

---

# Towards Skilled Population Curriculum for Multi-Agent Reinforcement Learning

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1       Recent advances in multi-agent reinforcement learning (MARL) allow agents to  
2       coordinate their behaviors in complex environments. However, common MARL  
3       algorithms still suffer from scalability and sparse reward issues. One promising  
4       approach to resolving them is *automatic curriculum learning* (ACL). ACL involves  
5       a *student* (curriculum learner) training on tasks of increasing difficulty controlled  
6       by a *teacher* (curriculum generator). Despite its success, ACL’s applicability is  
7       limited by (1) the lack of a general student framework for dealing with the varying  
8       number of agents across tasks and the sparse reward problem, and (2) the non-  
9       stationarity of the teacher’s task due to ever-changing student strategies. As a  
10      remedy for ACL, we introduce a novel automatic curriculum learning framework,  
11      Skilled Population Curriculum (SPC), which adapts curriculum learning to multi-  
12      agent coordination. Specifically, we endow the student with population-invariant  
13      communication and a hierarchical skill set, allowing it to learn cooperation and  
14      behavior skills from distinct tasks with varying numbers of agents. In addition, we  
15      model the teacher as a contextual bandit conditioned by student policies, enabling a  
16      team of agents to change its size while still retaining previously acquired skills. We  
17      also analyze the inherent non-stationarity of this multi-agent automatic curriculum  
18      teaching problem and provide a corresponding regret bound. Empirical results  
19      show that our method improves the performance, scalability and sample efficiency  
20      in several MARL environments. The source code and the video can be found at  
21      <https://sites.google.com/view/marl-spc/>.

## 22 1 Introduction

23      Multi-agent reinforcement learning (MARL) has long been a go-to tool in complex robotic and  
24      strategic domains [1, 2]. However, learning effective policies with sparse reward from scratch for  
25      large-scale multi-agent systems remains challenging. One of the challenges is the exponential growth  
26      of the joint observation-action space with an increasing number of agents. In addition, sparse reward  
27      signal requires a large number of training trajectories, posing difficulties in applying existing MARL  
28      algorithms directly to complex environments. As a result, these algorithms may produce agents that  
29      do not collaborate with each other, even when it would be of significant benefit [3, 4].

30      There are several lines of research related to the large-scale MARL problem with sparse reward,  
31      including reward shaping [5], curriculum learning [6], and learning from demonstrations [7]. Among  
32      these approaches, the curriculum learning paradigm, in which the difficulty of experienced tasks  
33      and the population of training agents progressively grow, shows particular promise. In *automatic*  
34      curriculum learning (ACL), a teacher (curriculum generator) learns to adjust the complexity and  
35      sequencing of tasks faced by a student (curriculum learner). Several works have even proposed *multi-*  
36      agent ACL algorithms, based on approximate or heuristic approaches to teaching, such as DyMA-CL

37 [8], EPC [9], and VACL [6]. However, these approaches rely on a framework of an off-policy student  
 38 with a replay buffer that is hard to decide the size of the replay buffer since the proportion of different  
 39 tasks matters. Also, they make a strong assumption that the value of the learned policy does not  
 40 change when agents switch to a different task. For example, In the football environment, when we  
 41 treat the score as the reward, the same state-action pairs of the team agents in different tasks might  
 42 lead to different returns. 3 learned agents could get more scores in a 3v1 match, while the same  
 43 three agents could get fewer scores in a 4v11 match with an unlearned random teammate. When  
 44 decomposing at the same state-action pairs, agents get different credit assignments. Moreover, the  
 45 teacher in these approaches still faces a non-stationarity problem due to the ever-changing student  
 46 strategies. Another class of larger-scale MARL solutions is hierarchical learning, which utilizes  
 47 temporal abstraction to decompose a task into a hierarchy of subtasks. This includes skill discovery  
 48 [10], option as response [11], role-based MARL [12], and two levels of abstraction [13]. However,  
 49 these approaches mostly focus on one specific task with a fixed number of agents and do not consider  
 50 the transferability of learned skills. In this paper, we provide our insight into this question:

51 *Whether an elaborate combination of principles from ACL and hierarchical learning can enable*  
 52 *complex cooperation with sparse reward in MARL?*

53 Specifically, we present a novel automatic curriculum learning algorithm, Skilled Population Curricu-  
 54 lum (SPC), that addresses the challenges of learning effective policies for large-scale multi-agent  
 55 systems with sparse reward. The core idea behind SPC, motivated by real-world team sports where  
 56 players often train their skills by gradually increasing the difficulty of tasks and the number of  
 57 coordinating players, is to encourage the student to learn skills from tasks with different numbers of  
 58 agents, akin to how team sports players train by gradually increasing the difficulty of tasks and the  
 59 number of coordinating players. To achieve this, SPC is implemented with three key components.  
 60 First, to solve the final complex cooperative tasks, we equip the contextual bandit teacher with an  
 61 RNN-based [14] imitation model to represent student policies and generate the bandit’s context.  
 62 Second, to handle the varying number of agents across these tasks and bypass the limitation of the  
 63 related studies, we utilize population-invariant communication in the student module is implemented  
 64 to handle varying number of agents across tasks. By treating each agent’s message as a word and  
 65 using a self-attention communication channel [15], SPC supports an arbitrary number of agents to  
 66 share messages. Third, to learn transferable skills in the sparse reward setting, a hierarchical skill  
 67 framework is used in the student module to learn transferable skills in the sparse reward setting,  
 68 where agents communicate on the high-level about a set of shared low-level policies. Empirical  
 69 results show that our method achieves state-of-the-art performance in several tasks in Multi-agent  
 70 Particle Environment (MPE) [16] and the challenging 5vs5 competition in Google Research Football  
 71 (GRF) [17].

## 72 2 Preliminaries

73 **Dec-POMDP.** A cooperative MARL problem can be formulated as a *decentralized par-*  
 74 *tially observable Markov decision process* (Dec-POMDP) [18], which is described as a tuple  
 75  $\langle n, \mathcal{S}, \mathcal{A}, P, R, \mathcal{O}, \Omega, \gamma \rangle$ , where  $n$  represents the number of agents.  $\mathcal{S}$  represents the space of global  
 76 states.  $\mathcal{A} = \{A_i\}_{i=1, \dots, n}$  denotes the space of actions of all agents.  $\mathcal{O} = \{O_i\}_{i=1, \dots, n}$  denotes  
 77 the space of observations of all agents.  $P : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$  denotes the state transition probability  
 78 function. All agents share the same reward as a function of the states and actions of the agents  
 79  $R : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ . Each agent  $i$  receives a private observation  $o_i \in O_i$  according to the observation  
 80 function  $\Omega(s, i) : \mathcal{S} \rightarrow O_i$ .  $\gamma \in [0, 1]$  denotes the discount factor.

81 **Multi-armed Bandit.** Multi-armed bandits (MABs) are a simple but very powerful framework that  
 82 repeatedly makes decisions under uncertainty. In this framework, a learner performs a sequence  
 83 of actions and immediately observes the corresponding reward after each action. The goal is to  
 84 maximize the total reward over a given set of  $K$  actions and a specific time horizon  $T$ . The measure  
 85 of success in MABs is often determined by the regret, which is the difference between the cumulative  
 86 reward of an MAB algorithm and the best-arm benchmark. One well-known MAB algorithm is the  
 87 Exp3 algorithm [19], which aims to increase the probability of selecting good arms and achieves a  
 88 regret of  $O(\sqrt{KT \log(K)})$  under a time-varying reward distribution. Another related concept is the  
 89 contextual bandit problem [20], where the learner makes decisions based on prior information as the  
 90 context.

### 91 3 Skilled Population Curriculum

92 In this section, we first provide a formal definition of the curriculum-enhanced Dec-POMDP frame-  
 93 work, which formulates the MARL with curriculum problem under the Dec-POMDP framework.  
 94 We then present our multi-agent ACL algorithm, Skilled Population Curriculum (SPC), as shown in  
 95 Fig. 1. In the following subsections, we establish the curriculum learning framework in Sec. 3.1, and  
 96 then present a contextual multi-armed bandit algorithm as the teacher to address the non-stationarity  
 97 in Sec. 3.2. Lastly, we introduce the student with transferable skills and population-invariant commu-  
 98 nication to tackle the varying number of agents and the sparse reward problem in Sec. 3.3.

#### 99 3.1 Problem Formulation

100 We consider environments from multi-agent automatic curriculum learning problems are equipped  
 101 with parameterized task spaces and thus can be modeled as curriculum-enhanced Dec-POMDPs.

102 **Definition 3.1** (Curriculum-enhanced Dec-POMDP). A curriculum-enhanced Dec-POMDP is defined  
 103 by a tuple  $\langle \Phi, \mathcal{M} \rangle$ , where  $\Phi$  and  $\mathcal{M}$  represent a task space and a Dec-POMDP, respectively. Given the  
 104 task  $\phi$ , the Dec-POMDP  $\mathcal{M}(\phi)$  is presented as  $\{n^\phi, \mathbf{S}^\phi, \mathbf{A}^\phi, P^\phi, r^\phi, O^\phi, \Omega^\phi, \gamma^\phi\}$ . The superscript  
 105  $\phi$  denotes that the Dec-POMDP elements are determined by the task  $\phi$ . Note that task  $\phi$  can be  
 106 a few parameters of the environment or task IDs in a finite task space. *In a curriculum-enhanced*  
 107 *Dec-POMDP, the objective is to improve the student’s performance on the target tasks through the*  
 108 *sequence of training tasks given by the teacher.*

109 Let  $\tau$  denote a trajectory whose unconditional distribution  $\Pr_\mu^{\pi, \phi}(\tau)$  (under a policy  $\pi$  and a task  $\phi$   
 110 with initial state distribution  $\mu(s_0)$ ) is  $\Pr_\mu^{\pi, \phi}(\tau) = \mu(s_0) \sum_{t=0}^{\infty} \pi(a_t | s_t) P^\phi(s_{t+1} | s_t, a_t)$ . We use  
 111  $p(\phi)$  to represent the distribution of target tasks and  $q(\phi)$  to represent the distribution of training tasks  
 112 at each task sampling step. We consider the joint agents’ policies  $\pi_\theta(a|s)$  and  $q_\psi(\phi)$  parameterized  
 113 by  $\theta$  and  $\psi$ , respectively. The overall objective to maximize in a curriculum-enhanced Dec-POMDP  
 114 is:

$$J(\theta, \psi) = \mathbb{E}_{\phi \sim p(\phi), \tau \sim \Pr_\mu^{\pi, \phi}} [R^\phi(\tau)] = \mathbb{E}_{\phi \sim q_\psi(\phi)} \left[ \frac{p(\phi)}{q_\psi(\phi)} V(\phi, \pi_\theta) \right] \quad (1)$$

115 where  $R^\phi(\tau) = \sum_t \gamma^t r^\phi(s_t, a_t; s_0)$  and  $V(\phi, \pi_\theta)$  represents the value function of  $\pi_\theta$  in Dec-  
 116 POMDP  $\mathcal{M}(\phi)$ . However, when optimizing  $q_\psi(\phi)$ , we cannot get the partial derivative  $\nabla_\psi J(\theta, \psi) =$   
 117  $\nabla_\psi \sum_\tau \frac{1}{q_\psi(\phi)} R^\phi(\tau) \Pr_\mu^{\pi, \phi}(\tau)$ <sup>1</sup> since the reward function and the transition probability function w.r.t  
 118 number of agents are non-parametric, non-differentiable, and discontinuous in most MARL scenarios.

119 Thus, we use the non-differentiable method, i.e., multi-armed bandit algorithms, to optimize  $q_\psi(\phi)$ ,  
 120 and use an RL algorithm (the student) in alternating periods to optimize  $\pi_\theta(a|s)$ . However, there are  
 121 three key challenges in solving this problem: (1) The teacher is facing a non-stationarity problem due  
 122 to the ever-changing student’s strategies. (2) The student will forget the old tasks and need to re-learn  
 123 them. Some tasks can be the prerequisites of other tasks, while some can be inter-independent and  
 124 parallel. (3) There is a lack of a general student framework to deal with the varying number of agents  
 125 across tasks and the sparse reward problem.

#### 126 3.2 Teacher as a Non-Stationary Contextual Bandit

127 As previously discussed, the teacher faces a non-stationarity problem due to the ever-changing  
 128 student’s strategies during the learning process. Specifically, as the student learns across different  
 129 tasks in different learning stages, the teacher will observe varying student performance when providing  
 130 the same task, resulting in a time-varying reward distribution for the teacher. In addition, the student  
 131 may forget previously learned policies. To mitigate this problem, the teacher should balance the  
 132 exploitation of tasks that have been found to benefit the student’s performance on the target tasks,  
 133 with the exploration of tasks that may not directly facilitate the student’s learning.

134 Fortunately, we notice that the non-stationarity stems from the student, which can be mitigated with  
 135 a contextual bandit which embeds the student policy into the context. As shown in Fig. 1 [Left](#),  
 136 the teacher utilizes the student’s policy representation as the context and chooses a task from the

<sup>1</sup> $p(\phi)$  is not in the partial derivative since it is a fixed distribution.

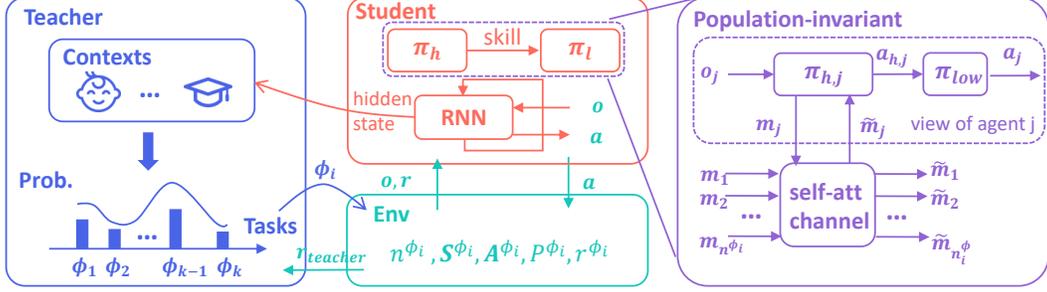


Figure 1: The overall framework of SPC. It consists of three parts: configurable environments, a teacher, and a student. **Left.** The teacher is modeled as a contextual multi-armed bandit. At each teacher timestep, the teacher chooses a training task from the distribution of bandit actions. **Mid.** The student is endowed with a hierarchical skill framework and population-invariant communication. It is trained with MARL algorithms on the training tasks. The student returns not only the hidden state of its RNN imitation model as contexts to the teacher, but also the average discounted cumulative rewards on the testing task. **Right.** The student learns hierarchical policies, with the population-invariant communication taking place at the high-level, implemented with a self-attention communication channel to handle the messages from a varying number of agents. The agents in the student share the same low-level policy.

137 distribution of training tasks. Specifically, we extend the Exp3 algorithm [19] by incorporating  
 138 contexts through a two-step online clustering process [21]. The context, represented by  $x$ , is the  
 139 student’s policy representation. The teacher’s action is a specific task, denoted by  $\phi$ , and the teacher’s  
 140 reward is the return of the student in the target tasks. The teacher’s algorithm is outlined in Alg. 1.  
 141 During the sampling stage (steps 1-5), the teacher selects a task for the student’s training. In the  
 142 training stage (steps 6-7), the teacher adjusts the parameters based on the evaluation reward received  
 143 from the student.

---

**Algorithm 1** Teacher Sampling and Training

---

**Input:** Context  $x$ , the number of Clusters  $N_c$ ,  $N_c$  instances of Exp3 with task distribution  $w(\phi_k, c)$  for  $k = 1, \dots, K$  and for  $c = 1, \dots, N_c$ , learning rate  $\alpha$ , a buffer maintaining the historical contexts

**Output:**  $\mathcal{M}(\phi) = \{n^\phi, S^\phi, A^\phi, P^\phi, r^\phi, O^\phi, \Omega^\phi, \gamma^\phi\}$ , the teacher bandit parameters

**Sampling**

1. Get the the context  $x$ , and save it to the buffer
2. Run the online cluster algorithm and get the index of the cluster center  $c(x)$
3. Let the active Exp3 instance be the instance with index  $c(x)$
4. Set the probability  $p(\phi_k, c(x)) = \frac{(1-\alpha)w(\phi_k, c(x))}{\sum_{j=1}^K w(\phi_j, c(x))} + \frac{\alpha}{K}$  for each task  $\phi_k$
5. Sample a new task according to the distribution of  $p_{\phi_k, c}$

**Training**

6. Get the return (discounted cumulative rewards) from student testing  $r$
  7. Update the active Exp3 instance by setting  $w(\phi_k, c(x)) = w(\phi_k, c(x))e^{\alpha r/K}$
- 

144 **3.2.1 Context Representation**

145 Upon analysis, it is essential to learn an effective representation for the student’s policy as the context.  
 146 One straightforward representation is to use the student parameters  $\theta$  directly as the context. However,  
 147 the number of parameters is too large to be used as the input of neural network if we change the  
 148 student’s architecture. Therefore, we propose an alternative method.

149 A principle for learning a good representation of a policy is *predictive representation*, which means  
 150 the representation should be accurate to predict policy actions given states. In accordance with this  
 151 principle, we utilize an imitation function through supervised learning. Supervised learning does  
 152 not require direct access to reward signals, making it an attractive approach for reward-agnostic  
 153 representation learning. Intuitively, the imitation function attempts to mimic low-level policy based

154 on historical behaviors. In practice, we use an RNN-based imitation function  $f_{im} : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ .  
 155 Since recurrent neural networks are theoretically Turing complete [22], their internal states can be  
 156 used as the representation of the student’s policy. We train this imitation function by using the  
 157 negative cross entropy objective  $\mathbb{E}[\log f_{im}(s, a)]$ .

### 158 3.2.2 Regret Analysis

159 In this subsection, we demonstrate that the proposed teacher algorithm has a regret bound of  
 160  $\mathbb{E}[R(T)] = O(T^{2/3}(LK \log T)^{1/3})$ , where  $T$  is the number of total rounds,  $L$  is the Lipschitz  
 161 constant, and  $K$  is the number of arms (the number of the teacher’s actions). The regret analysis  
 162 is used to justify the usage of the bandit algorithm in the non-stationary setting. The regret bound  
 163 represents the optimality of SPC, as the teacher’s reward is the return of the student in the target tasks.

164 First, we introduce the Lipschitz assumption about the generalization ability of the task space.

165 **Assumption 3.2** (Lipschitz continuity w.r.t the context). Without loss of generality, the contexts are  
 166 mapped into the  $[0, 1]$  interval, so that the expected rewards for the teacher are Lipschitz with respect  
 167 to the context.

$$|r(\phi | x) - r(\phi | x')| \leq L \cdot |x - x'|$$

for any arm  $\phi \in \Phi$  and any pair of contexts  $x, x' \in \mathcal{X}$  (2)

168 where  $L$  is the Lipschitz constant, and  $\mathcal{X}$  is the context space.

169 This assumption suggests that for any policy trained on a set of tasks, the rate at which performance  
 170 improves is not faster than the rate at which the policy changes. This is a realistic assumption, as we  
 171 cannot expect the student to achieve a significant improvement on a task with only a few training  
 172 steps under a new context. We use an existing contextual bandit algorithm for a limited number of  
 173 contexts [19] (see Appendix A) and Lemma 3.3 as a foundation for proving Theorem 3.4.

174 **Lemma 3.3.** *Alg. 2 has a regret bound of  $\mathbb{E}[R(T)] = O(\sqrt{TK|\mathcal{X}|\log K})$ .*

175 Lemma 3.3 introduces a square root dependence on  $|\mathcal{X}|$  if separate copies of Exp3 are run for  
 176 each context [19]. This motivates us to address the large context space by utilizing discretization  
 177 techniques.

178 **Theorem 3.4.** *Consider the Lipschitz contextual bandit problem with contexts in  $[0, 1]$ . The Alg. 1  
 179 yields regret  $\mathbb{E}[R(T)] = O(T^{2/3}(LK \ln T)^{1/3})$ .*

180 *Proof.* See Appendix B. □

181 In practice, the high-dimensional context space cannot be discretized using a uniform mesh in  $[0, 1]$   
 182 as in the proof of Theorem 3.4. To address this issue, we utilize the Balanced Iterative Reducing  
 183 and Clustering using Hierarchies (BIRCH) online clustering algorithm [21] to discretize the context  
 184 space. BIRCH is an efficient and easy-to-update algorithm that can effectively cluster large datasets.  
 185 In this case, it is used to cluster the high-dimensional RNN-based policy representation. The resulting  
 186 clusters can be seen as an approximation of a uniform mesh.

### 187 3.3 Student with Population-Invariant Skills

188 We propose a population-invariant skill framework to address the challenges of varying number of  
 189 agents and sparse reward problem. This framework allows agents to communicate via a self-attention  
 190 channel, enabling them to learn transferable skills across different tasks. The student module is  
 191 designed to be algorithm-agnostic and is orthogonal to any state-of-the-art MARL algorithm. While  
 192 there have been some efforts in the literature to address the varying number of agents [23, 24], these  
 193 approaches heavily rely on prior knowledge of the environments.

194 **Population-Invariant Teamwork Communication.** In order to enable the population-invariant  
 195 property and learn tactics among agents, we introduce communication. Leveraging the transformer  
 196 architecture’s capability to process inputs of varying lengths [15], we incorporate self-attention into  
 197 our communication mechanism. As illustrated in Fig. 1 Right, each agent  $j$  receives an observation  
 198  $o_j$  and encodes it into a message vector  $m_j = f(o_j)$  which is then sent through a self-attention  
 199 channel, where  $f$  is an observation encoder function.

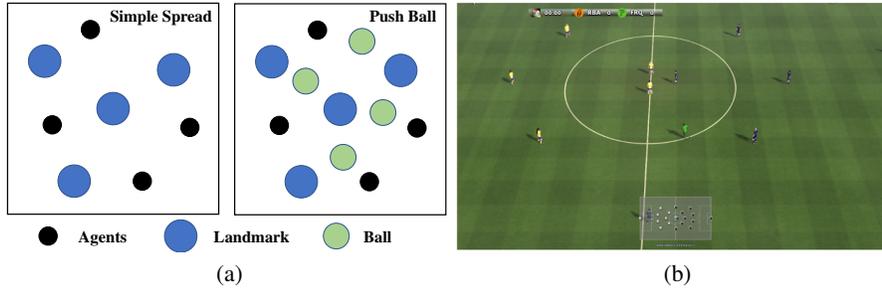


Figure 2: (a) Multi-agent Particle Environment. (b) Google Research Football.

200 The channel aggregates all messages and sends the new message vector,  $\tilde{m}_j$ , through the self-attention  
 201 mechanism. Concretely, given the channel input  $\mathbf{M} = [m_1, m_2, \dots, m_n] \in R^{n \times d_m}$ , and the  
 202 trainable weight of the channel  $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V \in R^{d_m \times d_m}$ , we obtain three distinct representations:  
 203  $\mathbf{Q} = \mathbf{M}\mathbf{W}_Q, \mathbf{K} = \mathbf{M}\mathbf{W}_K, \mathbf{V} = \mathbf{M}\mathbf{W}_V$ . Then the output messages are

$$\tilde{\mathbf{M}} = \text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax} \left( \frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_m}} \right) \mathbf{V} \quad (3)$$

204 where  $d_m$  is the dimension of the messages. As the dimensions of the trainable weight are independent  
 205 of the number of agents, our student models can leverage the population-invariant property to  
 206 effectively learn tactics.

207 **Transferable Hierarchical Skills.** As depicted in the dotted box in Fig. 1 Right, after receiving  
 208 the new messages  $\tilde{m}_j$  from the channel, each agent employs a high-level action (skill)  $a_{h,j} =$   
 209  $\pi_{h,j}(o_j, \tilde{m}_j)$  to execute the low-level policy  $a_j = \pi_{low}(o_j, a_{h,j})$ . In this work, we generalize the  
 210 high-level action (skill)  $a_{h,j}$  to a continuous embedding space, so that the skill can be either a latent  
 211 continuous vector as in DIAYN [25], or a categorical distribution for sampling discrete options [26].

212 **Implementation.** We implement the high- and low-level policies in the student with Proximal Policy  
 213 Optimization (PPO) [27]. Following the common practice proposed in [28], the high-level policy  
 214 for each agent is learned independently, whereas the low-level policies share parameters, as the  
 215 fundamental action pattern should be consistent among different agents. The low-level agents are  
 216 rewarded by the environment, while the high-level policy is trained to take actions at fixed intervals.  
 217 Within this interval, the cumulative low-level reward is used as the high-level reward. When using  
 218 a categorical distribution to enable discrete skills, we sample an ‘‘option’’ from the distribution and  
 219 provide the corresponding one-hot embedding to the low-level policy.

## 220 4 Related Work

221 **Automatic Curriculum Learning in MARL.** Curriculum learning is a training strategy that mimics  
 222 the human learning process by organizing tasks based on their difficulty level [29]. The selection of  
 223 tasks is formulated as a Curriculum Markov Decision Process (CMDP) [30]. Automatic Curriculum  
 224 Learning mechanisms aim to learn a task selection function based on past interactions, such as ADR  
 225 [31, 32], ALP-GMM [33], SPCL [34], GoalGAN [35], PLR [36, 37], SPDL [38], CURROT [39],  
 226 and graph-curriculum [40]. Recently, several MARL curriculum learning frameworks have been  
 227 proposed, such as open-ended evolution [41–43], population-based training [44, 45], meta-learning  
 228 [46, 47] and training with emergent curriculum [48, 49, 29]. In summary, these frameworks share a  
 229 common principle of an automatic curriculum that continually generates improved agents through  
 230 selection pressure among a population of self-optimizing agents.

231 **Hierarchical MARL and Communication.** Hierarchical reinforcement learning (HRL) has been  
 232 extensively studied to address the issue of sparse reward and facilitate transfer learning. Single-agent  
 233 HRL focuses on learning the temporal decomposition of tasks, either by learning subgoals [50–  
 234 54] or by discovering reusable skills [55–58]. Recent developments in hierarchical MARL have  
 235 been discussed in Sec. 1. In multi-agent settings, communication has been effective in promoting  
 236 cooperation among agents [59–65]. However, current approaches that extend HRL to multi-agent  
 237 systems or utilize communication are limited to a fixed number of agents and lack the ability to  
 238 transfer to different agent counts.

## 239 5 Experiments

240 To demonstrate the effectiveness of our approach, we conduct experiments on several tasks in two  
241 environments: Simple-Spread and Push-Ball in the Multi-agent Particle Environment (MPE) [16],  
242 and the challenging 5vs5 task of the Google Research Football (GRF) environment [17]. We aim to  
243 investigate the following research questions:

244 **Q1:** *Is curriculum learning necessary in complex large-scale MARL problems?* (Sec. 5.2)

245 **Q2:** *Can SPC outperform previous curriculum-based MARL methods? If so, which components of*  
246 *SPC contribute the most to performance gains?* (Sec. 5.3)

247 **Q3:** *Can SPC effectively learn a curriculum for the student?* (Sec. 5.4)

### 248 5.1 Environments, Baselines and Metric

249 **Environments.** In the GRF 5vs5 scenario, we control four agents, excluding the goalkeeper, to  
250 compete against the built-in AI opponents. Each agent observes a compact encoding, consisting of a  
251 115-dimensional vector that summarizes various aspects of the game, such as player coordinates, ball  
252 possession and direction, active players, and game mode. The available action set for an individual  
253 agent includes 19 discrete actions, such as idle, move, pass, shoot, dribble, etc. The GRF provides  
254 two types of rewards: scoring and checkpoints, to encourage agents to move the ball forward and  
255 make successful shots. Additionally, we include a shooting reward in the challenging GRF 5vs5  
256 task. We select several basic scenarios in GRF, including 3vs3, Pass-Shoot, 3vs1, and Empty-Goal as  
257 curriculum.

258 In MPE, we investigate Simple-Spread and Push-Ball (see Fig. 2a). In Simple-Spread, there are  $n$   
259 agents that need to cover all  $n$  landmarks. Agents are penalized for collisions and only receive a  
260 positive reward when all the landmarks are covered. In Push-Ball, there are  $n$  agents,  $n$  balls, and  $n$   
261 landmarks. The agents must push the balls to cover each landmark. A success reward is given after  
262 all the landmarks have been covered.

263 **Baselines.** We compare our approach to the following methods in Table 1 as baselines<sup>2</sup>:

264 **Metric.** To evaluate the performance of  
265 our approach in the GRF 5vs5 scenario, we  
266 use metrics beyond just the mean episode  
267 reward, as this alone may not accurately re-  
268 flect the agents’ performance. Specifically,  
269 we use the win rate and the average goal  
270 difference, which is calculated as the num-  
271 ber of goals scored by the MARL agents  
272 minus the number of goals scored by the  
273 opposing team.

Table 1: Baseline algorithms.

Categories	Methods
MARL (Q1)	QMIX [68] IPPO [69]
Curriculum-based (Q2)	IPPO with uniform task sampling VACL [6]
Ablation Study (Q3)	SPC with uniform task sampling SPC without HRL and COM

274 We evaluate the performance of MARL algorithms to justify the need for curriculum learning in  
275 complex large-scale MARL problems. To ensure a fair comparison, we modify VACL by removing  
276 the centralized critic for MPE tasks. Centralized Training Decentralized Execution methods is not  
277 included as baselines since they are not suitable for varying numbers (e.g., MADDPG/MAPPO’s  
278 critic requires a fixed size of input or QMIX’s mixing network also fixed size of the input).

279 In all experiments, we use individual Proximal Policy Optimization (IPPO) as the backend MARL  
280 algorithm. To ensure the robustness of our results, we conduct experiments on a 30-node cluster, with  
281 one node containing a 128-core CPU and four A100 GPUs. Each trial of the experiment is repeated  
282 over five seeds and runs for 1-2 days.

### 283 5.2 The Necessity of Curriculum Learning

284 Our experiments first show that in simple environments, such as MPE, students can directly learn  
285 to complete the task without the need for curriculum. For MPE experiments, we randomly select  
286 a starting state and the episode ends after a fixed number of maximum steps. Specifically, the task

<sup>2</sup>We also run CDS [66] and CMARL [67], but we have not included their performance because the goal difference reported in CMARL [67] is relatively low compared to our method.

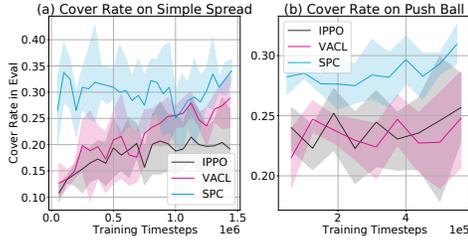


Figure 3: The evaluation performance of various methods on MPE.

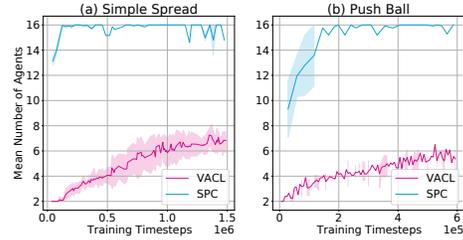


Figure 4: The changes in the number of agents on MPE.

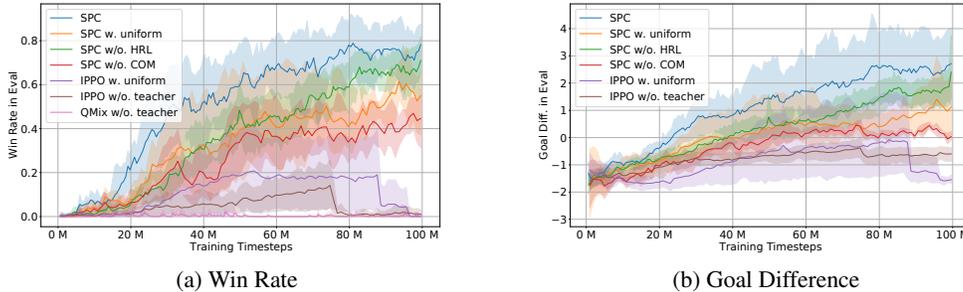


Figure 5: The evaluation performance of various methods on 5vs5 football competition. (p-value is less than 0.05 which means the results are statistically significant.)

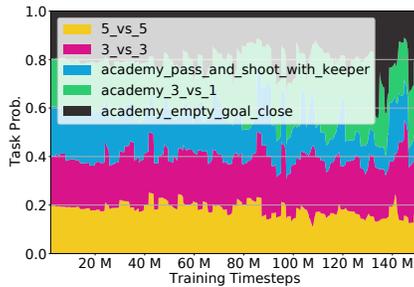
287 space consists of  $n$  agents, where  $n \in \{2, 4, 8, 16\}$ , and the maximum allowed steps is set to 25. All  
 288 evaluations are performed on the target task, with  $n = 16$ . IPPO is trained and evaluated directly on  
 289 the target task, and results in Fig. 3 demonstrate that it performs similarly to the VACL algorithm.  
 290 We plot the performance within a sliding window so that the starting point is not exactly from 0  
 291 timestep. VACL uses entity progression, which is a rule-based curriculum update mechanism so it  
 292 lacks the flexibility to switch the curriculum when relatively easy tasks can be learned quickly. The  
 293 reason for the performance jump is that SPC can switch to the largest population rapidly, which we  
 294 consider one advantage of SPC. Additionally, we observe that the SPC approach only achieves a  
 295 slightly higher coverage rate than the baseline methods. Furthermore, we investigate the probability  
 296 variation of different population sizes, shown in Fig. 4. We observe that the curriculum provided  
 297 by SPC is approaching the target task. These results suggest that in simple environments where the  
 298 student can learn to directly complete the task, curriculum learning may not be necessary.

299 When it comes to more complex scenarios, such as the 5vs5 task in GRF, our results demonstrate  
 300 that curriculum learning is a promising solution. As shown in Fig. 5a, without curriculum learning,  
 301 QMix and IPPO cannot perform well in the 5vs5 scenario, and IPPO is slightly better than QMix. In  
 302 Fig. 5b, we omit the curve of QMix as its mean score is low and affects the presentation of the figure.  
 303 The reason could be that QMix is an off-policy MARL algorithm, which would rely heavily on the  
 304 replay buffer. However, in such sparse reward scenarios, the replay buffer has much less effective  
 305 samples for QMix to learn. For example, the replay buffer would contain tons of zero-score samples,  
 306 leading to a non-promising performance. Meanwhile, IPPO, with its on-policy nature, is able to  
 307 achieve better sample efficiency and outperform off-policy algorithms like QMix in such scenarios.  
 308 Though MARL methods can achieve good performance in basic scenarios in GRF, they fail to solve  
 309 complex scenarios such as the 5vs5 task. Therefore, curriculum learning is a promising solution to  
 310 the complex large-scale MARL problem.

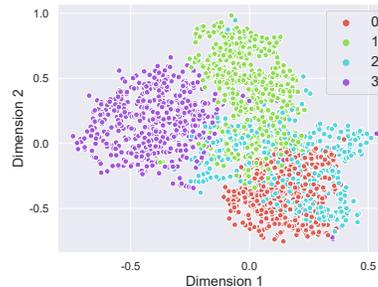
### 311 5.3 Performance and Ablation Study

312 Our study demonstrates that SPC outperforms VACL in MPE tasks. Instead of training with a  
 313 continuous relaxation of the population size variable as in VACL, our bandit teacher achieves a higher  
 314 success rate at test time, since the population size is a discrete variable in nature. Furthermore, the  
 315 curriculum provided by SPC is effective in exploring the task space and converge to the target task  
 316 when the task is relatively simple and curriculum is not necessary, as shown in Fig. 4.

317 In GRF experiments, we do not include VACL in our baselines in the GRF, as its implementation  
 318 relies heavily on prior knowledge of specific scenarios, such as the thresholds to divide the learning



(a) The task distribution of SPC during training.



(b) The visualization of contexts

Figure 6: Visualization of Learned Curriculum.

319 process. Fig. 6 indicates that SPC has higher win rate and goal difference than IPPO with uniform  
 320 task sampling in the 5vs5 competition. These experiments demonstrate that when the teacher is  
 321 rewarded by the student’s performance, a bandit-based teacher can exploit the student’s learning stage  
 322 and provide suitable training tasks.

323 In our ablation study, we examine the impact of two key components of our SPC algorithm: the  
 324 contextual multi-armed bandit teacher and the hierarchical structure of the student framework. By  
 325 replacing the former with uniform task sampling and removing the latter, As shown in Fig. 5a and  
 326 Fig. 5b, SPC can achieve a higher win rate and a greater score difference than SPC with uniform and  
 327 SPC without HRL. Furthermore, SPC with uniform task sampling outperforms IPPO with uniform  
 328 task sampling. This highlights the importance of HRL in the 5vs5 football competition, and suggests  
 329 that both the contextual multi-armed bandit and the hierarchical structure contribute equally to the  
 330 performance of SPC. When removing HRL and bandit, the performance degradation w.r.t. SPC are  
 331 similar. However, it should be noted that SPC with uniform task sampling has a larger variance in  
 332 performance than SPC without HRL, indicating that uniform sampling may introduce more undesired  
 333 tasks for student training. Overall, these results further justify the necessity of SPC in complex  
 334 large-scale MARL problems<sup>3</sup>.

### 335 5.4 Visualization of Learned Curriculum

336 We visualize the distribution of task sampling of SPC during training based on a selected trial as  
 337 shown in Fig. 6a. At the beginning of training, the task probability appears to be near-uniform, as  
 338 the teacher explores the task space and keeps track of the student’s learning status, acting as an  
 339 anti-forgetting mechanism. As training progresses, the probabilities change over time. For example,  
 340 the proportions of 3vs1 and Empty-Goal tasks gradually drop as the student becomes proficient in  
 341 these scenarios. We also visualize the distribution of contexts in Fig. 6b using t-SNE [70], where the  
 342 contexts are collected and stored in a buffer. We divide the contexts into four classes according to the  
 343 index, and different parts represent different contexts of the final student policy representation.

## 344 6 Discussion

345 **Conclusion.** We present Skilled Population Curriculum (SPC), a novel multi-agent ACL algorithm  
 346 that addresses scalability and sparse reward issues in multi-agent systems. SPC learns complex  
 347 behaviors from scratch by incorporating a population-invariant multi-agent communication framework  
 348 and using a hierarchical scheme for agents to learn skills. Moreover, SPC mitigates non-stationarity  
 349 by modeling the teacher as a contextual bandit, where the context is represented by the student’s  
 350 policy representation. Though our design choices focus on solving the GRF 5vs5 task, we believe  
 351 that analyzing and addressing these issues is crucial for further development in multi-agent ACL  
 352 algorithms. While SPC may be complex to implement due to its various components, we provide  
 353 clean and well-organized code for ease of use.

354 **Limitations.** We acknowledge that there are limitations of our algorithm. SPC is over-designed for  
 355 simple tasks since our objective is to solve difficult tasks. Also, it would be interesting to understand  
 356 the impact of varying number of agents on the dynamics of the environment.

<sup>3</sup>We also demonstrate the performance of SPC in the GRF 11vs11 full game (see Appendix C).

## References

- 357
- 358 [1] RoboCup. Robocup Federation Official Website. <https://www.robocup.org/>, 2019. Ac-  
359 cessed April 10, 2019.
- 360 [2] OpenAI. OpenAI Five. <https://openai.com/blog/openai-five/>, 2019. Accessed March  
361 4, 2019.
- 362 [3] Kaiqing Zhang, Zhuoran Yang, and Tamer Başar. Multi-agent reinforcement learning: A  
363 selective overview of theories and algorithms. *Handbook of Reinforcement Learning and*  
364 *Control*, pages 321–384, 2021.
- 365 [4] Yaodong Yang and Jun Wang. An overview of multi-agent reinforcement learning from game  
366 theoretical perspective. *arXiv preprint arXiv:2011.00583*, 2020.
- 367 [5] Yujing Hu, Weixun Wang, Hangtian Jia, Yixiang Wang, Yingfeng Chen, Jianye Hao, Feng Wu,  
368 and Changjie Fan. Learning to utilize shaping rewards: A new approach of reward shaping.  
369 *arXiv preprint arXiv:2011.02669*, 2020.
- 370 [6] Jiayu Chen, Yuanxin Zhang, Yuanfan Xu, Huimin Ma, Huazhong Yang, Jiaming Song, Yu Wang,  
371 and Yi Wu. Variational automatic curriculum learning for sparse-reward cooperative multi-agent  
372 problems. *Advances in Neural Information Processing Systems*, 34, 2021.
- 373 [7] Shiyu Huang, Wenze Chen, Longfei Zhang, Ziyang Li, Fengming Zhu, Deheng Ye, Ting  
374 Chen, and Jun Zhu. Tikick: Toward playing multi-agent football full games from single-agent  
375 demonstrations. *arXiv preprint arXiv:2110.04507*, 2021.
- 376 [8] Weixun Wang, Tianpei Yang, Yong Liu, Jianye Hao, Xiaotian Hao, Yujing Hu, Yingfeng Chen,  
377 Changjie Fan, and Yang Gao. From few to more: Large-scale dynamic multiagent curriculum  
378 learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages  
379 7293–7300, 2020.
- 380 [9] Qian Long, Zihan Zhou, Abhibav Gupta, Fei Fang, Yi Wu, and Xiaolong Wang. Evolutionary  
381 population curriculum for scaling multi-agent reinforcement learning. *arXiv preprint*  
382 *arXiv:2003.10423*, 2020.
- 383 [10] Jiachen Yang, Igor Borovikov, and Hongyuan Zha. Hierarchical cooperative multi-agent  
384 reinforcement learning with skill discovery. *arXiv preprint arXiv:1912.03558*, 2019.
- 385 [11] Alexander Sasha Vezhnevets, Yuhuai Wu, Remi Leblond, and Joel Z Leibo. Options as responses:  
386 Grounding behavioural hierarchies in multi-agent rl. *arXiv preprint arXiv:1906.01470*, 2019.
- 387 [12] Tonghan Wang, Tarun Gupta, Anuj Mahajan, Bei Peng, Shimon Whiteson, and Chongjie Zhang.  
388 Rode: Learning roles to decompose multi-agent tasks. *arXiv preprint arXiv:2010.01523*, 2020.
- 389 [13] Zhen-Jia Pang, Ruo-Ze Liu, Zhou-Yu Meng, Yi Zhang, Yang Yu, and Tong Lu. On reinforcement  
390 learning for full-length game of starcraft. In *Proceedings of the AAAI Conference on Artificial*  
391 *Intelligence*, 2019.
- 392 [14] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):  
393 1735–1780, 1997.
- 394 [15] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,  
395 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in Neural Information*  
396 *Processing Systems*, 30, 2017.
- 397 [16] Ryan Lowe, Yi Wu, Aviv Tamar, Jean Harb, Pieter Abbeel, and Igor Mordatch. Multi-agent  
398 actor-critic for mixed cooperative-competitive environments. *Advances in Neural Information*  
399 *Processing Systems*, 2017.
- 400 [17] Karol Kurach, Anton Raichuk, Piotr Stańczyk, Michal Zajkac, Olivier Bachem, Lasse Espeholt,  
401 Carlos Riquelme, Damien Vincent, Marcin Michalski, Olivier Bousquet, et al. Google research  
402 football: A novel reinforcement learning environment. *arXiv preprint arXiv:1907.11180*, 2019.

- 403 [18] Daniel S Bernstein, Robert Givan, Neil Immerman, and Shlomo Zilberstein. The complexity of  
404 decentralized control of Markov Decision Processes. *Mathematics of Operations Research*, 27  
405 (4):819–840, 2002.
- 406 [19] Peter Auer, Nicolo Cesa-Bianchi, Yoav Freund, and Robert E Schapire. The nonstochastic  
407 multiarmed bandit problem. *SIAM Journal on Computing*, 32(1):48–77, 2002.
- 408 [20] Elad Hazan and Nimrod Megiddo. Online learning with prior knowledge. In *International  
409 Conference on Computational Learning Theory*, pages 499–513. Springer, 2007.
- 410 [21] Tian Zhang, Raghu Ramakrishnan, and Miron Livny. Birch: An efficient data clustering method  
411 for very large databases. *ACM Aigmod Record*, 25(2):103–114, 1996.
- 412 [22] Heikki Hyötyniemi. Turing machines are recurrent neural networks. In *STeP '96/Publications  
413 of the Finnish Artificial Intelligence Society*, 1996.
- 414 [23] Shariq Iqbal, Christian A Schroeder De Witt, Bei Peng, Wendelin Böhmer, Shimon Whiteson,  
415 and Fei Sha. Randomized entity-wise factorization for multi-agent reinforcement learning. In  
416 *International Conference on Machine Learning*, pages 4596–4606. PMLR, 2021.
- 417 [24] Siyi Hu, Fengda Zhu, Xiaojun Chang, and Xiaodan Liang. Updet: Universal multi-agent rein-  
418 forcement learning via policy decoupling with transformers. *arXiv preprint arXiv:2101.08001*,  
419 2021.
- 420 [25] Benjamin Eysenbach, Abhishek Gupta, Julian Ibarz, and Sergey Levine. Diversity is all you  
421 need: Learning skills without a reward function. *arXiv preprint arXiv:1802.06070*, 2018.
- 422 [26] Pierre-Luc Bacon, Jean Harb, and Doina Precup. The option-critic architecture. In *Proceedings  
423 of the AAAI Conference on Artificial Intelligence*, volume 31, 2017.
- 424 [27] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal  
425 policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 426 [28] Wei Fu, Chao Yu, Zelai Xu, Jiaqi Yang, and Yi Wu. Revisiting some common practices in  
427 cooperative multi-agent reinforcement learning. *arXiv preprint arXiv:2206.07505*, 2022.
- 428 [29] Rémy Portelas, Cédric Colas, Lilian Weng, Katja Hofmann, and Pierre-Yves Oudeyer. Au-  
429 tomatic curriculum learning for deep RL: A short survey. *arXiv preprint arXiv:2003.04664*,  
430 2020.
- 431 [30] Sanmit Narvekar and Peter Stone. Learning curriculum policies for reinforcement learning.  
432 *arXiv preprint arXiv:1812.00285*, 2018.
- 433 [31] Ilge Akkaya, Marcin Andrychowicz, Maciek Chociej, Mateusz Litwin, Bob McGrew, Arthur  
434 Petron, Alex Paino, Matthias Plappert, Glenn Powell, Raphael Ribas, et al. Solving rubik’s cube  
435 with a robot hand. *arXiv preprint arXiv:1910.07113*, 2019.
- 436 [32] Bhairav Mehta, Manfred Diaz, Florian Golemo, Christopher J Pal, and Liam Paull. Active  
437 domain randomization. In *Conference on Robot Learning*, pages 1162–1176. PMLR, 2020.
- 438 [33] Rémy Portelas, Cédric Colas, Katja Hofmann, and Pierre-Yves Oudeyer. Teacher algorithms  
439 for curriculum learning of deep RL in continuously parameterized environments. In *Conference  
440 on Robot Learning*, pages 835–853, 2020.
- 441 [34] Lu Jiang, Deyu Meng, Qian Zhao, Shiguang Shan, and Alexander G Hauptmann. Self-paced  
442 curriculum learning. In *Twenty-Ninth AAAI Conference on Artificial Intelligence*, 2015.
- 443 [35] Carlos Florensa, David Held, Xinyang Geng, and Pieter Abbeel. Automatic goal generation  
444 for reinforcement learning agents. In *International Conference on Machine Learning*, pages  
445 1515–1528. PMLR, 2018.
- 446 [36] Minqi Jiang, Edward Grefenstette, and Tim Rocktäschel. Prioritized level replay. In *Internat-  
447 ional Conference on Machine Learning*, pages 4940–4950. PMLR, 2021.

- 448 [37] Minqi Jiang, Michael Dennis, Jack Parker-Holder, Jakob Foerster, Edward Grefenstette, and Tim  
449 Rocktäschel. Replay-guided adversarial environment design. *Advances in Neural Information*  
450 *Processing Systems*, 34:1884–1897, 2021.
- 451 [38] Pascal Klink, Carlo D’Eramo, Jan R Peters, and Joni Pajarinen. Self-paced deep reinforcement  
452 learning. *Advances in Neural Information Processing Systems*, 33:9216–9227, 2020.
- 453 [39] Pascal Klink, Haoyi Yang, Carlo D’Eramo, Jan Peters, and Joni Pajarinen. Curriculum rein-  
454 forcement learning via constrained optimal transport. In *International Conference on Machine*  
455 *Learning*, pages 11341–11358. PMLR, 2022.
- 456 [40] Maxwell Svetlik, Matteo Leonetti, Jivko Sinapov, Rishi Shah, Nick Walker, and Peter Stone.  
457 Automatic curriculum graph generation for reinforcement learning agents. In *Proceedings of*  
458 *the AAAI Conference on Artificial Intelligence*, volume 31, 2017.
- 459 [41] Wolfgang Banzhaf, Bert Baumgaertner, Guillaume Beslon, René Doursat, James A Foster,  
460 Barry McMullin, Vinicius Veloso De Melo, Thomas Miconi, Lee Spector, Susan Stepney, et al.  
461 Defining and simulating open-ended novelty: Requirements, guidelines, and challenges. *Theory*  
462 *in Biosciences*, 135(3):131–161, 2016.
- 463 [42] Joel Lehman, Kenneth O Stanley, et al. Exploiting open-endedness to solve problems through  
464 the search for novelty. In *ALIFE*, pages 329–336. Citeseer, 2008.
- 465 [43] Russell K Standish. Open-ended artificial evolution. *International Journal of Computational*  
466 *Intelligence and Applications*, 3(02):167–175, 2003.
- 467 [44] Max Jaderberg, Wojciech M Czarnecki, Iain Dunning, Luke Marris, Guy Lever, Antonio Garcia  
468 Castaneda, Charles Beattie, Neil C Rabinowitz, Ari S Morcos, Avraham Ruderman, et al.  
469 Human-level performance in 3d multiplayer games with population-based reinforcement learn-  
470 ing. *Science*, 364(6443):859–865, 2019.
- 471 [45] Siqi Liu, Guy Lever, Josh Merel, Saran Tunyasuvunakool, Nicolas Heess, and Thore Graepel.  
472 Emergent coordination through competition. *arXiv preprint arXiv:1902.07151*, 2019.
- 473 [46] Abhinav Gupta, Marc Lanctot, and Angeliki Lazaridou. Dynamic population-based meta-  
474 learning for multi-agent communication with natural language. *Advances in Neural Information*  
475 *Processing Systems*, 34:16899–16912, 2021.
- 476 [47] Rémy Portelas, Clément Romac, Katja Hofmann, and Pierre-Yves Oudeyer. Meta automatic  
477 curriculum learning. *arXiv preprint arXiv:2011.08463*, 2020.
- 478 [48] Bowen Baker, Ingmar Kanitscheider, Todor Markov, Yi Wu, Glenn Powell, Bob McGrew,  
479 and Igor Mordatch. Emergent tool use from multi-agent autotutorials. *arXiv preprint*  
480 *arXiv:1909.07528*, 2019.
- 481 [49] Joel Z Leibo, Edward Hughes, Marc Lanctot, and Thore Graepel. Autotutorials and the  
482 emergence of innovation from social interaction: A manifesto for multi-agent intelligence  
483 research. *arXiv preprint arXiv:1903.00742*, 2019.
- 484 [50] Ofir Nachum, Shixiang Shane Gu, Honglak Lee, and Sergey Levine. Data-efficient hierarchical  
485 reinforcement learning. *Advances in neural information processing systems*, 31, 2018.
- 486 [51] Ofir Nachum, Shixiang Gu, Honglak Lee, and Sergey Levine. Near-optimal representation  
487 learning for hierarchical reinforcement learning. *arXiv preprint arXiv:1810.01257*, 2018.
- 488 [52] Sainbayar Sukhbaatar, Emily Denton, Arthur Szlam, and Rob Fergus. Learning goal embeddings  
489 via self-play for hierarchical reinforcement learning. *arXiv preprint arXiv:1811.09083*, 2018.
- 490 [53] Suraj Nair and Chelsea Finn. Hierarchical foresight: Self-supervised learning of long-horizon  
491 tasks via visual subgoal generation. *arXiv preprint arXiv:1909.05829*, 2019.
- 492 [54] Rundong Wang, Runsheng Yu, Bo An, and Zinovi Rabinovich. I2hrl: Interactive influence-  
493 based hierarchical reinforcement learning. In *Proceedings of the Twenty-Ninth International*  
494 *Conference on International Joint Conferences on Artificial Intelligence*, pages 3131–3138,  
495 2021.

- 496 [55] Christian Daniel, Gerhard Neumann, and Jan Peters. Hierarchical relative entropy policy search.  
497 In *Artificial Intelligence and Statistics*, pages 273–281, 2012.
- 498 [56] Karol Gregor, Danilo Jimenez Rezende, and Daan Wierstra. Variational intrinsic control. *arXiv*  
499 *preprint arXiv:1611.07507*, 2016.
- 500 [57] Tanmay Shankar and Abhinav Gupta. Learning robot skills with temporal variational inference.  
501 In *Proceedings of the 37th International Conference on Machine Learning*. JMLR. org, 2020.
- 502 [58] Archit Sharma, Shixiang Gu, Sergey Levine, Vikash Kumar, and Karol Hausman. Dynamics-  
503 aware unsupervised discovery of skills. In *International Conference on Learning Representa-*  
504 *tions*, 2020.
- 505 [59] Jakob Foerster, Ioannis Alexandros Assael, Nando De Freitas, and Shimon Whiteson. Learning  
506 to communicate with deep multi-agent reinforcement learning. *Advances in neural information*  
507 *processing systems*, 29, 2016.
- 508 [60] Abhishek Das, Théophile Gervet, Joshua Romoff, Dhruv Batra, Devi Parikh, Mike Rabbat, and  
509 Joelle Pineau. Tarmac: Targeted multi-agent communication. In *International Conference on*  
510 *Machine Learning*, pages 1538–1546. PMLR, 2019.
- 511 [61] Sainbayar Sukhbaatar, Rob Fergus, et al. Learning multiagent communication with backpropa-  
512 gation. *Advances in neural information processing systems*, 29, 2016.
- 513 [62] Amanpreet Singh, Tushar Jain, and Sainbayar Sukhbaatar. Learning when to communicate at  
514 scale in multiagent cooperative and competitive tasks. *arXiv preprint arXiv:1812.09755*, 2018.
- 515 [63] Jiechuan Jiang and Zongqing Lu. Learning attentional communication for multi-agent coopera-  
516 tion. *Advances in neural information processing systems*, 31, 2018.
- 517 [64] Daewoo Kim, Sangwoo Moon, David Hostallero, Wan Ju Kang, Taeyoung Lee, Kyunghwan  
518 Son, and Yung Yi. Learning to schedule communication in multi-agent reinforcement learning.  
519 *arXiv preprint arXiv:1902.01554*, 2019.
- 520 [65] Rundong Wang, Xu He, Runsheng Yu, Wei Qiu, Bo An, and Zinovi Rabinovich. Learning  
521 efficient multi-agent communication: An information bottleneck approach. In *International*  
522 *Conference on Machine Learning*, pages 9908–9918. PMLR, 2020.
- 523 [66] Chenghao Li, Chengjie Wu, Tonghan Wang, Jun Yang, Qianchuan Zhao, and Chongjie  
524 Zhang. Celebrating diversity in shared multi-agent reinforcement learning. *arXiv preprint*  
525 *arXiv:2106.02195*, 2021.
- 526 [67] Siyang Wu, Tonghan Wang, Chenghao Li, and Chongjie Zhang. Containerized distributed  
527 value-based multi-agent reinforcement learning. *arXiv preprint arXiv:2110.08169*, 2021.
- 528 [68] Tabish Rashid, Mikayel Samvelyan, Christian Schroeder, Gregory Farquhar, Jakob Foerster,  
529 and Shimon Whiteson. Qmix: Monotonic value function factorisation for deep multi-agent  
530 reinforcement learning. In *International Conference on Machine Learning*, pages 4295–4304,  
531 2018.
- 532 [69] Christian Schroeder de Witt, Tarun Gupta, Denys Makoviichuk, Viktor Makoviychuk, Philip HS  
533 Torr, Mingfei Sun, and Shimon Whiteson. Is independent learning all you need in the StarCraft  
534 multi-agent challenge? *arXiv preprint arXiv:2011.09533*, 2020.
- 535 [70] Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine*  
536 *learning research*, 9(11), 2008.

537 **A Contextual Bandit for Limited Number of Contexts**

---

**Algorithm 2** A contextual bandit algorithm for a small number of contexts

---

- 1: **Initialization:** For each context  $x$ , create an instance  $\text{Exp3}_x$  of algorithm Exp3
  - 2: **for** round **do**
  - 3:   Invoke algorithm  $\text{Exp3}_x$  with  $x = x_t$
  - 4:   Play the action chosen by  $\text{Exp3}_x$
  - 5:   Return reward  $r_t$  to  $\text{Exp3}_x$
  - 6: **end for**
- 

538 **B Proof of Theorem 3.4**

539 **Theorem 3.4.** Consider the Lipschitz contextual bandit problem with contexts in  $[0, 1]$ . The Alg. 1  
 540 yields regret  $\mathbb{E}[R(T)] = O(T^{2/3}(LK \ln T)^{1/3})$ .

541 *Proof.* Let  $S_m$  be the  $\epsilon$ -uniform mesh on  $[0, 1]$ , that is, the set of all points in  $[0, 1]$  that are integer  
 542 multiples of  $\epsilon$ . We take  $\epsilon = 1/(d - 1)$  where the integer  $d$  is the number of points in  $S_m$ , which will  
 543 be adjusted later in the analysis.

544 We apply Alg. 2 to the context space  $S_m$ . Let  $f_{S_m}(x)$  be a mapping from context  $x$  to the closest  
 545 point in  $S_m$ :

$$f_{S_m}(x) = \min_{x' \in S_m} \left( \operatorname{argmin} |x - x'| \right)$$

546 In each round  $t$ , we replace the context  $x_t$  with  $f_{S_m}(x_t)$  and call  $\text{Exp3}_S$ . The regret bound  
 547 will have two components: the regret bound for  $\text{Exp3}_S$  and (a suitable notion of) the discretiza-  
 548 tion error. Formally, let us define the “discretized best response”  $\pi_{S_m}^* : \mathcal{X} \rightarrow \Phi$ :  $\pi_{S_m}^*(x) =$   
 549  $\pi^*(f_{S_m}(x))$  for each context  $x \in \mathcal{X}$ .

550 We define the total reward of an algorithm Alg is  $\text{Reward}(\text{Alg}) = \sum_{t=1}^T r_t$ . Then the regret of  
 551  $\text{Exp3}_S$  and the discretization error are defined as:

$$\begin{aligned} R_S(T) &= \text{Reward}(\pi_S^*) - \text{Reward}(\text{Exp3}_S) \\ \text{DE}(S) &= \text{Reward}(\pi^*) - \text{Reward}(\pi_S^*). \end{aligned}$$

552 It follows that regret is the sum  $R(T) = R_S(T) + \text{DE}(S)$ . We have  $\mathbb{E}[R_S(T)] = \mathcal{O}(\sqrt{TK \log K})$   
 553 from Lemma 3.3, so it remains to upper bound the discretization error and adjust the discretization  
 554 step  $\epsilon$ .

555 For each round  $t$  and the respective context  $x = x_t$ ,  $r(\pi_S^*(x) | f_S(x)) \geq r(\pi^*(x) | f_S(x)) \geq$   
 556  $r(\pi^*(x) | x) - \epsilon L$ . The first inequality is determined by the optimality of  $\pi_S^*$  and the second is  
 557 determined by Lipschitzness. Summing this up over all rounds  $t$ , we obtain  $\mathbb{E}[\text{Reward}(\pi_S^*)] \geq$   
 558  $\text{Reward}[\pi^*] - \epsilon LT$ .

559 Thus, the regret is that

$$\mathbb{E}[R(T)] \leq \epsilon LT + \mathcal{O}\left(\sqrt{\frac{1}{\epsilon} TK \log T}\right) = \mathcal{O}\left(T^{2/3}(LK \log T)^{1/3}\right) \quad (4)$$

560 For the last inequality, we want the two terms of the regret bound has the same asymptotic complexity.  
 561 So when  $\epsilon LT = \text{sqrt} \frac{1}{\epsilon} TK \log T$ , we can get  $\epsilon = \left(\frac{K \log T}{TL^2}\right)^{1/3}$ . So, we choose  $\epsilon = \left(\frac{K \log T}{TL^2}\right)^{1/3}$ .

562  $\square$

563 **C SPC on GRF 11vs11 Full Game**

564 We also conduct experiments on the GRF 11vs11 full game scenario with sparse reward. As shown in  
 565 Fig. 7, SPC achieves about 50% win rate against built-in AI in the target task after training with 200

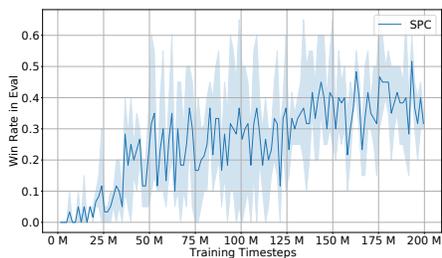


Figure 7: The performance of SPC on the 11v11 scenario.

566 million timesteps. This is non-trivial as this is one of the most challenging benchmarks for MARL  
 567 community, and most current MARL methods struggle to achieve progress without hand-crafted  
 568 engineering.

## 569 D Qualitatively Analysis On Low-Level Skills

570 We demonstrate game statistics under different high-level actions. For example, the times of shooting,  
 571 passing and running actions per game in GRF. These different low-level policies are induced by  
 572 the high-level actions. We evaluate these statistics by fixing one agent’s high-level actions and  
 573 maintaining other agents with SPC. The results in Table 2 are averaged over five runs in the 5vs5  
 574 scenario.

Table 2: Statistics of low-level skills.

	shooting per game	passing per game	running per game
skill 1	7.9 times	0.5 times	2254 time steps
skill 2	2.3 times	26.4 times	2149 time steps
skill 3	1.6 times	3.9 times	2875 time steps

575

## 576 E Comparing Different Teacher Algorithms on GRF Corner-5

577 To further illustrate the effectiveness of the SPC teacher module, we conduct experiments on the  
 578 corner-5 scenario on GRF, where the target task is to control five of the eleven players to obtain  
 579 a goal in the GRF Corner scenario. The experiments are designed to determine whether or not  
 580 the contextual bandit in SPC outperforms alternative curriculum learning methods to schedule the  
 581 number of agents in training. We compare SPC teacher against non-curriculum training (None),  
 582 uniform task sampling (Uniform), a state-of-the-art curriculum learning method (ALP-GMM), and a

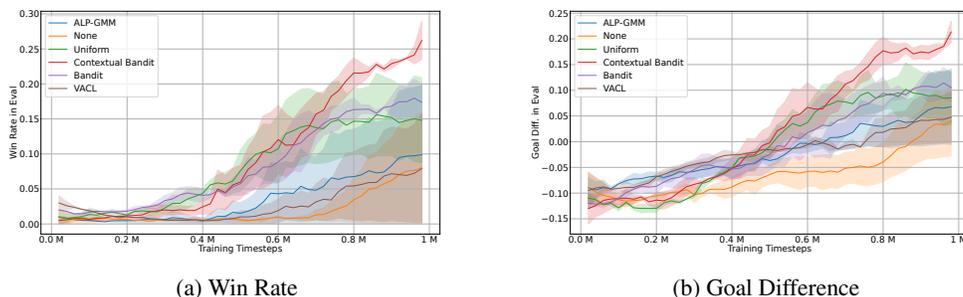


Figure 8: The evaluation performance of various teacher algorithms on the GRF corner-5 scenario.

583 multi-agent curriculum learning method (VACL). The training task space consists of  $n$  agents, where  
 584  $n \in \{1, 3, 5\}$ . All teachers have the same base architecture without transformer architecture and  
 585 HRL. We also investigate the ablation of the RNN-based contexts (see Contextual Bandit and Bandit).  
 586 Fig. 8 shows the benefit of SPC contextual bandit over other ACL methods after training with one  
 587 million timesteps.

## 588 F Implementation Details

589 We use the default implementation of Proximal Policy Optimization (PPO) in Ray RLlib, which  
 590 scales out using multiple workers for experience collection. This allows us to use a large amount of  
 591 rollouts from parallel workers during training to ameliorate high variance and aid exploration. We do  
 592 multiple rollouts in parallel with distributed workers and use parameter sharing for each agent. The  
 593 trainer broadcasts new weights to the workers after their synchronous sampling.

### 594 F.1 Google Research Football

595 We set five tasks for training the GRF 5vs5 scenario, including 5vs5, 3vs3, Pass-Shoot, 3vs1, and  
 596 Empty-Goal. In the Empty-Goal, one agent need to move forward and shoot with an empty goal.  
 597 In Pass-Shoot and 3vs3, two agents are controlled to play against a goalkeeper and three players,  
 598 with different position initialization. In 3vs1, three agents are controlled to play against a center-back  
 599 and a goalkeeper. In 5vs5, four agents are controlled to play against five players. Without loss of  
 600 generality, we initialize all player with fixed positions and roles as center midfielders.

601 We use both MLP and self-attention mechanism for the high-level policy, and use MLP for the  
 602 low-level policy. For high-level policy, the input is first projected to an embedding using two hidden  
 603 layers with 256 units each and ReLU activation, which is then fed into multi-head self-attention  
 604 (8 heads, 64 units each). The output is then projected to the actions and values using another fully  
 605 connected layer with 256 units. For low-level policy, we use MLP with two hidden layers with 256  
 606 units each, i.e., the default configuration of policy network in RLlib.

Table 3: SPC hyper-parameters.

(a) SPC hyper-parameters used in GRF.		(b) SPC hyper-parameters used in MPE.	
Name	Value	Name	Value
Discount rate	0.99	Discount rate	0.99
GAE parameter	1.0	GAE parameter	1.0
KL coefficient	0.2	KL coefficient	0.5
Rollout fragment length	1000	# of SGD iterations	10
Training batch size	100000	Learning rate	1e-4
SGD minibatch size	10000	Entropy coefficient	0.0
# of SGD iterations	60	Clip parameter	0.3
Learning rate	1e-4	Value function clip parameter	10.0
Entropy coefficient	0.0		
Clip parameter	0.3		
Value function clip parameter	10.0		

### 607 F.2 MPE

608 In MPE tasks, agents must cooperate through physical actions to reach a set of landmarks. Agents  
 609 observe the relative positions of other agents and landmarks, and are collectively rewarded based  
 610 on the proximity of any agent to each landmark. In other words, the agents have to cover all of the  
 611 landmarks. Further, the agents are penalized when colliding with each other. The agents need to infer  
 612 the landmark to cover and move there while avoid colliding with other agents.

613 The hyper-parameters of SPC in MPE are shown in Table 3b. In MPE, hyper-parameters such as  
 614 rollout fragment length, training batch size and SGD minibatch size are adjusted according to horizon

615 of the scenarios so that policy are updated after episodes are done. We use the same neural network  
616 architecture as in GRF, but with 128 units for all MLP hidden layers. Other omitted hyper-parameters  
617 follow the default configuration in RLlib PPO implementation.