal & Sha, 2018), we do some modifications on some algorithms for exper-  
 200 iments. Since deterministic policies are not possible, we use the Gumbel-Softmax reparametrization  
 201 trick for learning in discrete action spaces for both MADDPG(Lowe et al., 2017) and DDPG(Lillicrap  
 202 et al., 2015). The modified versions are referred to as MADDPG (Discrete) and DDPG (Discrete).  
 203 For a detailed description of this reparametrization, we use a soft actor-critic method (Haarnoja et al.,  
 204 2018) to optimize. We implement MADDPG with Soft Actor-Critic, named as MADDPG+SAC. Then  
 205 the baselines are (1)DDPG (Discrete) (2)MADDPG (Discrete) (3)MADDPG+SAC (4)MF-MARL  
 206 (5)ATOC.

207 Hyperparameters are tuned based on performance and kept constant across all variants of critic  
 208 architectures. All methods are re-implemented such that their approximate total number of parameters  
 209 (across agents) is close to our approach. These models are trained with eight random seeds each.

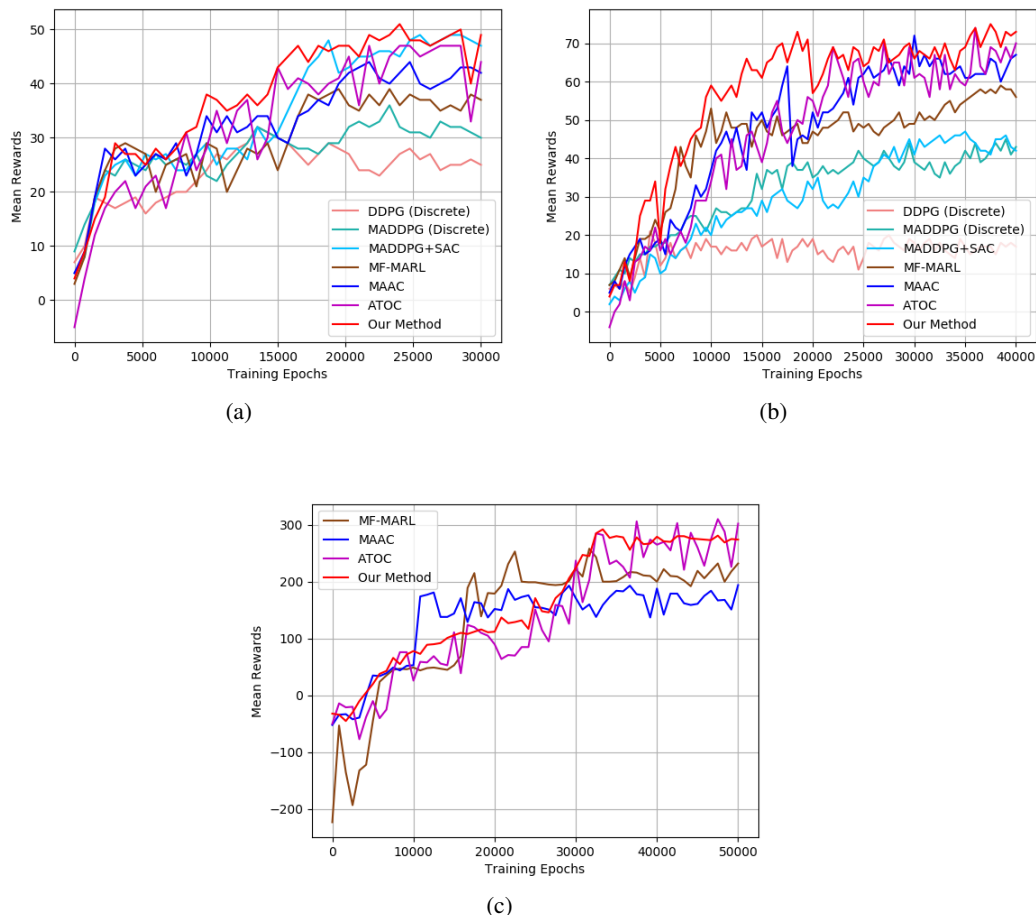


Figure 2: Results of our methods and others. In Coin Game(a) and CTC(b), methods involve DDPG(Discrete)4.2, MADDPG(Discrete)4.2, MAAC(Iqbal & Sha (2018)), MF-MARL(Iqbal & Sha (2018)) and ATOC(Jiang & Lu, 2018). In MAgent(c)(Zheng et al., 2018), we compare our method to MF-MARLYang et al. (2018), MAACIqbal & Sha (2018) and ATOC(Jiang & Lu, 2018).

### 210 4.3 RESULTS AND DISCUSSION

211 We first compare the average rewards attained by all approaches. We normalized by the range  
 212 of awards achieved in an environment, as the number of agents changes. The proposed approach  
 213 (SMAC) is competitive with other state-of-the-art approaches as shown in 4.3. In the Coin Game,  
 214 most algorithms show a pleasing result while the MARL method shows less poorly performance.  
 215 MAAC is competitive with our approach in both the Coin Game and the CTC environment. On the  
 216 other hand, DDPG(Discrete), MADDPG (Discrete), MADDPG+SAC and MARL don't perform well  
 217 on CTC. We infer that due to the simplicity of action modes and the limited scale of agents, it's not  
 218 hard for agents to learn tricks. Moreover, each agent's local observation provides enough information  
 219 to make a decent prediction of its expected rewards.

220 However, agents within MAgent(Zheng et al., 2018) dynamics over time so that it's not capable for  
 221 DDPG(Discrete), MADDPG (Discrete), MADDPG+SAC break down. Thus we compare our method  
 222 to MF-MARL(mean field-MARL, Yang et al. (2018)), MAAC(Iqbal & Sha (2018)) and ATOC(Jiang  
 223 & Lu, 2018). For all methods, rewards firstly are under zero, but along with the process of training,  
 224 the reward gradually grows and finally stop in different levels. In this game, subgroups of agents  
 225 are interacting and performing coordinated tasks with separate rewards while the components are

226 changing over time. Thus it exemplifies why dynamic attention can be beneficial. MAAC(Iqbal &  
227 Sha (2018) and ATOC(Jiang & Lu, 2018) take more iterations to reach a stationary state.

228 Further, we explore the improvements with growing scale as shown in Table 1. DDPG(Discrete)  
229 and MADDPG(Discrete) could not handle a high dimensional learning. MADDPG with SAC and  
230 MF-MARL(mean field-MARL, Yang et al. (2018) are barely satisfactory. But MAAC(Iqbal & Sha  
231 (2018)), ATOC(Jiang & Lu, 2018) and SMAC(ours) steadily performs when the number of agents  
232 increases. In future research, we will continue to improve the scalability when the number of agents  
233 further increases by sharing policies among agents and performing attention on sub-groups (of agents).  
234 We anticipate that in complicated scenarios, our method could work well.

## 235 5 CONCLUSIONS

236 We propose an algorithm, the SMAC(Scholastic-Actor-Critic) for training decentralized policies  
237 in multi-agent settings. We design a more powerful critic, the headmaster critic to learn how to  
238 group agents and when to communicate besides conventional ones. We also adapt useful advantage  
239 functions that avoid converging to non-optimal deterministic policies. We analyze the performance  
240 of the proposed approach compared the state-of-the-art methods on the Coin Game, CTC(Lerer &  
241 Peysakhovich, 2017), and MAgent(Zheng et al., 2018), concerning the number of agents. Thanks to  
242 the flexible setting, our results are promising in dynamic occasions with small training expenses. We  
243 intend to explore more to highly complex and dynamic environments.

## 244 REFERENCES

- 245 Abhishek Das, Satwik Kottur, José MF Moura, Stefan Lee, and Dhruv Batra. Learning cooperative  
246 visual dialog agents with deep reinforcement learning. In *Proceedings of the IEEE International  
247 Conference on Computer Vision*, pp. 2951–2960, 2017.
- 248 Richard Evans and Jim Gao. Deepmind ai reduces google data centre cooling bill by 40%. *DeepMind  
249 blog*, 20, 2016.
- 250 Jakob Foerster, Ioannis Alexandros Assael, Nando de Freitas, and Shimon Whiteson. Learning to  
251 communicate with deep multi-agent reinforcement learning. In *Advances in Neural Information  
252 Processing Systems*, pp. 2137–2145, 2016.
- 253 Jakob Foerster, Nantas Nardelli, Gregory Farquhar, Triantafyllos Afouras, Philip HS Torr, Pushmeet  
254 Kohli, and Shimon Whiteson. Stabilising experience replay for deep multi-agent reinforcement  
255 learning. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*,  
256 pp. 1146–1155. JMLR. org, 2017.
- 257 Jakob Foerster, Richard Y Chen, Maruan Al-Shedivat, Shimon Whiteson, Pieter Abbeel, and Igor  
258 Mordatch. Learning with opponent-learning awareness. In *Proceedings of the 17th International  
259 Conference on Autonomous Agents and MultiAgent Systems*, pp. 122–130. International Foundation  
260 for Autonomous Agents and Multiagent Systems, 2018a.
- 261 Jakob N Foerster, Gregory Farquhar, Triantafyllos Afouras, Nantas Nardelli, and Shimon Whiteson.  
262 Counterfactual multi-agent policy gradients. In *Thirty-Second AAAI Conference on Artificial  
263 Intelligence*, 2018b.
- 264 Jayesh K Gupta, Maxim Egorov, and Mykel Kochenderfer. Cooperative multi-agent control using  
265 deep reinforcement learning. In *International Conference on Autonomous Agents and Multiagent  
266 Systems*, pp. 66–83. Springer, 2017.
- 267 Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maxi-  
268 mum entropy deep reinforcement learning with a stochastic actor. *arXiv preprint arXiv:1801.01290*,  
269 2018.
- 270 He He, Jordan Boyd-Graber, Kevin Kwok, and Hal Daumé III. Opponent modeling in deep reinforce-  
271 ment learning. In *International Conference on Machine Learning*, pp. 1804–1813, 2016.
- 272 Shariq Iqbal and Fei Sha. Actor-attention-critic for multi-agent reinforcement learning. *arXiv preprint  
273 arXiv:1810.02912*, 2018.



- 274 Jiechuan Jiang and Zongqing Lu. Learning attentional communication for multi-agent cooperation.  
275 In *Advances in Neural Information Processing Systems*, pp. 7254–7264, 2018.
- 276 Vijay R Konda and John N Tsitsiklis. Actor-critic algorithms. In *Advances in neural information*  
277 *processing systems*, pp. 1008–1014, 2000.
- 278 Martin Lauer and Martin Riedmiller. An algorithm for distributed reinforcement learning in coop-  
279 erative multi-agent systems. In *In Proceedings of the Seventeenth International Conference on*  
280 *Machine Learning*. Citeseer, 2000.
- 281 Joel Z Leibo, Vinicius Zambaldi, Marc Lanctot, Janusz Marecki, and Thore Graepel. Multi-agent  
282 reinforcement learning in sequential social dilemmas. In *Proceedings of the 16th Conference on Au-*  
283 *tonomous Agents and MultiAgent Systems*, pp. 464–473. International Foundation for Autonomous  
284 Agents and Multiagent Systems, 2017.
- 285 Adam Lerer and Alexander Peysakhovich. Maintaining cooperation in complex social dilemmas  
286 using deep reinforcement learning. *arXiv preprint arXiv:1707.01068*, 2017.
- 287 Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa,  
288 David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *arXiv*  
289 *preprint arXiv:1509.02971*, 2015.
- 290 Ryan Lowe, Yi Wu, Aviv Tamar, Jean Harb, OpenAI Pieter Abbeel, and Igor Mordatch. Multi-agent  
291 actor-critic for mixed cooperative-competitive environments. In *Advances in Neural Information*  
292 *Processing Systems*, pp. 6379–6390, 2017.
- 293 Robert Duncan Luce and Howard Raiffa. *Games and decisions: Introduction and critical survey*.  
294 Wiley New York, 1958.
- 295 Laëticia Matignon, Guillaume J Laurent, and Nadine Le Fort-Piat. Hysteretic q-learning: an algorithm  
296 for decentralized reinforcement learning in cooperative multi-agent teams. In *2007 IEEE/RSJ*  
297 *International Conference on Intelligent Robots and Systems*, pp. 64–69. IEEE, 2007.
- 298 Luke Metz, Ben Poole, David Pfau, and Jascha Sohl-Dickstein. Unrolled generative adversarial  
299 networks. *arXiv preprint arXiv:1611.02163*, 2016.
- 300 Igor Mordatch and Pieter Abbeel. Emergence of grounded compositional language in multi-agent  
301 populations. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- 302 Junhyuk Oh, Valliappa Chockalingam, Satinder Singh, and Honglak Lee. Control of memory, active  
303 perception, and action in minecraft. *arXiv preprint arXiv:1605.09128*, 2016.
- 304 Shayegan Omidshafiei, Jason Pazis, Christopher Amato, Jonathan P How, and John Vian. Deep  
305 decentralized multi-task multi-agent reinforcement learning under partial observability. In *Pro-*  
306 *ceedings of the 34th International Conference on Machine Learning-Volume 70*, pp. 2681–2690.  
307 JMLR. org, 2017.
- 308 Peng Peng, Quan Yuan, Ying Wen, Yaodong Yang, Zhenkun Tang, Haitao Long, and Jun Wang.  
309 Multiagent bidirectionally-coordinated nets for learning to play starcraft combat games. *arXiv*  
310 *preprint arXiv:1703.10069*, 2, 2017.
- 311 Tuomas W Sandholm and Robert H Crites. Multiagent reinforcement learning in the iterated prisoner’s  
312 dilemma. *Biosystems*, 37(1-2):147–166, 1996.
- 313 David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche,  
314 Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering  
315 the game of go with deep neural networks and tree search. *nature*, 529(7587):484, 2016.
- 316 Sainbayar Sukhbaatar, Rob Fergus, et al. Learning multiagent communication with backpropagation.  
317 In *Advances in Neural Information Processing Systems*, pp. 2244–2252, 2016.
- 318 Sainbayar Sukhbaatar, Zeming Lin, Ilya Kostrikov, Gabriel Synnaeve, Arthur Szlam, and Rob  
319 Fergus. Intrinsic motivation and automatic curricula via asymmetric self-play. *arXiv preprint*  
320 *arXiv:1703.05407*, 2017.

- 321 Richard S Sutton, David A McAllester, Satinder P Singh, and Yishay Mansour. Policy gradient meth-  
322 ods for reinforcement learning with function approximation. In *Advances in neural information*  
323 *processing systems*, pp. 1057–1063, 2000.
- 324 Gerald Tesauro. Extending q-learning to general adaptive multi-agent systems. In *Advances in neural*  
325 *information processing systems*, pp. 871–878, 2004.
- 326 Hongwei Wang, Fuzheng Zhang, Jialin Wang, Miao Zhao, Wenjie Li, Xing Xie, and Minyi Guo.  
327 Ripplenet: Propagating user preferences on the knowledge graph for recommender systems. In *Pro-*  
328 *ceedings of the 27th ACM International Conference on Information and Knowledge Management*,  
329 pp. 417–426. ACM, 2018.
- 330 Yaodong Yang, Rui Luo, Minne Li, Ming Zhou, Weinan Zhang, and Jun Wang. Mean field multi-agent  
331 reinforcement learning. *arXiv preprint arXiv:1802.05438*, 2018.
- 332 Lianmin Zheng, Jiacheng Yang, Han Cai, Ming Zhou, Weinan Zhang, Jun Wang, and Yong Yu.  
333 Magent: A many-agent reinforcement learning platform for artificial collective intelligence. In  
334 *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.