

## MIT Open Access Articles

*Learning to Teach in Cooperative  
Multiagent Reinforcement Learning*

The MIT Faculty has made this article openly available. **Please share**  
how this access benefits you. Your story matters.

**Citation:** Omidshafie, Shayegan, Kim, Dong-Ki, Liu, Miao, Tesauro, Gerald, Riemer, Matthew et al. 2018. "Learning to Teach in Cooperative Multiagent Reinforcement Learning."

**Persistent URL:** <https://hdl.handle.net/1721.1/137981>

**Version:** Original manuscript: author's manuscript prior to formal peer review

**Terms of use:** Creative Commons Attribution-Noncommercial-Share Alike



# Learning to Teach in Cooperative Multiagent Reinforcement Learning

Shayegan Omidshafiei<sup>1,2</sup>  
shayegan@mit.edu

Dong-Ki Kim<sup>1,2</sup>  
dkkim93@mit.edu

Miao Liu<sup>2,3</sup>  
miao.liu1@ibm.com

Gerald Tesauro<sup>2,3</sup>  
gtesauro@us.ibm.com

Matthew Riemer<sup>2,3</sup>  
mdriemer@us.ibm.com

Christopher Amato<sup>4</sup>  
camato@ccs.neu.edu

Murray Campbell<sup>2,3</sup>  
mcam@us.ibm.com

Jonathan P. How<sup>1,2</sup>  
jhow@mit.edu

<sup>1</sup>LIDS, MIT    <sup>2</sup>MIT-IBM Watson AI Lab    <sup>3</sup>IBM Research    <sup>4</sup>CCIS, Northeastern University

## Abstract

Collective human knowledge has clearly benefited from the fact that innovations by individuals are taught to others through communication. Similar to human social groups, agents in distributed learning systems would likely benefit from communication to share knowledge and teach skills. The problem of teaching to improve agent learning has been investigated by prior works, but these approaches make assumptions that prevent application of teaching to general multiagent problems, or require domain expertise for problems they can apply to. This learning to teach problem has inherent complexities related to measuring long-term impacts of teaching that compound the standard multiagent coordination challenges. In contrast to existing works, this paper presents the first general framework and algorithm for intelligent agents to learn to teach in a multiagent environment. Our algorithm, Learning to Coordinate and Teach Reinforcement (LeCTR), addresses peer-to-peer teaching in cooperative multiagent reinforcement learning. Each agent in our approach learns both when and what to advise, then uses the received advice to improve local learning. Importantly, these roles are not fixed; these agents learn to assume the role of student and/or teacher at the appropriate moments, requesting and providing advice in order to improve teamwide performance and learning. Empirical comparisons against state-of-the-art teaching methods show that our teaching agents not only learn significantly faster, but also learn to coordinate in tasks where existing methods fail.

## Introduction

In social settings, innovations by individuals are taught to others in the population through communication channels (Rogers 2010), which not only improves final performance, but also the effectiveness of the entire learning process (i.e., rate of learning). There exist analogous settings where learning agents interact and adapt behaviors while interacting in a shared environment (e.g., autonomous cars and assistive robots). While any given agent may not be an expert during learning, it may have local knowledge that teammates may be unaware of. Similar to human social groups, these learning agents would likely benefit from communication to share knowledge and teach skills, thereby improving the effectiveness of system-wide learning. It is also desirable for agents in such systems to learn to teach one another, rather than rely on

hand-crafted teaching heuristics created by domain experts. The benefit of learned peer-to-peer teaching is that it can accelerate learning even without relying on the existence of “all-knowing” teachers. Despite these potential advantages, no algorithms exist for learning to teach in multiagent systems.

This paper targets the learning to teach problem in the context of cooperative Multiagent Reinforcement Learning (MARL). Cooperative MARL is a standard framework for settings where agents learn to coordinate in a shared environment. Recent works in cooperative MARL have shown final task performance can be improved by introducing inter-agent communication mechanisms (Sukhbaatar, Fergus, and others 2016; Foerster et al. 2016; Lowe et al. 2017). Agents in these works, however, merely communicate to coordinate in the given task, not to improve overall learning by teaching one another. By contrast, this paper targets a new multiagent paradigm in which agents learn to teach by communicating action advice, thereby improving final performance and accelerating teamwide learning.

The learning to teach in MARL problem has unique inherent complexities that compound the delayed reward, credit assignment, and partial observability issues found in general multiagent problems (Oliehoek and Amato 2016). As such, there are several key issues that must be addressed. First, agents must learn when to teach, what to teach, and how to learn from what is being taught. Second, despite coordinating in a shared environment, agents may be independent/decentralized learners with privacy constraints (e.g., robots from distinct corporations that cannot share full policies), and so must learn how to teach under these constraints. A third issue is that agents must estimate the impact of each piece of advice on their teammate’s learning progress. Delays in the accumulation of knowledge make this credit assignment problem difficult, even in supervised/unsupervised learning (Graves et al. 2017). Nonstationarities due to agent interactions and the temporally-extended nature of MARL compound these difficulties in our setting. These issues are unique to our learning to teach setting and remain largely unaddressed in the literature, despite being of practical importance for future decision-making systems. One of the main reasons for the lack of progress addressing these inherent challenges is the significant increase in the computational complexity of this new teaching/learning paradigm compared to multiagent problems that have previously been considered.

Our paper targets the problem of learning to teach in a multiagent team, which has not been considered before. Each agent in our approach learns both when and what to advise, then uses the received advice to improve local learning. Importantly, these roles are not fixed (see Fig. 1); these agents learn to assume the role of student and/or teacher at appropriate moments, requesting and providing advice to improve teamwide performance and learning. In contrast to prior works, our algorithm supports teaching of heterogeneous teammates and applies to settings where advice exchange incurs a communication cost. Comparisons conducted against state-of-the-art teaching methods show that our teaching agents not only learn significantly faster, but also learn to coordinate in tasks where existing methods fail.

### Background: Cooperative MARL

Our work targets cooperative MARL, where agents execute actions that jointly affect the environment, then receive feedback via local observations and a shared reward. This setting is formalized as a Decentralized Partially Observable Markov Decision Process (Dec-POMDP), defined as  $\langle \mathcal{I}, \mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{R}, \Omega, \mathcal{O}, \gamma \rangle$  (Oliehoek and Amato 2016);  $\mathcal{I}$  is the set of  $n$  agents,  $\mathcal{S}$  is the state space,  $\mathcal{A} = \times_i \mathcal{A}^i$  is the joint action space, and  $\Omega = \times_i \Omega^i$  is the joint observation space.<sup>1</sup> Joint action  $\mathbf{a} = \langle a^1, \dots, a^n \rangle$  causes state  $s \in \mathcal{S}$  to transition to  $s' \in \mathcal{S}$  with probability  $P(s'|s, \mathbf{a}) = \mathcal{T}(s, \mathbf{a}, s')$ . At each timestep  $t$ , joint observation  $\mathbf{o} = \langle o^1, \dots, o^n \rangle$  is observed with probability  $P(\mathbf{o}|s', \mathbf{a}) = \mathcal{O}(\mathbf{o}, s', \mathbf{a})$ . Given its observation history,  $h_t^i = (o_1^i, \dots, o_t^i)$ , agent  $i$  executes actions dictated by its policy  $a^i = \pi^i(h_t^i)$ . The joint policy is denoted by  $\pi = \langle \pi^1, \dots, \pi^n \rangle$  and parameterized by  $\theta$ . It may sometimes be desirable to use a recurrent policy representation (e.g., recurrent neural network) to compute an internal state  $h_t$  that compresses the observation history, or to explicitly compute a belief state (probability distribution over states); with abuse of notation, we use  $h_t$  to refer to all such variations of internal states/observation histories. At each timestep, the team receives reward  $r_t = \mathcal{R}(s_t, \mathbf{a}_t)$ , with the objective being to maximize value,  $V(s; \theta) = \mathbb{E}[\sum_t \gamma^t r_t | s_0 = s]$ , given discount factor  $\gamma \in [0, 1]$ . Let action-value  $Q^i(o^i, a^i; h^i)$  denote agent  $i$ 's expected value for executing action  $a^i$  given a new local observation  $o^i$  and internal state  $h^i$ , and using its policy thereafter. We denote by  $\vec{Q}(o^i; h^i)$  the vector of action-values (for all actions) given new observation  $o^i$ .

### Teaching in Cooperative MARL

This work explores multiagent teaching in a setting where no agent is necessarily an all-knowing expert. This section provides a high-level overview of the motivating scenario. Consider the cooperative MARL setting in Fig. 1, where agents  $i$  and  $j$  learn a joint task (i.e., a Dec-POMDP). In each learning iteration, these agents interact with the environment and collect data used by their respective learning algorithms,  $\mathbb{L}^i$  and  $\mathbb{L}^j$ , to update their policy parameters,  $\theta^i$  and  $\theta^j$ . This is the standard cooperative MARL problem, which we hereafter refer to as  $\mathcal{P}_{\text{Task}}$ : **the task-level learning problem**. For

<sup>1</sup>Superscript  $i$  denotes parameters for the  $i$ -th agent. Refer to the supplementary material for a notation list.

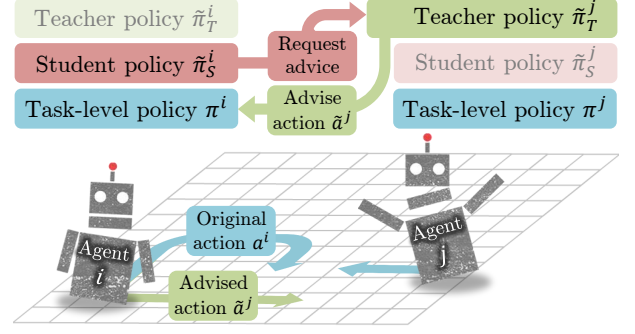


Figure 1: Overview of teaching via action advising in MARL. Each agent learns to execute the task using task-level policy  $\pi$ , to request advice using learned student policy  $\tilde{\pi}_S$ , and to respond with action advice using learned teacher policy  $\tilde{\pi}_T$ . Each agent can assume a student and/or teacher role at any time. In this example, agent  $i$  uses its student policy to request help, agent  $j$  advises action  $\tilde{a}^j$ , which the student executes instead of its originally-intended action  $a^i$ . By learning to transform the local knowledge captured in task-level policies into action advice, the agents can help one another learn.

example, task-level policy  $\pi^i$  is the policy agent  $i$  learns and uses to execute actions in the task. Thus, task-level policies summarize each agent's learned behavioral knowledge.

During task-level learning, it is unlikely for any agent to be an expert. However, each agent may have unique experiences, skill sets, or local knowledge of how to learn effectively in the task. Throughout the learning process, it would likely be useful for agents to advise one another using this local knowledge, in order to improve final performance and accelerate teamwide learning. Moreover, it would be desirable for agents to learn when and what to advise, rather than rely on hand-crafted and domain-specific advising heuristics. Finally, following advising, agents should ideally have learned effective task-level policies that no longer rely on teammate advice at every timestep. We refer to this new problem, which involves agents learning to advise one another to improve joint task-level learning, as  $\mathcal{P}_{\text{Advise}}$ : **the advising-level problem**.

The advising mechanism used in this paper is *action advising*, where agents suggest actions to one another. By learning to appropriately transform local knowledge (i.e., task-level policies) into action advice, teachers can affect students' experiences and their resulting task-level policy updates. Action advising makes few assumptions, in that learners need only use task-level algorithms  $\langle \mathbb{L}^i, \mathbb{L}^j \rangle$  that support off-policy exploration (enabling execution of action advice for policy updates), and that they receive advising-level observations summarizing teammates' learning progress (enabling learning of when/what to advise). Action advising has a good empirical track record (Torrey and Taylor 2013; Taylor et al. 2014; Fachantidis, Taylor, and Vlahavas 2017; da Silva, Glatt, and Costa 2017). However, existing frameworks have key limitations: the majority are designed for single-agent RL and do not consider multiagent learning; their teachers always advise optimal actions to students, mak-

ing decisions about when (not what) to teach; they also use heuristics for advising, rather than training teachers by measuring student learning progress. By contrast, agents in our paper learn to interchangeably assume the role of a student (advice requester) and/or teacher (advice responder), denoted  $S$  and  $T$ , respectively. Each agent learns task-level policy  $\pi$  used to actually perform the task, student policy  $\tilde{\pi}_S$  used to request advice during task-level learning, and teacher policy  $\tilde{\pi}_T$  used to advise a teammate during task-level learning.<sup>2</sup>

Before detailing the algorithm, let us first illustrate the multiagent interactions in this action advising scenario. Consider again Fig. 1, where agents are learning to execute a task (i.e., solving  $\mathcal{P}_{\text{Task}}$ ) while advising one another. While each agent in our framework can assume a student and/or teacher role at any time, Fig. 1 visualizes the case where agent  $i$  is the student and agent  $j$  is the teacher. At a given task-level learning timestep, agent  $i$ 's task-level policy  $\pi^i$  outputs an action ('original action  $a^i$ ' in Fig. 1). However, as the agents are still learning to solve  $\mathcal{P}_{\text{Task}}$ , agent  $i$  may prefer to execute an action that maximizes local learning. Thus, agent  $i$  uses its student policy  $\tilde{\pi}_S^i$  to decide whether to ask teammate  $j$  for advice. If this advice request is made, teammate  $j$  checks its teacher policy  $\tilde{\pi}_T^j$  and task-level policy  $\pi^j$  to decide whether to respond with action advice. Given a response, agent  $i$  then executes the advised action ( $\tilde{a}^i$  in Fig. 1) as opposed to its originally-intended action ( $a^i$  in Fig. 1). This results in a local experience that agent  $i$  uses to update its task-level policy. A reciprocal process occurs when the agents' roles are reversed. The benefit of advising is that agents can learn to use local knowledge to improve teamwide learning.

Similar to recent works that model the multiagent learning process (Hadfield-Menell et al. 2016; Foerster et al. 2018), we focus on the pairwise (two agent) case, targeting the issues of when/what to advise, then detail extensions to  $n$  agents. Even in the pairwise case, there exist issues unique to our learning to teach paradigm. First, note that the objectives of  $\mathcal{P}_{\text{Task}}$  and  $\tilde{\mathcal{P}}_{\text{Advise}}$  are distinct. Task-level problem,  $\mathcal{P}_{\text{Task}}$ , has a standard MARL objective of agents learning to coordinate to maximize final performance in the task. Learning to advise ( $\tilde{\mathcal{P}}_{\text{Advise}}$ ), however, is a higher-level problem, where agents learn to influence teammates' task-level learning by advising them. However,  $\mathcal{P}_{\text{Task}}$  and  $\tilde{\mathcal{P}}_{\text{Advise}}$  are also coupled, as advising influences the task-level policies learned. Agents in our problem must learn to advise despite the nonstationarities due to changing task-level policies, which are also a function of algorithms  $\langle \mathbb{L}^i, \mathbb{L}^j \rangle$  and policy parameterizations  $\langle \theta^i, \theta^j \rangle$ .

Learning to teach is also distinct from prior works that involve agents learning to communicate (Sukhbaatar, Ferguson, and others 2016; Foerster et al. 2016; Lowe et al. 2017). These works focus on agents communicating in order to coordinate in a task. By contrast, our problem focuses on agents learning how advising affects the underlying task-level learning process, then using this knowledge to accelerate learning even when agents are non-experts. Thus, the objectives of communication-based multiagent papers are disparate from ours, and the two approaches may even be combined.

<sup>2</sup>Tilde accents (e.g.,  $\tilde{\pi}$ ) denote advising-level properties.

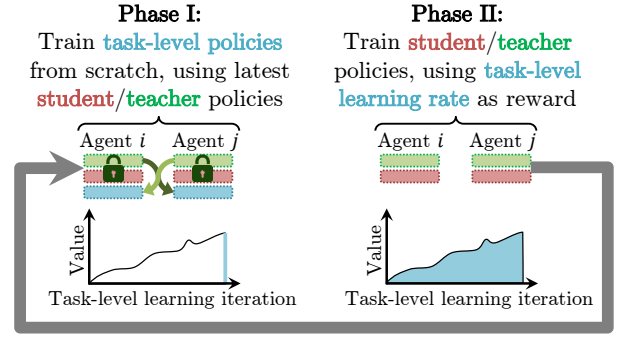


Figure 2: LeCTR consists of two iterated phases: task-level learning (Phase I), and advising-level learning (Phase II). In Phase II, advising policies are trained using rewards correlated to task-level learning (see Table 1). Task-level, student, and teacher policy colors above follows convention of Fig. 1.

## LeCTR: Algorithm for Learning to Coordinate and Teach Reinforcement

This section introduces our learning to teach approach, details how issues specific to our problem setting are resolved, and summarizes overall training protocol. Pseudocode is presented in the supplementary material due to limited space.

**Overview** Our algorithm, Learning to Coordinate and Teach Reinforcement (LeCTR), solves advising-level problem  $\tilde{\mathcal{P}}_{\text{Advise}}$ . The objective is to learn advising policies that augment agents' task-level algorithms  $\langle \mathbb{L}^i, \mathbb{L}^j \rangle$  to accelerate solving of  $\mathcal{P}_{\text{Task}}$ . Our approach involves 2 phases (see Fig. 2):

- Phase I: agents learn  $\mathcal{P}_{\text{Task}}$  from scratch using blackbox learning algorithms  $\langle \mathbb{L}^i, \mathbb{L}^j \rangle$  and latest advising policies.
- Phase II: advising policies are updated using advising-level rewards correlated to teammates' task-level learning.

No restrictions are placed on agents' task-level algorithms (i.e., they can be heterogeneous). Iteration of Phases I and II enables training of increasingly capable advising policies.

**Advising Policy Inputs & Outputs** LeCTR learns student policies  $\langle \tilde{\pi}_S^i, \tilde{\pi}_S^j \rangle$  and teacher policies  $\langle \tilde{\pi}_T^i, \tilde{\pi}_T^j \rangle$  for agents  $i$  and  $j$ , constituting a jointly-initiated advising approach that learns when to request advice and when/what to advise. It is often infeasible to learn high-level policies that directly map task-level policy parameters  $\langle \theta^i, \theta^j \rangle$  (i.e., local knowledge) to advising decisions: the agents may be independent/decentralized learners and the cost of communicating task-level policy parameters may be high; sharing policy parameters may be undesirable due to privacy concerns; and learning advising policies over the task-level policy parameter space may be infeasible (e.g., if the latter policies involve millions of parameters). Instead, each LeCTR agent learns advising policies over advising-level observations  $\tilde{o}$ . As detailed below, these observations are selected to provide information about agents' task-level state and knowledge in a more compact manner than full policy parameters  $\langle \theta^i, \theta^j \rangle$ .

Each LeCTR agent can be a student, teacher, or both simultaneously (i.e., request advice for its own state, while advising a teammate in a different state). For clarity, we

detail advising protocols when agents  $i$  and  $j$  are student and teacher, respectively (see Fig. 1). LeCTR uses distinct advising-level observations for student and teacher policies. Student policy  $\tilde{\pi}_S^i$  for agent  $i$  decides when to request advice using advising-level observation  $\tilde{o}_S^i = \langle o^i, \tilde{Q}^i(o^i; h^i) \rangle$ , where  $o^i$  and  $\tilde{Q}^i(o^i; h^i)$  are the agent’s task-level observation and action-value vectors, respectively. Through  $\tilde{o}_S^i$ , agent  $i$  observes a measure of its local task-level observation and policy state. Thus, agent  $i$ ’s student-perspective action is  $\tilde{a}_S^i = \tilde{\pi}_S^i(\tilde{o}_S^i) \in \{\text{request advice, do not request advice}\}$ .

Similarly, agent  $j$ ’s teacher policy  $\tilde{\pi}_T^j$  uses advising-level observation  $\tilde{o}_T^j = \langle o^j, \tilde{Q}^j(o^j; h^j), \tilde{Q}^j(o^i; h^i) \rangle$  to decide when/what to advise.  $\tilde{o}_T^j$  provides teacher agent  $j$  with a measure of student  $i$ ’s task-level state/knowledge (via  $o^i$  and  $\tilde{Q}^i(o^i; h^i)$ ) and of its own task-level knowledge given the student’s context (via  $\tilde{Q}^j(o^i; h^i)$ ). Using  $\tilde{\pi}_T^j$ , teacher  $j$  decides what to advise: either an action from student  $i$ ’s action space,  $\mathcal{A}^i$ , or a special no-advice action  $\tilde{a}_\emptyset$ . Thus, the teacher-perspective action for agent  $j$  is  $\tilde{a}_T^j = \tilde{\pi}_T^j(\tilde{o}_T^j) \in \mathcal{A}^i \cup \{\tilde{a}_\emptyset\}$ .

Given no advice, student  $i$  executes originally-intended action  $a^i$ . However, given advice  $\tilde{a}_T^j$ , student  $i$  executes action  $\beta^i(\tilde{a}_T^j)$ , where  $\beta^i(\cdot)$  is a local behavioral policy not known by  $j$ . The assumption of local behavioral policies increases the generality of LeCTR, as students may locally transform advised actions before execution.

Following advice execution, agents collect task-level experiences and update their respective task-level policies. A key feature is that LeCTR agents learn what to advise by training  $\langle \tilde{\pi}_T^i, \tilde{\pi}_T^j \rangle$ , rather than always advising actions they would have taken in students’ states. These agents may learn to advise exploratory actions or even decline to advise if they estimate that such advice will improve teammate learning.

**Rewarding Advising Policies** Recall in Phase II of LeCTR, advising policies are trained to maximize advising-level rewards that should, ideally, reflect the objective of accelerating task-level learning. Without loss of generality, we focus again on the case where agents  $i$  and  $j$  assume student and teacher roles, respectively, to detail these rewards. Since student policy  $\tilde{\pi}_S^i$  and teacher policy  $\tilde{\pi}_T^j$  must coordinate to help student  $i$  learn, they receive identical advising-level rewards,  $\tilde{r}_S^i = \tilde{r}_T^j$ . The remaining issue is to identify advising-level rewards that reflect learning progress.

**Remark 1.** *Earning task-level rewards by executing advised actions may not imply actual learning. Thus, rewarding advising-level policies with the task-level reward,  $r$ , received after advice execution can lead to poor advising policies.*

We evaluate many choices of advising-level rewards, which are summarized and described in Table 1. The unifying intuition is that each reward type corresponds to a different measure of the advised agent’s task-level learning, which occurs after executing an advised action. Readers are referred to the supplementary material for more details.

Note that at any time, task-level action  $a^i$  executed by agent  $i$  may either be selected by its local task-level policy, or by a teammate  $j$  via advising. In Phase II, pair  $\langle \tilde{\pi}_S^i, \tilde{\pi}_T^j \rangle$  is rewarded only if advising occurs (with zero advising reward

otherwise). Analogous advising-level rewards apply for the reverse student-teacher pairing  $j-i$ , where  $\tilde{r}_S^j = \tilde{r}_T^i$ . During Phase II, we train all advising-level policies using a joint advising-level reward  $\tilde{r} = \tilde{r}_T^i + \tilde{r}_T^j$  to induce cooperation.

Advising-level rewards are only used during advising-level training, and are computed using either information already available to agents or only require exchange of scalar values (rather than full policy parameters). It is sometimes desirable to consider advising under communication constraints, which can be done by deducting a communication cost  $c$  from these advising-level rewards for each piece of advice exchanged.

**Training Protocol** Recall LeCTR’s two phases are iterated to enable training of increasingly capable advising policies. In Phase I, task-level learning is conducted using agents’ blackbox learning algorithms and latest advising policies. At the task-level, agents may be independent learners with distinct algorithms. Advising policies are executed in a decentralized fashion, but their training in Phase II is centralized. Our advising policies are trained using the multiagent actor-critic approach of Lowe et al. (2017). Let joint advising-level observations, advising-level actions, and advising-level policies (i.e., ‘actors’) be, respectively, denoted by  $\tilde{o} = \langle \tilde{o}_S^i, \tilde{o}_S^j, \tilde{o}_T^i, \tilde{o}_T^j \rangle$ ,  $\tilde{a} = \langle \tilde{a}_S^i, \tilde{a}_S^j, \tilde{a}_T^i, \tilde{a}_T^j \rangle$ , and  $\tilde{\pi} = \langle \tilde{\pi}_S^i, \tilde{\pi}_S^j, \tilde{\pi}_T^i, \tilde{\pi}_T^j \rangle$ , with  $\theta$  parameterizing  $\tilde{\pi}$ . To induce  $\tilde{\pi}$  to learn to teach both agents  $i$  and  $j$ , we use a centralized action-value function (i.e., ‘critic’) with advising-level reward  $\tilde{r} = \tilde{r}_T^i + \tilde{r}_T^j$ . Critic  $\tilde{Q}(\tilde{o}, \tilde{a}; \tilde{\theta})$  is trained by minimizing loss,

$$\mathcal{L}(\tilde{\theta}) = \mathbb{E}_{\tilde{o}, \tilde{a}, \tilde{r}, \tilde{o}' \sim \tilde{\mathcal{M}}} [(\tilde{r} + \gamma \tilde{Q}(\tilde{o}', \tilde{a}'; \tilde{\theta}) - \tilde{Q}(\tilde{o}, \tilde{a}; \tilde{\theta}))^2] \Big|_{\tilde{a}' = \tilde{\pi}(\tilde{o}'),} \quad (1)$$

where  $\tilde{a}' = \tilde{\pi}(\tilde{o}')$  are next advising actions computed using the advising policies, and  $\tilde{\mathcal{M}}$  denotes advising-level replay buffer (Mnih et al. 2015). The policy gradient theorem (Sutton et al. 2000) is invoked on objective  $J(\tilde{\theta}) = \mathbb{E}[\sum_{t=k}^T \gamma^{t-k} \tilde{r}_t]$  to update advising policies using gradients,

$$\begin{aligned} \nabla_{\tilde{\theta}} J(\tilde{\theta}) &= \mathbb{E}_{\tilde{o}, \tilde{a} \sim \tilde{\mathcal{M}}} [\nabla_{\tilde{\theta}} \log \tilde{\pi}(\tilde{a} | \tilde{o}) \tilde{Q}(\tilde{o}, \tilde{a}; \tilde{\theta})] \\ &= \mathbb{E}_{\tilde{o}, \tilde{a} \sim \tilde{\mathcal{M}}} \left[ \sum_{\substack{\alpha \in \{i, j\} \\ \rho \in \{S, T\}}} \nabla_{\tilde{\theta}} \log \tilde{\pi}_\rho^\alpha(\tilde{a}_\rho^\alpha | \tilde{o}_\rho^\alpha) \nabla_{\tilde{a}_\rho^\alpha} \tilde{Q}(\tilde{o}, \tilde{a}; \tilde{\theta}) \right], \end{aligned} \quad (2)$$

where  $\tilde{\pi}_\rho^\alpha$  is agent  $\alpha$ ’s policy in role  $\rho$ .

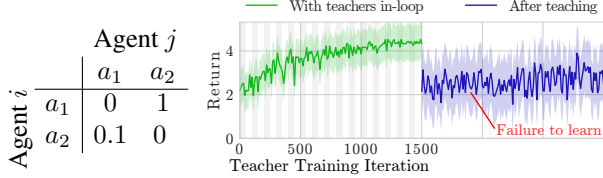
During training, the advising feedback nonstationarities mentioned earlier are handled as follows: in Phase I, task-level policies are trained online (i.e., no replay memory is used so impact of advice on task-level policies is immediately observed by agents); in Phase II, centralized advising-level learning reduces nonstationarities due to teammate learning, and reservoir sampling is used to further reduce advising reward nonstationarities (see supplementary material for details). Our overall approach stabilizes advising-level learning.

**Advising  $n$  Agents** In the  $n$  agent case, students must also decide how to fuse advice from multiple teachers. This is a complex problem requiring full investigation in future work; feasible ideas include using majority voting for advice fusion (as in da Silva, Glatt, and Costa (2017)), or asking a specific agent for advice by learning a ‘teacher score’ modulated based on teacher knowledge/previous teaching experiences.



Table 1: Summary of rewards used to train advising policies. Rewards shown are for the case where agent  $i$  is student and agent  $j$  teacher (i.e., flip the indices for the reverse case). Each reward corresponds to a different measure of task-level learning after the student executes an action advice and uses it to update its task-level policy. Refer to the supplementary material for more details.

Advising Reward Name	Description	Reward Value $\tilde{r}_T^j = \tilde{r}_S^i$
JVG: Joint Value Gain	Task-level value $V(s; \theta)$ improvement after learning	$V(s; \theta_{t+1}) - V(s; \theta_t)$
QTR: Q-Teaching Reward	Teacher’s estimate of best vs. intended student action	$\max_a Q_T(o^i, a; h^i) - Q_T(o^i, a^i; h^i)$
LG: Loss Gain	Student’s task-level loss $\mathcal{L}(\theta^i)$ reduction	$\mathcal{L}(\theta_t^i) - \mathcal{L}(\theta_{t+1}^i)$
LGG: Loss Gradient Gain	Student’s task-level policy gradient magnitude	$\ \nabla_{\theta^i} \mathcal{L}(\theta^i)\ _2^2$
TDG: TD Gain	Student’s temporal difference (TD) error $\delta^i$ reduction	$ \delta_t^i  -  \delta_{t+1}^i $
VEG: Value Estimation Gain	Student’s value estimate $\hat{V}(\theta^i)$ gain above threshold $\tau$	$\mathbb{1}(\hat{V}(\theta^i) > \tau)$



(a) Repeated game payoffs. Each agent reward choice  $\tilde{r}_T = r$ . LeCTR Phase I and has 2 actions ( $a_1, a_2$ ). II iterations are shown as background bands.

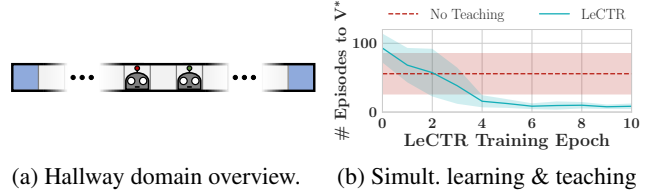
Figure 3: Repeated game. (b) shows a counterexample where using  $\tilde{r}_T = r$  yields poor advising, as teachers learn to advise actions that maximize reward (left half, green), but do not actually improve student task-level learning (right half, blue).

## Evaluation

We conduct empirical evaluations on a sequence of increasingly challenging domains involving two agents. In the ‘Repeated’ game domain, agents coordinate to maximize the payoffs in Fig. 3a over 5 timesteps. In ‘Hallway’ (see Fig. 4a), agents only observe their own positions and receive +1 reward if they reach opposite goal states; task-level actions are ‘move left/right’, states are agents’ joint grid positions. The higher-dimensional ‘Room’ game (see Fig. 5a) has the same state/observation/reward structure, but 4 actions (‘move up/right/down/left’). Recall student-perspective advising-level actions are to ‘ask’ or ‘not ask’ for advice. Teacher-perspective actions are to advise an action from the teammate’s task-level action space, or to decline advising.

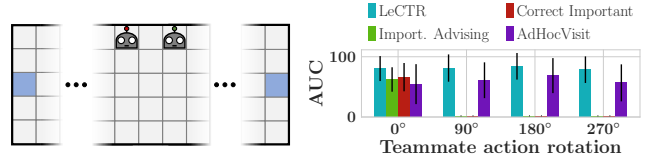
For the Repeated, Hallway, and Room games, respectively, each iteration of LeCTR Phase I consists of 50, 100, and 150 task-level learning iterations. Our task-level agents are independent Q-learners with tabular policies for the Repeated game and tile-coded policies (Sutton and Barto 1998) for the other games. Advising policies are neural networks with internal rectified linear unit activations. Refer to the supplementary material for hyperparameters. The advising-level learning nature of our problem makes these domains challenging, despite their visual simplicity; their complexity is comparable to domains tested in recent MARL works that learn over multiagent learning processes (Foerster et al. 2018), which also consider two agent repeated/gridworld games.

**Counterexample demonstrating Remark 1** Fig. 3b



(a) Hallway domain overview. (b) Simult. learning & teaching

Figure 4: Hallway game. (a) Agents receive +1 reward by navigating to opposite states in 17-grid hallway. (b) LeCTR accelerates learning & teaching compared to no-teaching.



(a) Room domain overview. (b) Teaching heterogeneous agents.

Figure 5: Room game. (a) Agents receive +1 reward by navigating to opposite goals in 17x5 grid. (b) LeCTR outperforms prior approaches when agents are heterogeneous.

shows results given a poor choice of advising-level reward,  $\tilde{r}_T = r$ , in the Repeated game. The left plot (in green) shows task-level return received due to both local policy actions *and* advised actions, which increases as teachers learn. However, in the right plot (blue) we evaluate how well task-level policies perform by themselves, after they have been trained using the final advising-level policies. The poor performance of the resulting task-level policies indicates that advising policies learned to maximize their own rewards  $\tilde{r}_T = r$  by always advising optimal actions to students, thereby disregarding whether task-level policies actually learn. No exploratory actions are advised, causing poor task-level performance after advising. This counterexample demonstrates that advising-level rewards that reflect student learning progress, rather than task-level reward  $r$ , are critical for useful advising.

**Comparisons to existing teaching approaches** Table 2 shows extensive comparisons of existing heuristics-based teaching approaches, no teaching (independent Q-learning), and LeCTR with all advising rewards introduced in Table 1. We use the VEG advising-level reward in the final version of

Table 2:  $\bar{V}$  and Area under the Curve (AUC) for teaching algorithms. Best results in bold (computed via a  $t$ -test with  $p < 0.05$ ). Independent Q-learning correspond to the no-teaching case. <sup>†</sup>Final version LeCTR uses the VEG advising-level reward.

Algorithm	Repeated Game		Hallway Game		Room Game	
	$\bar{V}$	AUC	$\bar{V}$	AUC	$\bar{V}$	AUC
Independent Q-learning (No Teaching)	2.75 $\pm$ 2.12	272 $\pm$ 210	0.56 $\pm$ 0.35	36 $\pm$ 24	0.42 $\pm$ 0.33	22 $\pm$ 25
Ask Important (Amir et al. 2016)	1.74 $\pm$ 1.89	178 $\pm$ 181	0.53 $\pm$ 0.36	39 $\pm$ 27	0.00 $\pm$ 0.00	0 $\pm$ 0
Ask Uncertain (Clouse 1996)	1.74 $\pm$ 1.89	170 $\pm$ 184	0.00 $\pm$ 0.00	0 $\pm$ 0	0.00 $\pm$ 0.00	0 $\pm$ 0
Early Advising (Torrey and Taylor 2013)	0.45 $\pm$ 0.00	45 $\pm$ 1	0.00 $\pm$ 0.00	0 $\pm$ 0	0.00 $\pm$ 0.00	0 $\pm$ 0
Import. Advising (Torrey and Taylor 2013)	0.45 $\pm$ 0.00	45 $\pm$ 1	0.67 $\pm$ 0.07	39 $\pm$ 17	0.57 $\pm$ 0.03	48 $\pm$ 8
Early Correcting (Amir et al. 2016)	0.45 $\pm$ 0.00	45 $\pm$ 1	0.00 $\pm$ 0.00	0 $\pm$ 0	0.00 $\pm$ 0.00	0 $\pm$ 0
Correct Important (Torrey and Taylor 2013)	0.45 $\pm$ 0.00	45 $\pm$ 1	0.67 $\pm$ 0.07	39 $\pm$ 16	0.56 $\pm$ 0.00	51 $\pm$ 7
AdHocVisit (da Silva, Glatt, and Costa 2017)	2.49 $\pm$ 2.04	244 $\pm$ 199	0.57 $\pm$ 0.34	38 $\pm$ 24	0.43 $\pm$ 0.33	22 $\pm$ 26
AdHocTD (da Silva, Glatt, and Costa 2017)	1.88 $\pm$ 1.94	184 $\pm$ 189	0.49 $\pm$ 0.37	26 $\pm$ 24	0.39 $\pm$ 0.33	26 $\pm$ 29
LeCTR (with JVG)	<b>4.16<math>\pm</math>1.17</b>	<b>405<math>\pm</math>114</b>	0.25 $\pm$ 0.37	21 $\pm$ 31	0.11 $\pm$ 0.27	6 $\pm$ 21
LeCTR (with QTR)	<b>4.52<math>\pm</math>0.00</b>	<b>443<math>\pm</math>3</b>	0.21 $\pm$ 0.35	12 $\pm$ 22	0.20 $\pm$ 0.32	11 $\pm$ 22
LeCTR (with TDG)	3.36 $\pm$ 1.92	340 $\pm$ 138	0.19 $\pm$ 0.34	15 $\pm$ 25	0.26 $\pm$ 0.34	25 $\pm$ 32
LeCTR (with LG)	<b>3.88<math>\pm</math>1.51</b>	375 $\pm$ 132	0.13 $\pm$ 0.29	13 $\pm$ 22	0.37 $\pm$ 0.34	30 $\pm$ 32
LeCTR (with LGG)	<b>4.41<math>\pm</math>0.69</b>	<b>430<math>\pm</math>53</b>	0.22 $\pm$ 0.35	27 $\pm$ 29	0.56 $\pm$ 0.27	56 $\pm$ 23
LeCTR <sup>†</sup>	<b>4.52<math>\pm</math>0.00</b>	<b>443<math>\pm</math>3</b>	<b>0.77<math>\pm</math>0.00</b>	<b>71<math>\pm</math>3</b>	<b>0.68<math>\pm</math>0.07</b>	<b>79<math>\pm</math>16</b>

our LeCTR algorithm, but show all advising reward results for completeness. We report both final task-level performance after teaching,  $\bar{V}$ , and also area under the task-level learning curve (AUC) as a measure of rate of learning; higher values are better for both. Single-agent approaches requiring an expert teacher are extended to the MARL setting by using teammates’ policies (pre-trained to expert level) as each agent’s teacher. In the Repeated game, LeCTR attains best performance in terms of final value and rate of learning (AUC). Existing approaches always advise the teaching agent’s optimal action to its teammate, resulting in suboptimal returns. In the Hallway and Room games, approaches that tend to over-advise (e.g., Ask Uncertain, Early Advising, and Early Correcting) perform poorly. AdHocVisit and AdHocTD fare better, as their probabilistic nature permits agents to take exploratory actions and sometimes learn optimal policies. Importance Advising and Correct Important heuristics lead agents to suboptimal (distant) goals in Hallway and Room, yet attain positive value due to domain symmetries.

LeCTR outperforms all approaches when using the VEG advising-level reward (Table 2). While the JVG advising-level reward seems an intuitive measure of learning progress due to directly measuring task-level performance, its high variance in situations where the task-level value is sensitive to policy initialization sometimes destabilizes training. JVG is also expensive to compute, requiring game rollouts after each advice exchange. LG and TDG perform poorly due to the high variance of task-level losses used to compute them. We hypothesize that VEG performs best as its thresholded binary advising-level reward filters the underlying noisy task-level losses for teachers. A similar result is reported in recent work on teaching of supervised learners, where threshold-based advising-level rewards have good empirical performance (Fan et al. 2018). Fig. 4b shows improvement of LeCTR’s advising policies due to training, measured by the number of task-level

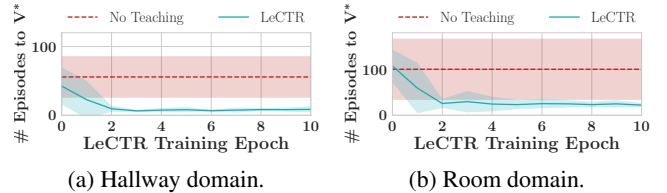
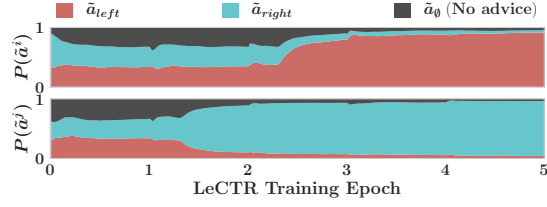


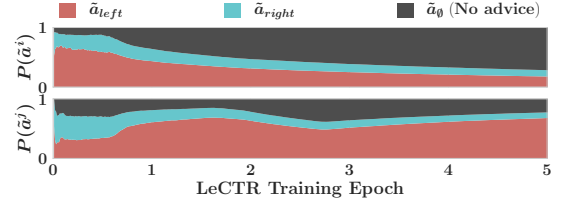
Figure 6: LeCTR accelerates multiagent transfer learning.

episodes needed to converge to the max value reached,  $V^*$ . LeCTR outperforms the rate of learning for the no-teaching case, stabilizing after roughly 4-6 training epochs.

**Teaching for transfer learning** Learning to teach can also be applied to multiagent transfer learning. We first pre-train task-level policies in the Hallway/Room tasks (denote these  $\mathbb{T}_1$ ), flip agents’ initial positions, then train agents to use teammates’  $\mathbb{T}_1$  task-level policies to accelerate learning in flipped task  $\mathbb{T}_2$ . Results for Hallway and Room are shown in Figs. 6a and 6b, respectively, where advising accelerates rate of learning using prior task knowledge. Next, we test *transferability of advising policies* themselves (i.e., use advising policies trained for one task to accelerate learning in a brand new, but related, task). We fix (no longer train) advising policies from the above transfer learning test. We then consider 2 variants of Room: one with the domain (including initial agent positions) flipped vertically ( $\mathbb{T}_3$ ), and one flipped vertically *and* horizontally ( $\mathbb{T}_4$ ). We evaluate the fixed advising policies (trained to transfer from  $\mathbb{T}_1 \rightarrow \mathbb{T}_2$ ) on transfer from  $\mathbb{T}_3 \rightarrow \mathbb{T}_4$ . Learning without advising on  $\mathbb{T}_4$  yields AUC  $24 \pm 26$ , while using the fixed advising policy for transfer  $\mathbb{T}_3 \rightarrow \mathbb{T}_4$  attains AUC  $68 \pm 17$ . Thus, learning is accelerated even when using pre-trained advising policies. While transfer learning typically involves more significant differences in tasks, these preliminary results motivate future



(a) No communication cost,  $c = 0$ , agents advise opposite actions.



(b) With  $c = 0.5$ , one agent leads & advises actions opposite its own.

Figure 7: Hallway game, impact of communication cost  $c$  on advising policy behaviors. First and second rows show probabilities of action advice,  $P(\tilde{a}^i)$  and  $P(\tilde{a}^j)$ , for agents  $i$  and  $j$ , respectively, as their advising policies are trained using LeCTR.

work on applications of advising for MARL transfer learning.

**Advising heterogeneous teammates** We consider heterogeneous variants of the Room game where one agent,  $j$ , uses rotated versions of its teammate  $i$ 's action space; e.g., for rotation  $90^\circ$ , agent  $i$ 's action indices correspond to (up/right/down/left), while  $j$ 's to (left/up/right/down). Comparisons of LeCTR and the best-performing existing methods are shown in Fig. 5b for all rotations. Prior approaches (Importance Advising and Correct Important) work well for homogeneous actions ( $0^\circ$  rotation). However, they attain 0 AUC for heterogeneous cases, as agents always advise action indices corresponding to their local action spaces, leading teammates to no-reward regions. AdHocVisit works reasonably well for all rotations, by sometimes permitting agents to explore. LeCTR attains highest AUC for all rotations.

**Effect of communication cost on advice exchange** We evaluate impact of communication cost on advising by deducting cost  $c$  from advising rewards for each piece of advice exchanged. Fig. 7 shows a comparison of action advice probabilities for communication costs  $c = 0$  and  $c = 0.5$  in the Hallway game. With no cost ( $c = 0$  in Fig. 7a), agents learn to advise each other opposite actions ( $\tilde{a}_{\text{left}}$  and  $\tilde{a}_{\text{right}}$ , respectively) in addition to exploratory actions. As LeCTR's VEG advising-level rewards are binary (0 or 1), two-way advising nullifies positive advising-level rewards, penalizing excessive advising. Thus, when  $c = 0.5$  (Fig. 7b), advising becomes unidirectional: one agent advises opposite exploratory actions of its own, while its teammate tends not to advise.

## Related Work

Effective diffusion of knowledge has been studied in many fields, including inverse reinforcement learning (Ng and Russell 2000), apprenticeship learning (Abbeel and Ng 2004), and learning from demonstration (Argall et al. 2009), wherein students discern and emulate key demonstrated behaviors. Works on curriculum learning (Bengio et al. 2009) are also related, particularly automated curriculum learning (Graves et al. 2017). Though Graves et al. focus on single student supervised/unsupervised learning, they highlight interesting measures of learning progress also used here. Several works meta-learn active learning policies for supervised learning (Bachman, Sordoni, and Trischler 2017; Fang, Li, and Cohn 2017; Pang, Dong, and Hospedales 2018; Fan et al. 2018). Our work also uses advising-level meta-learning, but in the regime of MARL, where agents must

learn to advise teammates without destabilizing coordination.

In action advising, a student executes actions suggested by a teacher, who is typically an expert always advising the optimal action (Torrey and Taylor 2013). These works typically use state importance value  $I(s, \hat{a}) = \max_a Q(s, a) - Q(s, \hat{a})$  to decide when to advise, estimating the performance difference between the student's best action versus intended/worst-case action  $\hat{a}$ . In student-initiated approaches such as Ask Uncertain (Clouse 1996) and Ask Important (Amir et al. 2016), the student decides when to request advice using heuristics based on  $I(s, \hat{a})$ . In teacher-initiated approaches such as Importance Advising (Torrey and Taylor 2013), Early Correcting (Amir et al. 2016), and Correct Important (Torrey and Taylor 2013), the teacher decides when to advise by comparing student policy  $\pi_S$  to expert policy  $\pi_T$ . Q-Teaching (Fachantidis, Taylor, and Vlahavas 2017) learns when to advise by rewarding the teacher  $I(s, \hat{a})$  when it advises. See the supplementary material for details of these approaches.

While most works on information transfer target single-agent settings, several exist for MARL. These include imitation learning of expert demonstrations (Le et al. 2017), cooperative inverse reinforcement learning with a human and robot (Hadfield-Menell et al. 2016), and transfer to parallel learners in tasks with similar value functions (Taylor et al. 2013). To our knowledge, AdHocVisit and AdHocTD (da Silva, Glatt, and Costa 2017) are the only action advising methods that do not assume expert teachers; teaching agents always advise the action they would have locally taken in the student's state, using state visit counts as a heuristic to decide when to exchange advice. Wang et al. (2018) uses da Silva, Glatt, and Costa's teaching algorithm with minor changes.

## Contribution

This work introduced a new paradigm for learning to teach in cooperative MARL settings. Our algorithm, LeCTR, uses agents' task-level learning progress as advising policy feedback, training advisors that improve the rate of learning without harming final performance. Unlike prior works (Torrey and Taylor 2013; Taylor et al. 2014; Zimmer, Viappiani, and Weng 2014), our approach avoids hand-crafted advising policies and does not assume expert teachers. Due to the many complexities involved, we focused on the pairwise problem, targeting the issues of *when* and *what* to teach. A natural avenue for future work is to investigate the  $n$ -agent setting, extending the ideas presented here where appropriate.



## Acknowledgements

Research funded by IBM (as part of the MIT-IBM Watson AI Lab initiative) and a Kwanjeong Educational Foundation Fellowship. The authors thank Dr. Kasra Khosoussi for fruitful discussions early in the paper development process.

## References

- Abbeel, P., and Ng, A. Y. 2004. Apprenticeship learning via inverse reinforcement learning. In *Proceedings of the twenty-first international conference on Machine learning*, 1. ACM.
- Amir, O.; Kamar, E.; Kolobov, A.; and Grosz, B. J. 2016. Interactive teaching strategies for agent training. *International Joint Conferences on Artificial Intelligence*.
- Argall, B. D.; Chernova, S.; Veloso, M.; and Browning, B. 2009. A survey of robot learning from demonstration. *Robotics and autonomous systems* 57(5):469–483.
- Bachman, P.; Sordoni, A.; and Trischler, A. 2017. Learning algorithms for active learning. In *International Conference on Machine Learning*, 301–310.
- Bengio, Y.; Louradour, J.; Collobert, R.; and Weston, J. 2009. Curriculum learning. In *Proceedings of the 26th annual international conference on machine learning*, 41–48. ACM.
- Clouse, J. A. 1996. On integrating apprentice learning and reinforcement learning.
- da Silva, F. L.; Glatt, R.; and Costa, A. H. R. 2017. Simultaneously learning and advising in multiagent reinforcement learning. In *Proceedings of the 16th Conference on Autonomous Agents and MultiAgent Systems*, 1100–1108. International Foundation for Autonomous Agents and Multiagent Systems.
- Fachantidis, A.; Taylor, M. E.; and Vlahavas, I. 2017. Learning to teach reinforcement learning agents. *Machine Learning and Knowledge Extraction* 1(1):2.
- Fan, Y.; Tian, F.; Qin, T.; Li, X.-Y.; and Liu, T.-Y. 2018. Learning to teach. In *International Conference on Learning Representations*.
- Fang, M.; Li, Y.; and Cohn, T. 2017. Learning how to active learn: A deep reinforcement learning approach. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, 595–605.
- Foerster, J.; Assael, I. A.; de Freitas, N.; and Whiteson, S. 2016. Learning to communicate with deep multi-agent reinforcement learning. In *Advances in Neural Information Processing Systems*, 2137–2145.
- Foerster, J. N.; Chen, R. Y.; Al-Shedivat, M.; Whiteson, S.; Abbeel, P.; and Mordatch, I. 2018. Learning with opponent-learning awareness. In *Proceedings of the 17th Conference on Autonomous Agents and MultiAgent Systems*. International Foundation for Autonomous Agents and Multiagent Systems.
- Graves, A.; Bellemare, M. G.; Menick, J.; Munos, R.; and Kavukcuoglu, K. 2017. Automated curriculum learning for neural networks. In *International Conference on Machine Learning*, 1311–1320.
- Hadfield-Menell, D.; Russell, S. J.; Abbeel, P.; and Dragan, A. 2016. Cooperative inverse reinforcement learning. In *Advances in neural information processing systems*, 3909–3917.
- Jang, E.; Gu, S.; and Poole, B. 2016. Categorical reparameterization with gumbel-softmax. *arXiv preprint arXiv:1611.01144*.
- Kingma, D. P., and Ba, J. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Le, H. M.; Yue, Y.; Carr, P.; and Lucey, P. 2017. Coordinated multi-agent imitation learning. In *International Conference on Machine Learning*, 1995–2003.
- Lowe, R.; Wu, Y.; Tamar, A.; Harb, J.; Abbeel, O. P.; and Mordatch, I. 2017. Multi-agent actor-critic for mixed cooperative-competitive environments. In *Advances in Neural Information Processing Systems*, 6382–6393.
- Mnih, V.; Kavukcuoglu, K.; Silver, D.; Rusu, A. A.; Veness, J.; Bellemare, M. G.; Graves, A.; Riedmiller, M.; Fidjeland, A. K.; Ostrovski, G.; et al. 2015. Human-level control through deep reinforcement learning. *Nature* 518(7540):529.
- Ng, A. Y., and Russell, S. J. 2000. Algorithms for inverse reinforcement learning. In *Proceedings of the Seventeenth International Conference on Machine Learning*, 663–670. Morgan Kaufmann Publishers Inc.
- Oliehoek, F. A., and Amato, C. 2016. *A concise introduction to decentralized POMDPs*, volume 1. Springer.
- Pang, K.; Dong, M.; and Hospedales, T. 2018. Meta-learning transferable active learning policies by deep reinforcement learning.
- Rogers, E. M. 2010. *Diffusion of innovations*. Simon and Schuster.
- Sukhbaatar, S.; Fergus, R.; et al. 2016. Learning multiagent communication with backpropagation. In *Advances in Neural Information Processing Systems*, 2244–2252.
- Sutton, R. S., and Barto, A. G. 1998. *Reinforcement learning: An introduction*, volume 1. MIT press Cambridge.
- Sutton, R. S.; McAllester, D. A.; Singh, S. P.; and Mansour, Y. 2000. Policy gradient methods for reinforcement learning with function approximation. In *Advances in neural information processing systems*, 1057–1063.
- Taylor, A.; Dusparic, I.; Galván-López, E.; Clarke, S.; and Cahill, V. 2013. Transfer learning in multi-agent systems through parallel transfer. In *Workshop on Theoretically Grounded Transfer Learning at the 30th International Conf. on Machine Learning (Poster)*, volume 28, 28. Omnipress.
- Taylor, M. E.; Carboni, N.; Fachantidis, A.; Vlahavas, I.; and Torrey, L. 2014. Reinforcement learning agents providing advice in complex video games. *Connection Science* 26(1):45–63.
- Torrey, L., and Taylor, M. 2013. Teaching on a budget: Agents advising agents in reinforcement learning. In *Proceedings of the 2013 international conference on Autonomous agents and multi-agent systems*, 1053–1060. International Foundation for Autonomous Agents and Multiagent Systems.
- Wang, Y.; Lu, W.; Hao, J.; Wei, J.; and Leung, H.-F. 2018. Efficient convention emergence through decoupled reinforcement social learning with teacher-student mechanism. In *Proceedings of the 17th International Conference on Autonomous Agents and MultiAgent Systems*, 795–803. International Foundation for Autonomous Agents and Multiagent Systems.
- Zimmer, M.; Viappiani, P.; and Weng, P. 2014. Teacher-student framework: a reinforcement learning approach. In *AAMAS Workshop Autonomous Robots and Multirobot Systems*.

## Supplementary Material

### Details of Advising-level Rewards

Recall in Phase II of LeCTR, advising policies are trained to maximize advising-level rewards that should, ideally, reflect the objective of accelerating task-level learning. Selection of an appropriate advising-level reward is, itself, non-obvious. Due to this, we considered a variety of advising-level rewards, each corresponding to a different measure of task-level learning after the student executes an action advice and uses it to update its task-level policy. Advising-level rewards  $\tilde{r}_T^j$  below are detailed for the case where agent  $i$  is student and agent  $j$  teacher (i.e., flip the indices for the reverse case). Recall that the shared reward used to jointly train all advising policies is  $\tilde{r} = \tilde{r}_T^i + \tilde{r}_T^j$ .

- **Joint Value Gain (JVG):** Let  $\theta_t$  and  $\theta_{t+1}$ , respectively, denote agents' joint task-level policy parameters before and after learning from an experience resulting from action advice. The JVG advising-level reward measures improvement in task-level value due to advising, such that,

$$\tilde{r}_T^j = V(s; \theta_{t+1}) - V(s; \theta_t). \quad (3)$$

This is, perhaps, the most intuitive choice of advising-level reward, as it directly measures the gain in task-level performance due to advising. However, the JVG reward has high variance in situations where the task-level value is sensitive to policy initialization, which sometimes destabilizes training. Moreover, the JVG reward requires a full evaluation of task-level performance after each advising step, which can be expensive due to the game rollouts required.

- **Q-Teaching Reward (QTR):** The QTR advising-level reward extends Q-Teaching (Fachantidis, Taylor, and Vlahavas 2017) to MARL by using

$$\tilde{r}_T^j = I_T(o^i, a^i; h^i) = \max_a Q_T(o^i, a; h^i) - Q_T(o^i, a^i; h^i), \quad (4)$$

each time advising occurs. The motivating intuition for QTR is that teacher  $j$  should have higher probability of advising when they estimate that the student's intended action,  $a^i$ , can be outperformed by a different action (the arg max action).

- **TD Gain (TDG):** For temporal difference (TD) learners, the TDG advising-level reward measures improvement of student  $i$ 's task-level TD error due to advising,

$$\tilde{r}_T^j = |\delta_t^i| - |\delta_{t+1}^i|, \quad (5)$$

where  $\delta_t^i$  is  $i$ 's TD error at timestep  $t$ . For example, if agents are independent Q-learners at the task-level, then,

$$\delta = r + \max_{a'} Q(o', a'; \theta^i, h^i) - Q(o, a; \theta^i, h^i). \quad (6)$$

The motivating intuition for the TDG advising-level reward is that actions that are anticipated to reduce student  $i$ 's task-level TD error should be advised by the teacher  $j$ .

- **Loss Gain (LG):** The LG advising-level reward applies to many loss-based algorithms, measuring improvement of the task-level loss function used by student learner  $i$ ,

$$\tilde{r}_T^j = \mathcal{L}(\theta_t^i) - \mathcal{L}(\theta_{t+1}^i). \quad (7)$$

For example, if agents are independent Q-learners using parameterized task-level policies at the task-level, then

$$\mathcal{L}(\theta^i) = [r + \gamma \max_{a'} Q(o', a'; \theta^i, h^i) - Q(o, a; \theta^i, h^i)]^2. \quad (8)$$

The motivating intuition for the LG reward is similar to the TDG, in that teachers should advise actions they anticipate to decrease student's task-level loss function.

- **Loss Gradient Gain (LGG):** The LGG advising-level reward is an extension of the gradient prediction gain (Graves et al. 2017), which measures the magnitude of student parameter updates due to teaching,

$$\tilde{r}_T^j = \|\nabla_{\theta^i} \mathcal{L}(\theta^i)\|_2^2. \quad (9)$$

The intuition here is that larger task-level parameter updates may be correlated to learning progress.

- **Value Estimation Gain (VEG):** VEG rewards teachers when student's local value function estimates exceed a threshold  $\tau$ ,

$$\tilde{r}_T^j = \mathbb{1}(\hat{V}(\theta^i) > \tau), \quad (10)$$

using  $\hat{V}(\theta^i) = \max_{a^i} Q(o^i, a^i; \theta^i, h^i)$  and indicator function  $\mathbb{1}(\cdot)$ . The motivation here is that the student's value function approximation is correlated to its estimated performance as a function of its local experiences. A convenient means of choosing  $\tau$  is to set it as a fraction of the value estimated when no teaching occurs.

### Details of Heuristics-based Advising Approaches

Existing works on action advising typically use the state importance value  $I_\rho(s, \hat{a}) = \max_a Q_\rho(s, a) - Q_\rho(s, \hat{a})$  to decide when to advise, where  $\rho = S$  for student-initiated advising,  $\rho = T$  for teacher-initiated advising,  $Q_\rho$  is the corresponding action-value function, and  $\hat{a}$  is the student's intended action if known (or the worst-case action otherwise).  $I_\rho(s, \hat{a})$  estimates the performance difference of best versus intended student action in state  $s$ . The following is a summary of prior advising approaches:

- The **Ask Important** heuristic (Amir et al. 2016) requests advice whenever  $I_S(s, \hat{a}) \geq k$ , where  $k$  is a threshold parameter.
- **Ask Uncertain** requests when  $I_S(s, \hat{a}) < k$  (Clouse 1996), where  $k$  is a threshold parameter.
- **Early Advising** advises until advice budget depletion.
- **Importance Advising** advises when  $I_T(s, a) \geq k$  (Torrey and Taylor 2013), where  $k$  is a threshold parameter.
- **Early Correcting** advises when  $\pi_S(s) \neq \pi_T(s)$  (Amir et al. 2016).
- **Correct Important** advises when  $I_T(s) \geq k$  and  $\pi_S(s) \neq \pi_T(s)$  (Torrey and Taylor 2013), where  $k$  is a threshold parameter.
- **Q-Teaching** (Fachantidis, Taylor, and Vlahavas 2017) learns when to advise by rewarding the teacher  $I_T(s, \hat{a})$  when advising occurs. Constrained by a finite advice budget, Q-Teaching has advising performance similar to Importance Advising, with the advantage of not requiring a tuned threshold  $k$ .

Pairwise combinations of student- and teacher-initiated approaches can be used to constitute a jointly-initiated approach (Amir et al. 2016), such as ours. As shown in our experiments, application of single-agent teaching approaches yields poor performance in MARL games.

**Optimal Action Advising** Note that in the majority of prior approaches, the above heuristics are used to decide *when* to advise. To address the question of *what* to advise, these works typically assume that teachers have expert-level knowledge and always advise optimal action to students.

Optimal action advising has a strong empirical track record in single-agent teaching approaches (Torrey and Taylor 2013; Zimmer, Viappiani, and Weng 2014; Amir et al. 2016). In such settings, the assumed homogeneity of the teacher and student’s optimal policies indeed leads optimal action advice to improve student learning (i.e., when the expert teacher’s optimal policy is equivalent to student’s optimal policy). In the context of multiagent learning, however, this advising strategy has primarily been applied to games where behavioral homogeneity does not substantially degrade team performance (da Silva, Glatt, and Costa 2017). However, there exist scenarios where multiple agents learn best by exhibiting behavioral diversity (e.g., by exploring distinct regions of the state-action space), or where agents have heterogeneous capabilities/action/observation spaces altogether (e.g., coordination of 2-armed and 3-armed robots, robots with different sensors, etc.). Use of optimal action advising in cooperative multiagent tasks can lead to suboptimal joint return, particularly when the optimal policies for agents are heterogeneous. We show this empirically in several of our experiments.

In contrast to earlier optimal advising approaches, our LeCTR algorithm applies to the above settings in addition to the standard homogeneous case; this is due to LeCTR’s ability to learn a policy over not only when to advise, but also what to advise. As shown in our experiments, while existing probabilistic advising strategies (e.g., AdHocTD and AdHocVisit) attain reasonable performance in heterogeneous action settings, they do so passively by permitting students to sometimes explore their local action spaces. By contrast, LeCTR agents attain even better performance by actively learning what to advise within teammates’ action spaces; this constitutes a unique strength of our approach.

## Architecture, Training Details, and Hyperparameters

At the teaching level, our advising-level critic is parameterized by a 3-layer multilayer perceptron (MLP), consisting of internal rectified linear unit (ReLU) activations, linear output, and 32 hidden units per layer. Our advising-level actors (advice request/response policies) use a similar parameterization, with the softmax function applied to outputs for discrete advising-level action probabilities. Recurrent neural networks may also be used in settings where use of advising-level observation histories yields better performance, though we did not find this necessary in our domains. As in Lowe et al. (2017), we use the Gumbel-Softmax estimator (Jang, Gu, and Poole 2016) to compute gradients for the teaching policies over discrete advising-level actions (readers are referred

to their paper for additional details).

Policy training is conducted with the Adam optimization algorithm (Kingma and Ba 2014), using a learning rate of  $1e-3$ . We use  $\gamma = 0.95$  at the task-level and  $\tilde{\gamma} = 0.99$  at the advising-level level to induce long-horizon teaching policies. Similar to Graves et al. (2017), we use reservoir sampling to adaptively rescale advising-level rewards with time-varying and non-normalized magnitudes (all except VEG) to the interval  $[-1, 1]$ . Refer to Graves et al. for details on how this is conducted.

## Experimental Procedures

In Table 2,  $\bar{V}$  is computed by running each algorithm until convergence of task-level policies  $\pi = \langle \pi^i, \pi^j \rangle$ , and computing the mean value obtained by the final joint policy  $\pi$ . The area under the learning curve (AUC) is computed by intermittently evaluating the resulting task-level policies  $\pi$  throughout learning; while teacher advice is used *during* learning, the AUC is computed by evaluating the resulting  $\pi$  *after* advising (i.e., in absence of teacher advice actions, such that AUC measures actual student learning progress). All results and uncertainties are reported using at least 10 independent runs, with most results using over 20 independent runs. In Table 2, best results in bold are computed using a Student’s  $t$ -test with significance level  $\alpha = 0.05$ .

## Notation

The following summarizes the notation used throughout the paper. In general: superscripts  $i$  denote properties for an agent  $i$  (e.g.,  $a^i$ ); bold notation denotes joint properties for the team (e.g.,  $\mathbf{a} = \langle a^i, a^j \rangle$ ); tilde accents denote properties at the advising-level (e.g.,  $\tilde{a}^i$ ); and bold characters with tilde accent denote joint advising-level properties (e.g.,  $\tilde{\mathbf{a}} = \langle \tilde{a}^i, \tilde{a}^j \rangle$ ).

Symbol	Definition
$\mathbb{T}$	Task (a Dec-POMDP)
$\mathbb{L}$	Task-level learning algorithm
$\pi$	Joint task-level policy
$\pi^i$	Agent $i$ 's task-level policy
$\theta^i$	Agent $i$ 's task-level policy parameters
$V$	Task-level value
$Q^i$	Agent $i$ 's action value function
$\tilde{Q}^i$	Agent $i$ 's action value vector (i.e., vector of action-values for all actions)
$\mathcal{S}$	State space
$s$	State
$\mathcal{T}$	State transition function
$\mathcal{A}$	Joint action space
$\mathbf{a}$	Joint action
$\mathcal{A}^i$	Agent $i$ 's action space
$a^i$	Agent $i$ 's action
$\Omega$	Joint observation space
$\mathbf{o}$	Joint observation
$\Omega^i$	Agent $i$ 's observation space
$o^i$	Agent $i$ 's observation
$\mathcal{O}$	Observation function
$\mathcal{R}$	Reward function
$r$	Task-level reward
$\gamma$	Discount factor
$\mathcal{M}$	Task-level experience replay memory
$\delta$	Task-level temporal difference error
$\rho$	Agent's advising-level role, where $\rho = S$ for student role, $\rho = T$ for teacher role
$\tilde{\pi}_S^i$	Agent $i$ 's advice request policy
$\tilde{\pi}_T^i$	Agent $i$ 's advice response policy
$\tilde{\pi}$	Joint advising policy $\tilde{\pi} = \langle \tilde{\pi}_S^i, \tilde{\pi}_S^j, \tilde{\pi}_T^i, \tilde{\pi}_T^j \rangle$
$\tilde{\theta}$	Joint advising-level policy parameters
$\tilde{a}_S^i$	Agent $i$ 's advice request action
$\tilde{a}_\emptyset$	Special no-advice action
$\tilde{a}_T^i$	Agent $i$ 's advice response $\tilde{a}_T^i \in \mathcal{A}^j \cup \{\tilde{a}_\emptyset\}$
$\tilde{\mathbf{a}}$	Joint advising action $\tilde{\mathbf{a}} = \langle \tilde{a}_S^i, \tilde{a}_S^j, \tilde{a}_T^i, \tilde{a}_T^j \rangle$
$\beta^i$	Agent $i$ 's behavioral policy
$\tilde{o}_S^i$	Agent $i$ 's student-perspective advising obs.
$\tilde{o}_T^i$	Agent $i$ 's teacher-perspective advising obs.
$\tilde{\mathbf{o}}$	Joint advising obs. $\tilde{\mathbf{o}} = \langle \tilde{o}_S^i, \tilde{o}_S^j, \tilde{o}_T^i, \tilde{o}_T^j \rangle$
$\tilde{r}$	Advising-level reward $\tilde{r} = \tilde{r}_T^i + \tilde{r}_T^j$
$\tilde{\gamma}$	Advising-level discount factor
$\tilde{\mathcal{M}}$	Advising-level experience replay memory

## LeCTR Algorithm

### Algorithm 1 Get advising-level observations

```

1: function GETADVISEOBS( $\mathbf{o}, \theta$ )
2:   for agents  $\alpha \in \{i, j\}$  do
3:     Let  $-\alpha$  denote  $\alpha$ 's teammate.
4:      $\tilde{o}_S^\alpha = \langle o^\alpha, Q^\alpha(o^\alpha; h^\alpha) \rangle$ 
5:      $\tilde{o}_T^\alpha = \langle o^{-\alpha}, Q^{-\alpha}(o^{-\alpha}; h^{-\alpha}), Q^\alpha(o^{-\alpha}; h^{-\alpha}) \rangle$ 
6:   end for
7:   return  $\tilde{\mathbf{o}} = \langle \tilde{o}_S^i, \tilde{o}_S^j, \tilde{o}_T^i, \tilde{o}_T^j \rangle$ 
8: end function

```

### Algorithm 2 LeCTR Algorithm

```

1: for Phase II episode  $\tilde{e} = 1$  to  $\tilde{E}$  do
2:   Initialize task-level policy parameters  $\theta$ 
3:   for Phase I episode  $e = 1$  to  $E$  do
4:      $\mathbf{o} \leftarrow$  initial task-level observation
5:     for task-level timestep  $t = 1$  to  $t_{end}$  do
6:        $\tilde{\mathbf{o}} \leftarrow$  GETADVISEOBS( $\mathbf{o}, \theta$ )
7:       for agents  $\alpha \in \{i, j\}$  do
8:         Exchange advice  $\tilde{a}^\alpha$  via advising policies
9:         if No advising occurred then
10:           Select action  $a^\alpha$  via local policy  $\pi^\alpha$ 
11:         end if
12:       end for
13:        $\mathbf{a} \leftarrow \langle a^i, a^j \rangle, \tilde{\mathbf{a}} \leftarrow \langle \tilde{a}^i, \tilde{a}^j \rangle$ 
14:        $r, \mathbf{o}' \leftarrow$  Execute action  $\mathbf{a}$  in task
15:        $\theta'^i \leftarrow \mathbb{L}^i(\theta^i, \langle o^i, a^i, r, o'^i \rangle)$ 
16:        $\theta'^j \leftarrow \mathbb{L}^j(\theta^j, \langle o^j, a^j, r, o'^j \rangle)$ 
17:        $\tilde{\mathbf{o}}' \leftarrow$  GETMETA OBS( $\mathbf{o}', \theta'$ )
18:        $\tilde{r}_T^i, \tilde{r}_T^j \leftarrow$  Compute advising-level rewards
19:       Store  $\langle \tilde{\mathbf{o}}, \tilde{\mathbf{a}}, \tilde{r} = \tilde{r}_T^i + \tilde{r}_T^j, \tilde{\mathbf{o}}' \rangle$  in buffer  $\tilde{\mathcal{M}}$ 
20:     end for
21:   end for
22:   Update advising-level critic by minimizing loss,

```

$$\mathcal{L}(\tilde{\theta}) = \mathbb{E}_{\tilde{\mathbf{o}}, \tilde{\mathbf{a}}, \tilde{r}, \tilde{\mathbf{o}}' \sim \tilde{\mathcal{M}}} [(\tilde{r} + \gamma \tilde{Q}(\tilde{\mathbf{o}}', \tilde{\mathbf{a}}'; \tilde{\theta}) - \tilde{Q}(\tilde{\mathbf{o}}, \tilde{\mathbf{a}}; \tilde{\theta}))^2] \Big|_{\tilde{\mathbf{a}}' = \tilde{\pi}(\tilde{\mathbf{o}}')}$$

```

23:   for agents  $\alpha \in \{i, j\}$  do
24:     for roles  $\rho \in \{S, T\}$  do
25:       Update advising policy parameters  $\tilde{\theta}_\rho^\alpha$  via,

```

$$\nabla_{\tilde{\theta}} J(\tilde{\theta}) = \mathbb{E}_{\tilde{\mathbf{o}}, \tilde{\mathbf{a}} \sim \tilde{\mathcal{M}}} \left[ \sum_{\substack{\alpha \in \{i, j\} \\ \rho \in \{S, T\}}} \nabla_{\tilde{\theta}_\rho^\alpha} \log \tilde{\pi}_\rho^\alpha(\tilde{a}_\rho^\alpha | \tilde{o}_\rho^\alpha) \nabla_{\tilde{a}_\rho^\alpha} \tilde{Q}(\tilde{\mathbf{o}}, \tilde{\mathbf{a}}; \tilde{\theta}) \right],$$

```

26:   end for
27: end for
28: end for

```