



# HAMMER: Multi-level coordination of reinforcement learning agents via learned messaging

Nikunj Gupta<sup>1</sup> · G. Srinivasaraghavan<sup>2</sup> · Swarup Mohalik<sup>3</sup> · Nishant Kumar<sup>4</sup> · Matthew E. Taylor<sup>5,6</sup>

Received: 15 November 2021 / Accepted: 20 September 2023

© The Author(s), under exclusive licence to Springer-Verlag London Ltd., part of Springer Nature 2023

## Abstract

Cooperative multi-agent reinforcement learning (MARL) has achieved significant results, most notably by leveraging the representation-learning abilities of deep neural networks. However, large centralized approaches quickly become infeasible as the number of agents scale, and fully decentralized approaches can miss important opportunities for information sharing and coordination. Furthermore, not all agents are equal—in some cases, individual agents may not even have the ability to send communication to other agents or explicitly model other agents. This paper considers the case where there is a single, powerful, *central agent* that can observe the entire observation space, and there are multiple, low-powered *local agents* that can only receive local observations and are not able to communicate with each other. The central agent's job is to learn what message needs to be sent to different local agents based on the global observations, not by centrally solving the entire problem and sending action commands, but by determining what additional information an individual agent should receive so that it can make a better decision. In this work, we present our MARL algorithm HAMMER, describe where it would be most applicable, and implement it in the cooperative navigation and multi-agent walker domains. Empirical results show that (1) learned communication does indeed improve system performance, (2) results generalize to heterogeneous local agents, and (3) results generalize to different reward structures.

**Keywords** Multi-agent reinforcement learning · Learning to communicate · Heterogeneous agent learning

## 1 Introduction

The field of multi-agent reinforcement learning (MARL) combines ideas from single-agent reinforcement learning (SARL), game theory, and multi-agent systems. Cooperative MARL calls for simultaneous learning and interaction of multiple agents in the same environment to achieve shared

goals. Applications like distributed logistics [66], package delivery [47], and disaster rescue [38] are domains that can be modeled naturally using this framework. However, even cooperative MARL suffers from several complications inherent to multi-agent systems, including non-stationarity [3], a potential need for coordination [3], the curse of dimensionality [49], and global exploration [32].

Multi-agent reasoning has been extensively studied in MARL [32, 37] in both centralized and decentralized settings. While very small systems could be completely centralized, decentralized implementation becomes indispensable as the number of agents increase, and to cope with the exponential growth in the joint observation and action spaces. However, it often suffers from synchronization issues [32] and complex teammate modeling [1]. Moreover, independent learners may have to optimize their own, or the global, reward from only local, private observations [58]. In contrast, centralized approaches can leverage global information and mitigate non-stationarity through full awareness of all teammates.

---

✉ Nikunj Gupta  
ng2531@nyu.edu

<sup>1</sup> New York University, New York, USA

<sup>2</sup> International Institute of Information Technology, Bangalore, India

<sup>3</sup> Ericsson Research, Bangalore, India

<sup>4</sup> Indian Institute of Technology (Banaras Hindu University), Varanasi, India

<sup>5</sup> University of Alberta, Edmonton, Canada

<sup>6</sup> Alberta Machine Intelligence Institute (Amii), Edmonton, Canada

In MARL, communication has been shown to be an important aspect, especially in tasks requiring coordination. For instance, agents tend to locate each other or the landmarks more easily using shared information in navigation tasks [14]. Communication can also influence the final outcomes in group strategy coordination [17, 65]. There have been significant achievements using explicit communication in video games like StarCraft II [39] as well as in mobile robotic teams [31], smart-grid control [40], and autonomous vehicles [4]. Communication can be in the form of sharing experiences among the agents [67], sharing low-level information like gradient updates via communication channels [10] or sometimes directly advising appropriate actions using a pretrained agent (teacher) [60] or even learning teachers [36, 50].

Inspired by the advantages of centralized learning and communication for synchronization, we propose multi-level coordination among intelligent agents via messages learned by a separate agent to ease the localized learning of task-related policies. We propose a single *central agent* designed to learn high-level messages based on complete knowledge of all the local agents in the environment. These messages are communicated to the *independent learners* who are free to use or discard them while learning local policies to achieve a set of shared goals. By introducing centralization in this manner, the supplemental agent can play the role of a *facilitator* of learning for the rest of the team. Furthermore, the independent learners need not be as powerful—they do not need to communicate with other agents or explicitly monitor/model other agents in the system—as required if they must train to communicate or model other agents alongside learning task-specific policies. Rather, it might not be feasible to assume such powerful agents in the environment. This motivates HAMMER, where the independent learners can be simple agents interacting with the environment based on their private local observations. Only one central agent needs to be powerful enough to monitor agents and augment their observations with learned messages. Warehouse management and traffic lights management are two example applications that could fit these assumptions well. This novel setting is explained further in details in Sect. 3.

A hierarchical approach to MARL is not new—we will contrast with other existing methods in Sect. 2. However, the main insight of our algorithm is to learn to communicate relevant pieces of information from a global perspective to help agents with limited capabilities improve their performance. Potential applications include autonomous warehouse management [8] and traffic light control [28], where there can be a centralized monitor. After we introduce our algorithm, HAMMER, we will show results in two very different simulated environments to showcase its versatility. OpenAI’s multi-agent cooperative navigation lets agents learn in a continuous state space with

discrete actions and global team rewards. In contrast, Stanford’s multi-agent walker environment has a continuous action space and agents can receive only local rewards.

The main contributions of this paper are to explain a novel and important setting that combines agents with different abilities and knowledge (Sect. 3), introduce the HAMMER algorithm that addresses this setting (Sect. 4), and then empirically demonstrate that HAMMER can make significant improvements in learning decision policies for agents (Sect. 6) in two multi-agent domains.

## 2 Background and related work

This section will provide a summary of background concepts necessary to understand the paper and a selection of related work.

### 2.1 Single-agent reinforcement learning

Single-agent reinforcement learning (SARL) can be formalized in terms of Markov Decision Processes (MDPs). An MDP is defined as a tuple  $\langle S, A, P, R, \gamma \rangle$ , where  $S$  is the set of states,  $A$  is the set of available actions for the agent,  $P: S \times A \times S \rightarrow [0, 1]$  is the state transition function,  $R: S \times A \rightarrow \mathcal{R}$  is the reward function, and  $\gamma \in (0, 1]$  is a discount factor. Actions, selected by a policy  $\pi: S \times A \rightarrow [0, 1]$ , are taken and the agent tries to maximize the return, which is the expectation over the sum of discounted future rewards:  $\mathbf{E}_\pi \left[ \sum_{t=0}^{\infty} \gamma^t R(s_t, a_t) \right]$ , where  $t$  is the time step.

The goal of solving a MDP is generally to find a policy that guarantees a maximum reward. To maximize this expectation, one class of RL algorithms aim to learn a Q-value (state-action) function and the Bellman optimality equation for the same can be defined as

$$Q^*(s \in S, a \in A) = \sum_{s' \in S} P(s, a, s') \left[ R(s, a) + \gamma \max_{a' \in A} Q^*(s', a') \right]$$

It provably converges to the optimal function  $Q^*$  under certain conditions. It has also become popular to use nonlinear function approximators to scale to huge state spaces—like in DQN [33]. DQN learns the Q-value function corresponding to the optimal policy by minimizing the following loss:

$$\mathbf{L}(\theta) = \mathbf{E}_{s,a,r,s'} \left[ \left( Q^*(s, a | \theta) - (r + \gamma \max_{a'} Q_t^*(s', a')) \right)^2 \right],$$

where  $Q_t$  is the target Q-function, and its parameters are updated only periodically—this helps in stable learning.

Based on the aforementioned Q-value function, gradient of the policy can be defined as:

$$\nabla_{\theta} \mathbf{J}(\theta) = \mathbf{E}_{s \in S, a \in \pi_{\theta}} [\nabla_{\theta} \log(\pi_{\theta}(a, s)) Q^{\pi}(s, a)]$$

A popular choice for solving RL tasks is to use policy gradients, where the parameters of the policy  $\theta$  are directly updated to maximize an objective  $\mathbf{J}(\theta)$  by moving in the direction of  $\nabla \mathbf{J}(\theta)$ . In our work, we make use of Proximal Policy Optimization (PPO) [46], which reduces challenges like exhibiting high variance gradient estimates, being sensitive to the selection of step-size, progressing slowly, or encountering catastrophic drops in performance. Moreover, it is relatively easy to implement. Also, its objective function, well-suited for updates using stochastic gradient descent, can be defined as follows:

$$L^{CLIP}(\theta) = \mathbf{E}_t [\min(r_t(\theta)A_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon)A_t)],$$

where  $r_t$  is the ratio of probability under the new and old policies, respectively,  $\left(\frac{\pi_{\theta}(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)}\right)$ ,  $A_t$  is the estimated advantage at time  $t$ , and  $\epsilon$  is a hyperparameter.

## 2.2 Multi-agent reinforcement learning

Single-agent reinforcement learning can be generalized to competitive, cooperative, or mixed multi-agent settings. We focus on the fully cooperative multi-agent setting, which can be described as a multi-agent extension to MDPs (Markov games). A Markov game for  $n$  local agents, and an additional centralized agent in our case, can be defined as the tuple  $\langle S, U, O_1, \dots, O_n, A_1, \dots, A_n, P, R, \gamma \rangle$ , where  $S$  is the set of all configurations of environment states for all agents,  $U$  is the set of all actions for the central agent,  $O_1, \dots, O_n$  represent the observations of each local agent,  $A_1, \dots, A_n$  correspond to the set of actions available to each local agent, and  $P$  is the state transition function.  $\gamma$  is the discount factor. The state transitions in a multi-agent case are a result of the joint action of all the agents  $U \times A_1 \times \dots \times A_n$ . The policies join together to form a joint policy  $\pi: S \times U \times \mathbf{A}$ . There can be two possible reward structures for the agents. First, they could share a common team reward signal,  $R: S \times A \rightarrow \mathfrak{R}$ , defined as a function of the state  $s \in S$  where  $A: A_1 \times \dots \times A_n$  is the joint action taken by the agents. In the case of such shared rewards, agents aim to directly maximize the returns for the team. Second, each agent  $i$  could receive its own reward  $R_i: O_i \times A_i \rightarrow \mathfrak{R}$ . A localized reward structure means that agents maximize their own individual expected discounted return  $(\sum_{t=0}^{\infty} \gamma^t r_i^t)$ .

## 2.3 Relevant prior work

Recent works in both SARL and MARL have employed deep learning methods to tackle the high dimensionality of the observation and action spaces [12, 27, 33, 39, 57].

Several works in the past have taken advantage of hierarchical approaches to MARL. The MAXQ algorithm

was designed to provide for a hierarchical break-down of the reinforcement learning problem by decomposing the value function for the main problem into a set of value functions for the sub-problems [7]. Tang et al. [59] used temporal abstraction to let agents learn high-level coordination and independent skills at different temporal scales together. Kumar et al. [20] present another framework benefiting from temporal abstractions to achieve coordination among agents with reduced communication complexities. Factored  $Q$  value functions, which represent the joint value but decompose it as the sum of a number of local components (each involving only a subset of the agents), have recently shown successful in deep MARL [5, 41]. These works show positive results from combining centralized and decentralized approaches in different manners and are therefore closely related to our work. Vezhnevets et al. [63] introduce Feudal networks in hierarchical reinforcement learning and Ma & Wu [30] propose a novel MARL approach that leverages ideas from feudal hierarchy. Both of these works employ a Manager-Worker framework that is related to HAMMER. However, there are some key differences. In their case, the manager directly interacts with the environment to receive the team reward and accordingly distributes it among the workers (analogous to setting their goals), whereas in our work, the central agent interacts indirectly and receives the same reward as the local agents. Here, the central agent is only allowed to influence the actions of independent learners rather than set their goals explicitly. One further distinction is that they partly pretrain their workers before introducing the manager into the scene. In contrast, our results show that when the central agent and the independent learners simultaneously learn, they achieve better performance.

Some works have developed architectures that use centralized learning but ensure decentralized execution. COMA [12] used a centralized critic to estimate the Q-function along with a counterfactual advantage for the decentralized actors in MARL. The VDN [54] architecture trained individual agents by learning to decompose the team value functions into agent-wise value functions in a centralized manner. QMIX [42] employs a mixing network to factor the joint action-value into a monotonic nonlinear combination of individual value functions for each agent. Another family of works [9, 43] have developed novel decentralized coordination approaches to efficiently use constrained computational and communication resources. Another work, MADDPG [29], extends deep deterministic policy gradients (DDPG) [27] to MARL. They learn a centralized critic for each agent and continuous policies for the actors and allow explicit communication among agents. Even though the prior research works mentioned here address a similar setting and allow for using extra globally accessible information like in our work, they mainly aim at decentralized execution

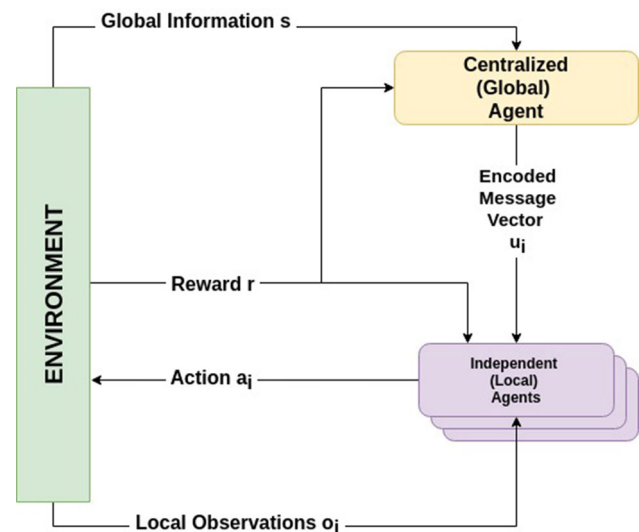
which applies to domains where global view is unavailable. In contrast, we target domains where a global view is accessible to a single agent, even during execution.

There has been considerable progress in learning by communication in cooperative settings involving partially observable environments. Reinforced Inter-Agent Learning (RIAL) and Differentiable Inter-Agent Learning (DIAL) [10] use neural networks to output communication messages in addition to the agent's Q-values. RIAL used a shared network to learn a single policy whereas DIAL used gradient sharing during learning and communication actions during execution. Both methods use discrete communication channels. On the other hand, CommNet [53], used continuous vectors, enabled multiple communication cycles per time step and the agents were allowed to freely enter and exit the environment. Lazaridou et al. [23] and Mordatch et al. [35] trained the agents to develop an emergent language for communication. Furthermore, standard techniques used in deep learning, such as dropout [52], have inspired works where messages of other agents are dropped out during learning to work well even in conditions with only limited communication feasible [19]. However, in all these works, the goal is to learn inter-agent communication alongside local policies that suffer from the bottleneck of simultaneously achieving effective communication and global collaboration [48]. They also face difficulty in extracting essential and high-quality information for exchange among agents [48]. Further, unlike HAMMER, these works expect that more sophisticated agents are available in the environment in terms of communication capabilities or the ability to run complex algorithms to model other agents present—which might not always be feasible.

One of the popular ways to carry out independent learning is by emergent behaviors [24, 25, 57], where each agent learns its own private policy and assumes all other agents to be a part of the environment. This method disregards the underlying assumptions of single-agent reinforcement learning, particularly the Markov property. Although this may achieve good results [32], it can also fail due to non-stationarity [22, 62]. Self-play can be a useful concept in such cases [2, 51, 61], but it is still susceptible to failures through the loss of past knowledge [21, 26, 44]. Gupta et al. [15] extend three SARL algorithms, Deep Q Network (DQN) [34], Deep Deterministic Policy Gradient (DDPG) [27], and Trust Region Policy Optimization (TRPO) [45], to cooperative multi-agent settings.

### 3 Setting

This section details a novel setting in which multiple agents—the central agent and the local agents—with different capabilities and knowledge are combined (see



**Fig. 1** Our cooperative MARL setting: a single global agent sends messages to help multiple independent local agents act in an environment

Fig. 1). The HAMMER algorithm is designed for this type of cooperative multi-agent environment.

Consider a warehouse setting where lots of small, simple, robots fetch and stock items on shelves, as well as bring them to packing stations. If the local agents could communicate among themselves, they could run a distributed reasoning algorithm, but this would require more sophisticated robots and algorithms. Or, if the observations and actions could be centralized, one very powerful agent could determine the joint action of all the agents, but this would not scale well and would require a very powerful agent (to address an exponential growth in the observation and actions spaces with the number of agents). Section 6 will make this more concrete with two multi-agent tasks. Now, assume there is an additional central agent in the team, that has a global perspective, unlike the local agents, who receive only local observations. Further, the central agent is more powerful—not only does it have access to more information, but it can also communicate messages to all local agents. The local agents can only transmit their observations and actions to the central agent and receive messages—local agents rely on the communicated messages to know about other agents in the environment. The central agent must learn to encapsulate available information into small messages to facilitate local agents.

Having described an overview of the setting, we can take a closer look at the inputs, outputs, and roles of the agents in the heterogeneous system as described in Fig. 1. Mathematically, HAMMER's setting can be defined by the following tuple

$$\langle \mathcal{N} \cup \{\mathcal{H}\}, S, \{O^i\}_{i \in \mathcal{N}}, \{A^i\}_{i \in \mathcal{N}}, \mathcal{U}, P, R \rangle$$

The centralized agent  $\mathcal{H}$  receives a global observation  $s \in S$  on every time step and outputs a unique message (equivalently, executes an action),  $u_i \in \mathcal{U}$ , to each of the local agents, where  $i(\in \mathcal{N})$  is the agent identifier and  $\mathcal{N}$  is the set of local agents. Its global observation  $s$  is the union of all the local observations  $o_i \in O_i$  and actions of the independent learners  $a_i \in A_i$ —can either be obtained from the environment or transmitted to it by the local agents at every time step.  $u_i$  encodes a message vector that a local agent can use to make better decisions. Local agents receive a partial observation  $o_i$ , and a private message  $u_i$  from the central agent. Based on  $o_i$  and  $u_i$ , at each time step, all  $n$  local agents will choose their actions simultaneously, forming a joint action  $(A_1 \times \dots \times A_n)$  and causing a change in the environment state according to the state transition function  $P$ .

Upon changing the dynamics of the environment, a reward  $r \in R$ —which could be team-based or localized—is sent back to the local agents, using which they must learn how to act. If the messages from the central agent were not useful, local agents could learn to ignore them. Every time the central agent communicates a message  $u_i$  to a local agent, it learns from the same reward as is obtained by that local agent on performing an action  $a_i$  in the environment. In other words, the central agent does not directly interact with the environment to get feedback, instead, it learns to output useful messages by looking at how the independent agents performed in the environment using the messages it communicated. In domains with localized rewards for

agents, the central agent gets a tangible feedback for sent messages, whereas, in the case of team rewards, it needs to learn using comparatively abstract feedback. In Sect. 6, we show that HAMMER generalizes to both the reward structures.

## 4 The HAMMER Algorithm

This section introduces HAMMER, the *Heterogeneous Agents Mastering Messaging to Enhance Reinforcement learning* algorithm, designed for the cooperative MARL setting discussed above.

There are multiple local and independent agents in an environment that is given tasks. Depending on the domain, they may take discrete or continuous actions to interact with the environment. In HAMMER, we introduce a single, relatively powerful central agent into the environment, capable of 1) obtaining a global view of all the other agents present in an environment, and 2) transmitting messages to them. It learns a separate policy and aims to support the local team by communicating short messages to them. It is designed to use both a global or local reward structure. As a result, local agents' private observations will now have additional messages sent to them by the central agent and they can choose to use or discard this information while learning their policies.

---

### Algorithm 1: HAMMER

---

```

1 Initialize Actor-Critic Network for independent agents (shared network) (IA),
   Actor-Critic Network for the central agent (CA) or a multi-layered perceptron if
   gradients from IA are backpropagated to CA, and two experience replay memory
   buffers (B and B');
2 for episode  $e = 1$  to TOTAL_EPISODES do
3   Fetch combined initial random observations  $s = [o_1, \dots, o_i]$  from environment ( $o_i$ 
   is agent  $i$ 's local observation);
4   Input: Concat ( $s = [o_1, \dots, o_i]$ )  $\rightarrow$  CA ;
5   for time step  $t = 1$  to TOTAL_STEPS do
6     for each agent  $n_i$  do
7       Output: message vector  $u_i \leftarrow$  CA, for agent  $n_i$ ;
8       Pass each message vector  $u_i$  through a regularization unit ;
9       Input: Concat ( $o_i \in s, u_i$ )  $\rightarrow$  IA ;
10      Output: local action  $a_i \leftarrow$  IA ;
11    end
12    Perform sampled actions in environment and get next set of observations  $s'$ 
    and rewards  $r_i$  for each agent ;
13    Add experiences in B and B' for CA and IA respectively;
14    if update interval reached then
15      Sample random minibatch  $b \in B$  and  $b' \in B'$ ;
16      Train IA on  $b'$  using stochastic policy gradients ;
17      Train CA on  $b$  or directly by backpropagating gradients from IA ;
18    end
19  end
20 end

```

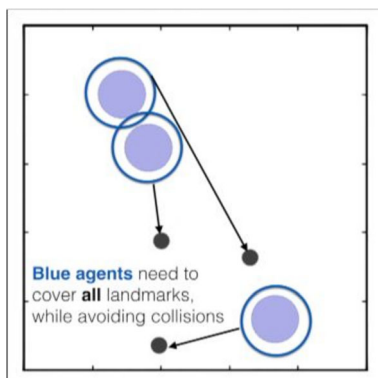
---

As described by Algorithm 1, in every iteration, the centralized agent receives the union of private observations of all local agents (line 3). It encodes its input and outputs an individual message vector for each agent (line 7). These messages from the central agent are sent to the independent learners, augmenting their private partial observations obtained from the environment (line 9). Then, they output an action (line 10) affecting the environment (line 12). Reward for the joint action performed in the environment is returned (line 12) and is utilized as feedback by all the agents to adjust their parameters and learn private policies (lines 15–17).

There are multiple techniques for training the central agent to learn how to communicate. One strategy could be to employ any RL algorithm to train HAMMER's central agent to learn its policy and use the local agents' reward as a gradient signal. A second strategy would be to push gradients from the local agent's network to HAMMER's central agent network, by directly connecting the latter's outputted communication actions to the input of the local agent's network. This strategy is inspired from a number of other relevant works [11, 53]. Another strategy could be to allow messages  $u_i$  output by the central agent to first be processed by a regularization unit, like  $RU(u_i) = \text{Logistic}(\mathbf{N}(u_i, \sigma))$ , where  $\sigma$  is the standard deviation of the noise added to the channel and then, be passed to the local agent's network [6, 11, 16]. In Sect. 6, we compare our results for all of these strategies.

## 5 Tasks and implementation details

This section details the two multi-agent environments used to evaluate HAMMER. In addition to releasing our code after acceptance, we fully detail our approach so that results are replicable.



**Fig. 2** The cooperative navigation environment is composed of blue agents and black (stationary) landmarks the agents must cover, while avoiding collisions

### 5.1 Cooperative navigation

Cooperative navigation is one of the Multi-Agent Particle Environments [29]. It is a two-dimensional cooperative multi-agent task with a continuous observation space and a discrete action space consisting of  $n$  movable agents and  $n$  fixed landmarks. Figure 2 shows the case where  $n = 3$ . Agents occupy physical space (i.e., are not point masses), perform physical actions in the environment, and have the same action and observation spaces. The agents must learn to cover all the landmarks while avoiding collisions, without explicit communication. The global reward signal, seen by all agents, is based on the proximity of any agent to each landmark and the agents are penalized upon colliding with each other. The team reward can be defined by the equation:

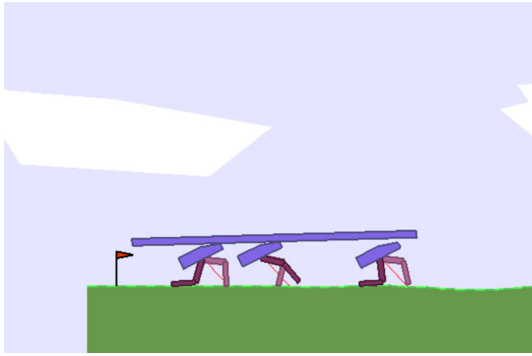
$$R = \left[ - \sum_{n=1, l=1}^{N, L} \min(\text{dist}(a_n, l)) \right] - c,$$

where  $N$  is the number of agents and  $L$  is the number of landmarks in the environment. The function  $\text{dist}()$  calculates the distance in terms of the agents' and landmarks'  $(x_i, y_i)$  positions in the environment. The number of collisions,  $c$ , among the agents and is set as a penalty of -1 for each time two agents collide. The action set is discrete, corresponding to moving in the four cardinal directions or remaining motionless. Each agent's observation includes the relative positions of other agents and landmarks within the frame. Note that the local observations do not convey the velocity (movement direction) of other agents. Consistent with past work, the initial positions of all the agents and the landmark locations are randomized at the start of each episode, and each episode ends after 25 time steps.

To test HAMMER, we modify this task so that a centralized agent receives the union of local agents' observations at every time step. We also conduct relevant ablation studies on modified versions of this environment (exhibiting graceful degradation to the system, particularly to what the local agents can observe) to understand the contribution of HAMMER's central agent to the overall system. We are most interested in this environment because of the motivating robotic warehouse example.

### 5.2 Multi-agent walker

Multi-agent walker is a more complex, continuous control benchmark locomotion task [15]. A package is placed on top of  $n$  pairs of robot legs which can be controlled. The agents must learn to move the package as far to the right as possible, without dropping it. The package is large enough (it stretches across all of the walkers) that the agents must cooperate. The environment demands high inter-agent



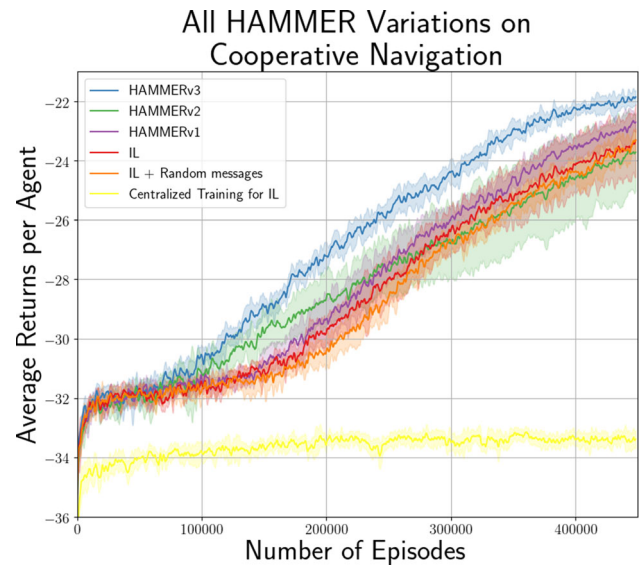
**Fig. 3** Multi-Agent Walker environment; multiple robots work together to act in continuous action spaces to transport a package over varying terrain to a destination without falling

coordination for them to successfully navigate a complex terrain while keeping the package balanced. This environment supports both team and individual rewards, but we focus on the latter, to demonstrate the effectiveness of HAMMER in individual feedback settings (as team rewards were already explored in the cooperative navigation task). Each walker is penalized with  $-10$  if it falls and all walkers receive a reward of  $-100$  if the package falls. The agents also receive a small shaped reward of  $-5$  times the change in their head angle to keep the head straight. Throughout the episode, a positive reward, proportional to the change in the package distance, is awarded to each walker. However, there is no supplemental reward for reaching the destination. By default, the episode ends if any walker falls, the package falls, after 500 time steps, or if the package successfully reaches the destination. Each agent receives 31 real-valued numbers, representing information about noisy Lidar measurements of the terrain and displacement information related to the neighboring agents. The action space is a 4-dimensional continuous vector, representing the torques in the walker's two joints of both legs. Figure 3 is an illustration of the environment with  $n = 3$  agents.

Similar to the Cooperative Navigation domain, this environment is tweaked so that all local agents in it can transmit their private observations to the central agent. This environment was chosen primarily to test our approach in a continuous action domain and in cases where individual agents receive their own local rewards instead of global team rewards. It is also significantly more challenging than the cooperative navigation domain.

## 6 Results

This section describes the experiments conducted to test HAMMER's potential of encapsulating and communicating learned messages and speeding up independent learning on



**Fig. 4** HAMMER agents outperform independent PPO learners in cooperative navigation. Using a similar framework as HAMMER, but providing the local agents with random messages, causes degraded performance (as expected). HAMMER also significantly outperforms centrally learned policy for independent agents

the two environments—cooperative navigation and multi-agent walker—whose details are described in the previous section. All the curves are averaged over five independent trials. Additional ablative studies were performed on the modified versions of cooperative navigation environment.

### 6.1 Cooperative navigation results

We investigated in detail the learning of independent learners in the cooperative navigation environment under different situations. Learning curves used for evaluation are plots of the average reward per local agent as a function of episodes. Moreover, to smooth the curves for readability, a moving average with a window size of 1500 episodes was used for each of the cases.

At first, we let the local agents learn independently, without any aid from other sources. The corresponding curve (red), as shown in Fig. 4, can also act as our baseline to evaluate learning curves obtained when the learners are equipped with additional messages. As described earlier, the experimental setup is consistent with earlier work [15].

Next, we used a central agent to learn messages, as described by HAMMER, and communicated them to the local learners to see if the learning improved. As described earlier (Sect. 4), HAMMER's central agent is trained using more than one strategies and we compare the results here. First, the central agent learned to communicate by employing PPO, and used the local agents' team-based reward as its feedback (HAMMERV1). HAMMERV1 performed only slightly better than independent learners alone, over a

training period of 500,000 episodes, as can be seen in Fig. 4 (purple). This training strategy does provide a small amount of improvement, which could be because of two types of problems: (*Case 1*) The central agent may emit a relevant message to a local agent, but the local agent is unable to intercept it correctly resulting in a poor reward. In such a case, the central agent gets penalized, even though it performed its part well. (*Case 2*) The central agent emits relevant messages to some local agents, and non-relevant messages to others. In such a case, the central agent is penalized in spite of emitting relevant messages for some agents.

Now, to avoid the problems encountered in the previous case, gradients from the local agent’s network were pushed to HAMMER’s central agent network, by directly connecting the outputted communication vectors to the input of the local agent’s network (HAMMERV2). Letting gradients flow in this manner, gives the central agent richer feedback, thus reducing the required amount of learning by trial and error, and easing the discovery of effective real-valued messages. Since these messages function as any other network activation, gradients can be passed back to the central agent’s network, allowing end-to-end backpropagation across the entire framework. As shown in Fig. 4, HAMMERV2 (green) performs significantly better than unaided independent learners as well as HAMMERV1. However, HAMMER agents seem to slow after 250,000 episodes. We speculate that a larger action space (real-valued messages) causes the central agent learns to “over-encode” the global information available to it. Having included this extra information along with the relevant pieces would have masked the private observations of the local agents, slowing their progress.

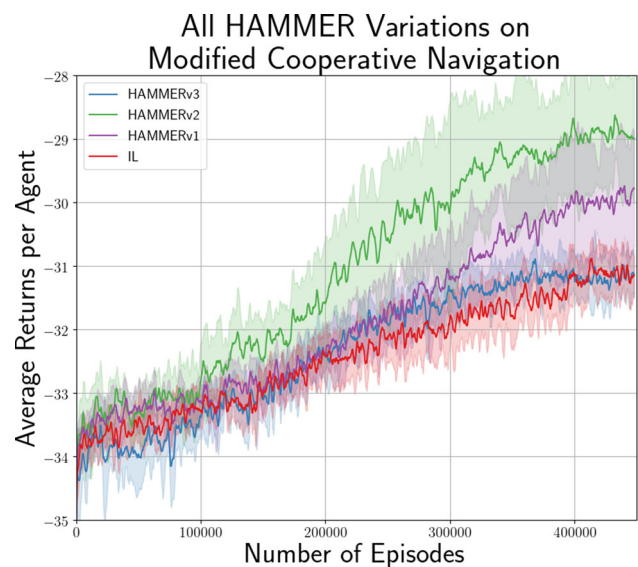
Finally, addressing the problem of “over-encoding” faced in the previous case, the emitted messages were first processed by a regularization unit— $RU(u_i) = \text{Logistic}(\mathbf{N}(u_i, \sigma))$ , to encourage discretization of the messages (by pushing the activations of the communication vectors into two different modes during training, i.e., where the noise is minimized), and then passed to the local agent’s network (HAMMERV3). Figure 4 illustrates that HAMMERV3 (blue) outperforms all the other cases. Hence, using a regularization unit to limit the central agent’s action space was essential for learning a better communication policy. Using a noisy channel in this manner has been useful in other relevant works too—differentiable inter-agent learning [11], training document models [16] and performing classification [6]. For our experiments, we used a value of  $\sigma = 0.2$ .

From these experiments, our results show that: (1) HAMMER agents were able to learn much faster when compared to independent local agents, and (2) the central agent was able to successfully learn to produce useful smaller

messages. Recall that the total global observation vector has  $18 \times 3 = 54$  real-valued numbers (for 3 agents and 3 landmarks in the environment), and our message uses only 1.

To help evaluate the communication quality, random messages of the same length were generated and communicated to the local agents to see if the central agent was indeed learning relevant or useful messages. As expected, random messages induce a larger observation space with meaningless values and degrade the performance of independent learners (Fig. 4, orange curve). This also supports the claim that the central agent is learning much better messages to communicate (rather than sending random values) as it outperforms the independent learning of agents provided with random messages. Note here that the local agents are free to learn to ignore unhelpful messages and independent agents’ learning while receiving random messages only degrades performance slightly.

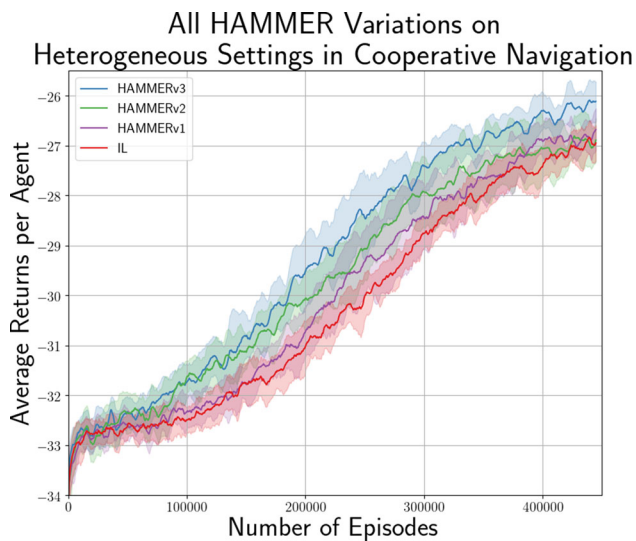
To confirm that the central agent is not simply forwarding a compressed version of the global information vector to all the agents, we let the independent agents learn using complete information about all the other agents in the environment. We call this Concatenated Observations Independent Learning (COIL). In other words, all the agents witness a joint observation space and have to learn to take individual actions in the environment. We make sure that the sequence of concatenating all the agents’ local



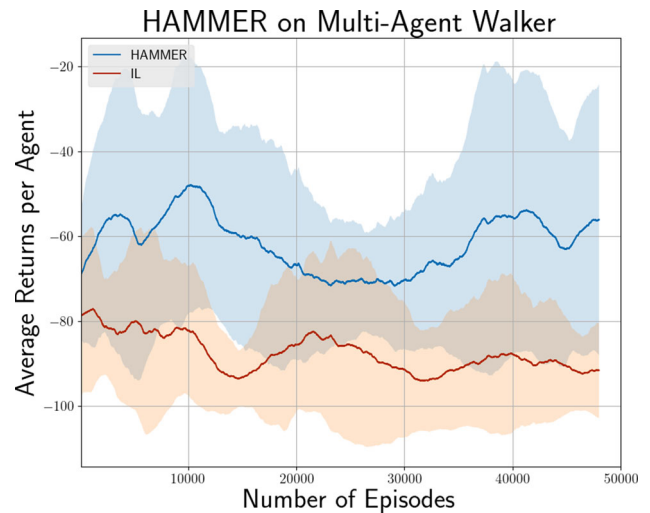
**Fig. 5** Results on a modified cooperative navigation environment—disallowing the local agents to be able to observe each other, hence necessitating the need for communication via the central agent to coordinate and cover the landmarks in the environment cooperatively. HAMMER’s ability to perform in this setting shows that the central agent is indeed learning effective communication to help the local agents coordinate

observations remains constant, hence ensuring that the agents can learn to differentiate them. The performance drops drastically in this case (Fig. 4, yellow curve). This suggests that the central agent in HAMMER is learning to encapsulate only partial but relevant information as messages and communicates them to facilitate the learning of local agents. Complete information would have become too overwhelming for the localized agents to learn how to act and hence slowed the progress, as is clear by the curve shown in Fig. 4.

**Ablation Studies.** First, we perform an ablation experiment to validate that HAMMER’s central agent is indeed facilitating independent agents’ learning and helping them coordinate by sending out relevant messages. For this, we modify the environment preventing agents from seeing each other. Local agents be unable to observe other agents and would have to rely on the central agent for coordinating and covering respective landmarks. Results of HAMMER in this scenario, compared to unaided independent learners, are in Fig. 5. As expected, it became difficult for the independent learners to learn the cooperation-intensive task on their own. On the other hand, HAMMER agents learned substantially better policies faster. Note here that HAMMERV2 in this case performs better than HAMMERV3. We speculate that this is due to local agents’ greater dependency on central agent’s communicated messages for coordinating, and hence, HAMMER requiring a larger spectrum to encode and communicate more information into its messages.



**Fig. 6** HAMMER applied on another modified version of the cooperative navigation environment such that the local agents in the environment were heterogeneous in nature—one of the local agents was unable to observe the other agents, whereas the other two agents could. Results show that HAMMER was able to generalize to these settings too



**Fig. 7** HAMMER considerably improves the performance over independent learners in the multi-agent walker task too

In our second ablation study, we validate whether HAMMER would generalize to a setting with heterogeneous local agents. This experiment uses an environment where one of the local agents was unable to observe other agents, while the other two agents still received their original observations. All local agents still learn and use the same policy, but one of them cannot see the others. As expected, independent learners had increased difficulty in learning to coordinate (Fig. 6) and were outperformed by all variants of learning strategies for HAMMER. We conclude that HAMMER generalizes to heterogeneous settings too.

In summary, HAMMER successfully summarizes relevant global knowledge into small real-valued messages to individual learners that outperforms other cases—(1) unaided independent learning, (2) independent agents supplied with random messages, and (3) COIL training of independent agents. Specifically, HAMMERV3 (i.e., HAMMER trained by directly passing message gradients from the local agent’s network to the central agent using back-propagation, and its messages preprocessed using a regularization unit) outperforms the other training strategies, learning a significantly better policy than local agents learning independently.

## 6.2 Multi-agent walker results

This section shows that HAMMER also works in a multi-agent task with continuous control and individual rewards. Figure 7 shows that HAMMERV3 agents perform considerably better than unaided independent local agents. Results here are averaged over eight independent trials, with an additional 2000-episode moving window to increase readability. HAMMER performing better in a domain like Multi-Agent Walker confirms the generalization of the approach

to continuous action spaces and different reward structures.<sup>1</sup>

The results for the two test domains show that (1) heterogeneous agents successfully learn messaging and enable multi-level coordination among the independent agents to enhance reinforcement learning using HAMMER approach, (2) HAMMER generalizes to heterogeneous local agents in the environment, (3) the approach works well in both discrete and continuous action spaces, and (4) the approach performed well with both individual rewards and global team rewards.

## 7 Analyzing the communication

Results discussed in the Sect. 6 show that by adding another agent for communicating messages, HAMMER agents achieve a higher reward. To detect whether communication is emerging at all, showing this improvement in reward could be a sufficient indicator. But it provides only a coarse measure of the agent's learned communication abilities. It is possible that agents may obtain a similar reward by coordinating using learned convention and not communication [35]. This means that communication could be an alternative way of optimization. It would be useful to further quantify the degree of communication as it can provide more insights into agent behavior and become useful for humans to monitor agents' behavior for fault detection, assessing performance, or even building trust. We now attempt to examine whether there is any semantic meaning or any influence of HAMMERV2-generated messages communicated to the local agents before their actions when performing the cooperative navigation task.

Detecting useful communication is not as simple as looking at the decrease in reward on removing the communication channel at test time because neural networks can be very sensitive to their input distribution [56]. A change in this distribution (for example setting the messages to 0) may cause the agents to fail, even if the messages contain no useful information. In the simplest case, we can measure the effect that HAMMER's messages have on the actions of the local agents. Hence, we start by visualizing the action trajectories of local agents and noting the messages communicated to them by the central agent. As an example, the first row of Fig. 8 shows an episode as a collection of typical sequences of events in the environment shown over time where local agents operate under the message-guidance of HAMMER. As expected, these HAMMER-

supported local agents cooperatively cover all the landmarks in the environment. We now experiment by overriding the original messages with varying human-generated values and communicate these new values to the local agents. Consequently, a substantial change in the action trajectories of all the local agents was noticed (Fig. 8). First, we manipulate the message values for only one local agent (Agent 1). On sending values close to -6 to Agent 1 throughout the episode, it navigated to landmark L1 leaving landmark L2 uncovered (second row of Fig. 8). Second, on modifying the message values to 3 for Agent 1 throughout the episode, it attempted to navigate to landmark L3 instead of landmark L2 (third row of Fig. 8). Similarly, manipulating the message values for Agent 2 and Agent 3 caused a change in the actions taken by the local agents in the episode. This confirms that HAMMER messages are indeed influencing the local agents' behavior. On carefully manipulating the values of these messages, we were able to navigate the local agents to desired landmarks (the fourth row of Fig. 8 shows an example of all agents navigating to landmark L1). These experiments were done on numerous episodes and for multiple independently trained HAMMER models, all of which brought us to the same inferences (with different message values).

In conclusion, our experiments confirm that communication by HAMMER influences local agents' actions. Moreover, these messages seem to have an effect further in the future than a single time step, and hence the language used by the central agent seems to be temporally extended. There are other aspects of the environment that HAMMER could have learned to signal about: sending a message to control the velocity of the local agents, to share an observation it has made, or to reveal the sequence of actions it or the other local agents have taken in the past. New metrics need to be developed to evaluate these possibilities and we leave it to the future for now.

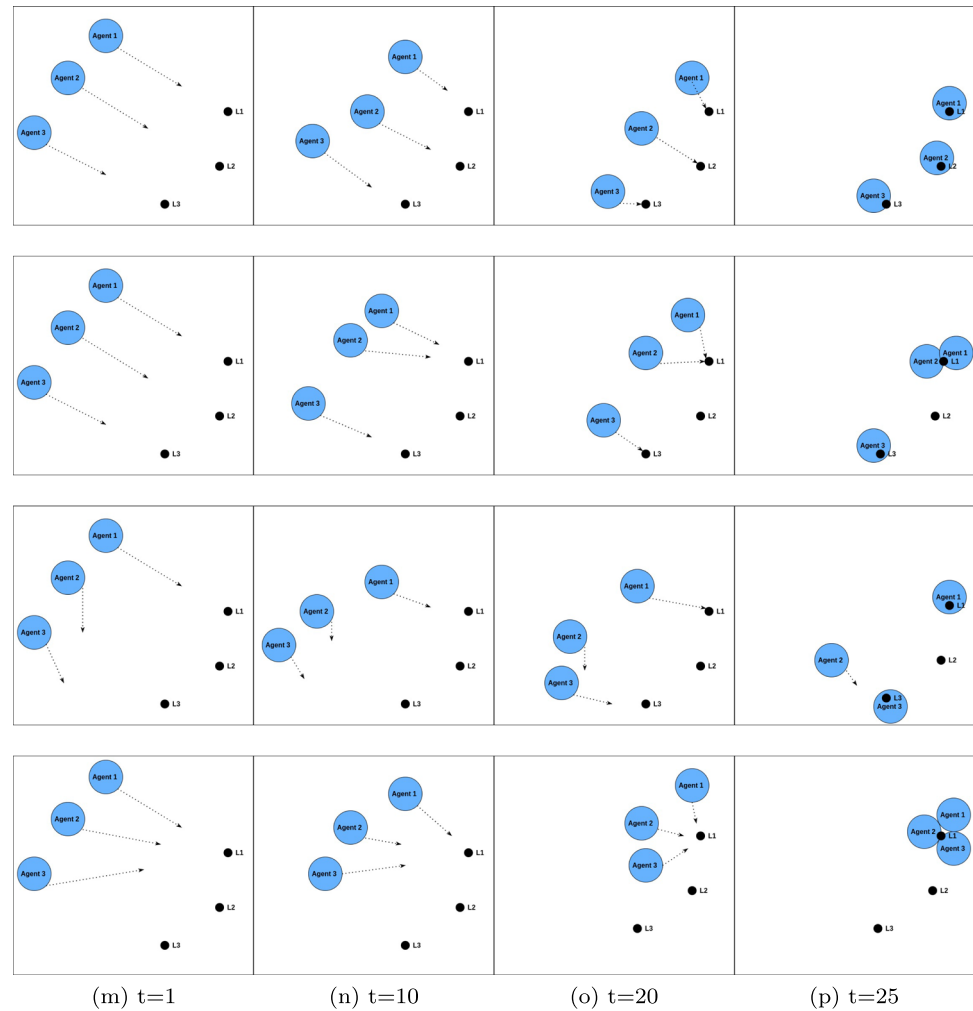
## 8 Analyzing HAMMER's scalability

Results discussed previously are limited to 3 agents in the environment. To see how well HAMMER scales with the number of agents, this section analyzes the learning performances of HAMMERV2 and HAMMERV3 in the cooperative navigation environment with a varying number of local agents.<sup>2</sup> Furthermore, for simpler experiments (discussed earlier in Sect. 6), we showed that both HAMMERV2 and HAMMERV3 outperform independent learners, however, the

<sup>1</sup> We also note here that although the curves show that HAMMER agents outperform IL, the learning curves do become quite flat. We speculate that because the multi-agent walker task is a difficult task, agents may require much larger actor-critic networks, larger messages in HAMMER, and further hyperparameter tuning to avoid such stagnancy.

<sup>2</sup> HAMMERV1 uses a weak feedback signal and it is expected not to be able to scale well with more number of agents. We thus focus on analyzing the scalability of the better performing HAMMERV2 and HAMMERV3.

**Fig. 8** HAMMER's Communication Analysis: A collection of typical sequences of events in the cooperative navigation environment shown over time. Each row represents a different experiment on a single episode (for demonstration) with varying messages given to trained HAMMER agents. Four checkpoints from a 25 time step long episode are shown here for brevity



best performing variant seems to depend on the complexity of the task at hand. This analysis also helps us identify one best solution instead of different variations whose performance vary depending on the settings.

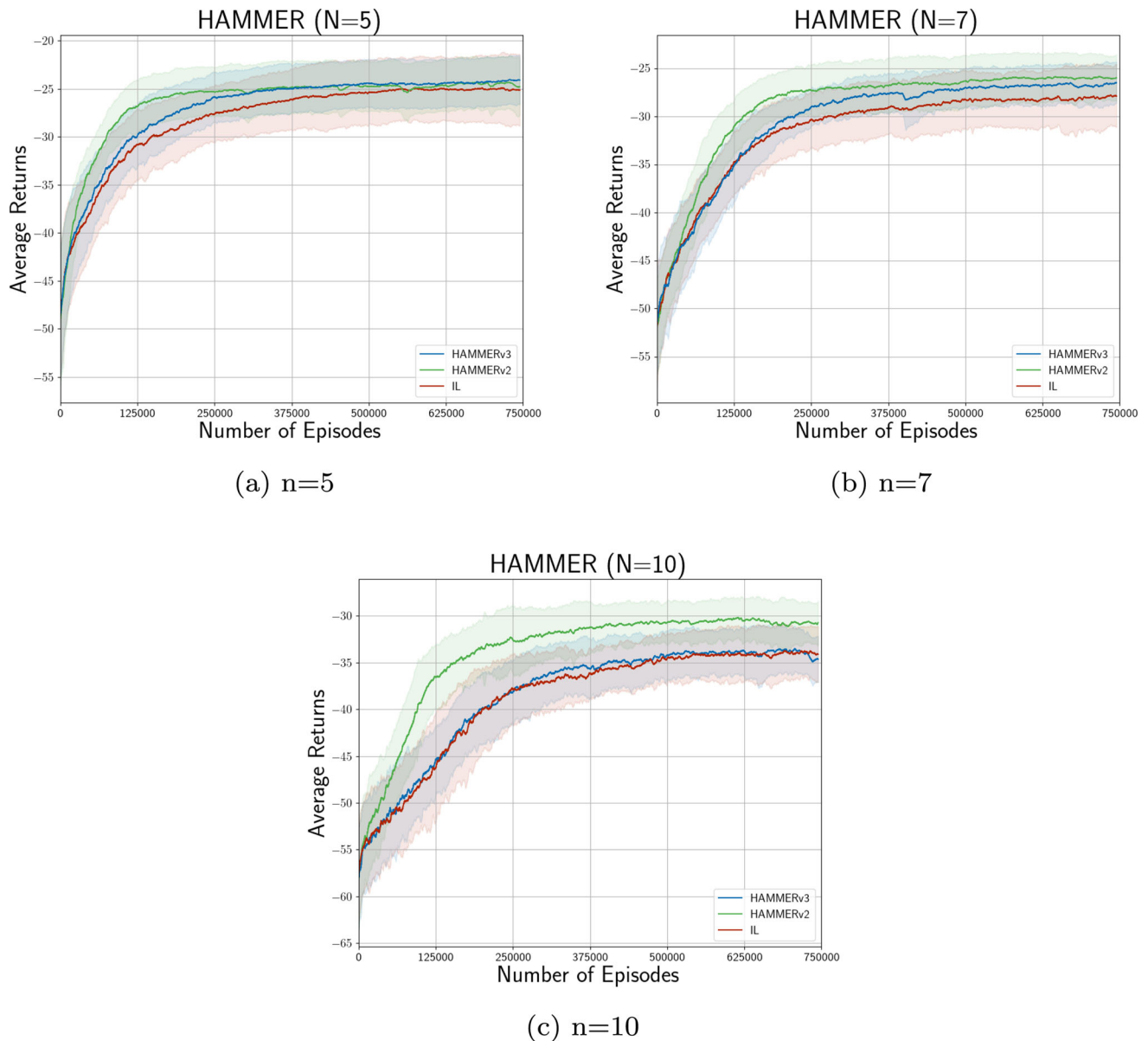
Figure 9 compares the learning performances of HAMMERV2 and HAMMERV3 in the cooperative navigation environment with 5, 7, and 10 local agents. HAMMERV2 can be identified as a clear winner with the increasing number of agents, however, HAMMERV3's performance begins to get closer to that of independent learners, especially in the case of  $n=10$  agents in the environment. We speculate that a discrete message of length one is not enough for HAMMER agents to encode enough information from the environment and use for efficient coordination. We suspect that larger messages of discrete bits can possibly help HAMMERV3's performance. We leave analyzing HAMMER's performance with varying lengths of message vectors to future work. Nevertheless, we propose HAMMERV2 as most appropriate choice out of the three variants that we discuss in this article for its high adaptability (speculatively due the availability of an infinite spectrum to encode information

via real-valued messages) and consistently better performance compared to independent learning agents in varying settings.

In conclusion, the curves show that HAMMER scales to a larger number of agents too. Moreover, through this analysis, we also propose HAMMERV2 as a generic solution that should be used for consistently superior performances across different settings with varying complexities.

## 9 Conclusion

This paper presented HAMMER, an approach used to achieve multi-level coordination among heterogeneous agents by learning useful messages from a global view of the environment or the multi-agent system. Using HAMMER, we addressed challenges like non-stationarity and inefficient coordination among agents in MARL by introducing a powerful all-seeing central agent in addition to the independent learners in the environment. Our results in two domains, cooperative navigation and multi-agent walker,



**Fig. 9** HAMMER's Scalability Analysis: Learning performances of HAMMER are compared with that of the Independent Learning (IL) baseline for  $n = 5, 7,$  and  $10$  agents in the cooperative navigation environment. Results here show that HAMMER scales to environments

with a larger number of agents. This analysis also helps identify HAMMERv2 as the best choice to be used for high adaptability and a consistent performance in varying environmental settings

showed that HAMMER can be generalized to discrete and continuous action spaces with both global team rewards and localized personal rewards. HAMMER also proved to generalize to heterogeneous local agents in the environment. We believe that the key reasons for the success of HAMMER were two-fold. First, we leveraged additional global information like global states, actions and rewards in the system. Second, centralizing observations have its own benefits, including helping to bypass problems like non-

stationarity in multi-agent systems and avoiding getting stuck in local optima [18, 64]. Several works related to ours (discussed earlier) demand powerful agents in the environment—ones that can transmit to other agents and/or are capable of modeling other agents in the system. This might not always be feasible, thus motivating HAMMER, where only one central agent needs to be powerful while the independent learners can be simple agents interacting with the environment based on their private local

observations augmented with learned messages. Warehouse management and traffic lights management are two example applications that could fit these assumptions well.

There are several directions for future work from here. Both the domains used in this work involved low-dimensional observation spaces and seem to offer substantial global coverage on combining local observations. HAMMER's results in multi-agent settings with tighter coupling and further complex interactions involved among agents such as in autonomous driving in SMARTS [68] or heterogeneous multi-agent battles in StarCraft [39] could help better appreciate the significance of the method. Also, though the results of HAMMER's three variants outperform independent learners in all scenarios, they seem to vary in their performance in different settings, and further investigation can be done to propose one generic solution. Additionally, more complex hierarchies could be used, such as by making several central agents available in the system. In this work, we performed an initial analysis of message vectors communicated by HAMMER, but additional work remains to better understand if and how HAMMER tailors messages for the local agents using its global observation. It would also be interesting to further study how RL learns to encode information in messages, understanding what the encoding means, and analyzing the effect of larger message lengths. It would also be interesting to test the scalability of HAMMER to a larger number of local agents in the environment. As part of the ablation studies for this work, we empirically justify the generalization of HAMMER in heterogeneous settings by blinding one of the agents. Further exciting heterogeneous settings, for instance, agents with different properties (such as velocity) can support assessing HAMMER's generalization.

Lastly, in our setting, communication is free—future work could consider the case where it was costly and attempt to trade off the number of messages sent with the learning speed of HAMMER. If the central agent had the ability to communicate small amounts of data occasionally, would it still be able to provide a significant improvement?<sup>3</sup>

## Appendix

### Implementation details

In both domains, two networks were used—one for the central agent and the second for the independent agents. PPO was used to update the decision-making policy for the local agents. For the actor and critic networks of these

individual agents, 2 fully-connected multi-layer perceptron layers were used to process the input layer and to produce the output from the hidden state. Three different variants were used to train the central agent. The first variant, HAMMERV1, had its own actor-critic network and learned its policy to communicate using PPO. In this case, the central agent's rewards were the same as those received by the local agents from the environment. In the second variant, HAMMERV2, the central agent used a multi-layer perceptron to output real-valued messages for the local agents and learned directly via the gradients passed back to it from the local agent network using backpropagation. In the final variant, HAMMERV3, instead of directly passing the messages to the local agent, we preprocessed the messages using a regularization unit, RU (as described earlier in Sect. 4). The central agent produces a vector of floating points scaled to  $[-1, 1]$  and is consistent across domains. A tanh activation function is used everywhere except for the local agent's output. Local agents in cooperative navigation use a soft-max over the output, corresponding to the probability of executing each of the five discrete actions. Local agents in the multi-agent walker task use 4-output nodes, each of which model a multivariate Gaussian distribution over torques.

Trials were run for 500,000 episodes in the cooperative navigation and 50,000 in the multi-agent walker environments. The training times of baselines and HAMMER were similar—12 h and 14 h, respectively, for 500,000 episodes on cooperative navigation and 14 h and 25 h, respectively, for 50,000 episodes in multi-agent walker. We set  $\gamma = 0.95$  for all experiments. After having tried multiple learning rates  $\{0.01, 0.001, 0.002, 0.005, 3 \times 10^{-4}, 1 \times 10^{-2}, 3 \times 10^{-3}\}$  and training batch sizes in PPO—we found  $3 \times 10^{-3}$  and 800, respectively, to work best for both the centralized agent and the independent agents in both domains. Five independent trials with different random seeds were performed on cooperative navigation environment, and eight for multi-agent walker, to establish statistical significance. The clip parameter for PPO was set to 0.2.

The independent learners follow the formerly proposed idea of allowing all of them to share the parameters of a single policy, hence enabling a single policy to be learned with experiences from all the agents simultaneously [11, 15]. This still ensures different behavior among agents because each of them receives different observations. Parameter sharing has been successfully applied in several other multi-agent deep reinforcement learning settings [13, 39, 42, 53, 55]. We also note here that in COIL, we do not allow sharing of parameters among the independent agents. The objective of COIL was to confirm that the central HAMMER agent is not simply forwarding

<sup>3</sup> As part of Nikunj Gupta's Master's Thesis titled "Fully Cooperative Multi-Agent Reinforcement Learning".

complete information about all agents to all of them. Hence, here the independent agents learn separate policies (i.e., by not sharing the same network parameters) using complete information about all the other agents in the environment. This experimental setup is inherently different from that used by other related works' implementation of typical "centralized training".

While we focus on having a shared network for independent agents and using PPO methods for policy learning, HAMMER can be implemented with other MADRL algorithms. Moreover, even though we test with a single central agent in this paper, it is entirely possible that multiple central agents could better assist independent learners. This is left to future works.

**Acknowledgements** This work commenced at Ericsson Research Laboratory Bangalore, and most of the follow-up work was done at the International Institute of Information Technology-Bangalore.<sup>3</sup> Part of this work has taken place in the Intelligent Robot Learning (IRL) Laboratory at the University of Alberta, which is supported in part by research grants from Alberta Innovates; the Alberta Machine Intelligence Institute (Amii); a Canada CIFAR AI Chair, Amii; Compute Canada; Huawei; Mitacs; and NSERC. We would like to thank Laura Petrich, Shahil Mawjee and anonymous reviewers for comments and suggestions on this paper.

## Declarations

**Conflict of interest** The authors declare that they have no conflict of interest.

**Data availability** Data sharing not applicable to this article as no datasets were generated or analyzed during the current study.

## References

- Albrecht SV, Stone P (2018) Autonomous agents modelling other agents: a comprehensive survey and open problems. *Artif Intell* 258:66–95
- Bowling M, Burch N, Johanson M, Tammelin O (2015) Heads-up limit hold'em poker is solved. *Science* 347(6218):145–149
- Busoniu L, Babuska R, De Schutter B (2008) A comprehensive survey of multiagent reinforcement learning. *IEEE Trans Syst Man Cybern Part C (Applications and Reviews)* 38(2):156–172
- Cao Y, Yu W, Ren W, Chen G (2012) An overview of recent progress in the study of distributed multi-agent coordination. *IEEE Trans Ind Inform* 9(1):427–438
- Castellini J, Oliehoek FA, Savani R, Whiteson S (2019) The representational capacity of action-value networks for multi-agent reinforcement learning. arXiv preprint [arXiv:1902.07497](https://arxiv.org/abs/1902.07497)
- Courbariaux M, Hubara I, Soudry D, El-Yaniv R, Bengio Y (2016) Binarized neural networks: Training deep neural networks with weights and activations constrained to  $\pm 1$  or  $\pm 1$ . arXiv preprint [arXiv:1602.02830](https://arxiv.org/abs/1602.02830)
- Dietterich TG (2000) Hierarchical reinforcement learning with the maxq value function decomposition. *J Artif Intell Res* 13:227–303
- Enright JJ, Wurman PR (2011) Optimization and coordinated autonomy in mobile fulfillment systems. In: Workshops at the twenty-fifth AAAI conference on artificial intelligence, Citeseer
- Farinelli A, Rogers A, Petcu A, Jennings NR (2008) Decentralised coordination of low-power embedded devices using the max-sum algorithm
- Foerster J, Assael IA, De Freitas N, Whiteson S (2016a) Learning to communicate with deep multi-agent reinforcement learning. In: Advances in neural information processing systems, pp 2137–2145
- Foerster J, Assael IA, de Freitas N, Whiteson S (2016b) Learning to communicate with deep multi-agent reinforcement learning. In: Lee DD, Sugiyama M, Luxburg UV, Guyon I, Garnett R (eds) Advances in Neural Information Processing Systems 29, Curran Associates, Inc., pp 2137–2145
- Foerster J, Farquhar G, Afouras T, Nardelli N, Whiteson S (2017a) Counterfactual multi-agent policy gradients. arXiv preprint [arXiv:1705.08926](https://arxiv.org/abs/1705.08926)
- Foerster J, Nardelli N, Farquhar G, Afouras T, Torr PH, Kohli P, Whiteson S (2017b) Stabilising experience replay for deep multi-agent reinforcement learning. arXiv preprint [arXiv:1702.08887](https://arxiv.org/abs/1702.08887)
- Fox D, Burgard W, Kruppa H, Thrun S (2000) A probabilistic approach to collaborative multi-robot localization. *Auton Robots* 8(3):325–344
- Gupta JK, Egorov M, Kochenderfer M (2017) Cooperative multi-agent control using deep reinforcement learning. In: International Conference on Autonomous Agents and Multiagent Systems, Springer, pp 66–83
- Hinton G, Salakhutdinov R (2011) Discovering binary codes for documents by learning deep generative models. *Topics Cogn Sci* 3(1):74–91
- Ito T, Zhang M, Robu V, Fatima S, Matsuo T, Yamaki H (2010) Innovations in agent-based complex automated negotiations, vol 319. Springer
- Johanson M, Waugh K, Bowling M, Zinkevich M (2011) Accelerating best response calculation in large extensive games. *IJCAI* 11:258–265
- Kim W, Cho M, Sung Y (2019) Message-dropout: An efficient training method for multi-agent deep reinforcement learning. Proceedings of the AAAI Conference on Artificial Intelligence 33:6079–6086
- Kumar S, Shah P, Hakkani-Tur D, Heck L (2017) Federated control with hierarchical multi-agent deep reinforcement learning. arXiv preprint [arXiv:1712.08266](https://arxiv.org/abs/1712.08266)
- Lanctot M, Zambaldi V, Gruslys A, Lazaridou A, Tuyls K, Pérolat J, Silver D, Graepel T (2017) A unified game-theoretic approach to multiagent reinforcement learning. In: Advances in neural information processing systems, pp 4190–4203
- Laurent GJ, Matignon L, Fort-Piat L et al (2011) The world of independent learners is not Markovian. *Int J Knowl Based Intell Eng Syst* 15(1):55–64
- Lazaridou A, Peysakhovich A, Baroni M (2016) Multi-agent cooperation and the emergence of (natural) language. arXiv preprint [arXiv:1612.07182](https://arxiv.org/abs/1612.07182)
- Leibo JZ, Zambaldi V, Lanctot M, Marecki J, Graepel T (2017) Multi-agent reinforcement learning in sequential social dilemmas. arXiv preprint [arXiv:1702.03037](https://arxiv.org/abs/1702.03037)
- Leibo JZ, Pérolat J, Hughes E, Wheelwright S, Marblestone AH, Duéñez-Guzmán E, Sunehag P, Dunning I, Graepel T (2018) Malthusian reinforcement learning. arXiv preprint [arXiv:1812.07019](https://arxiv.org/abs/1812.07019)
- Leibo JZ, Hughes E, Lanctot M, Graepel T (2019) Autocurricula and the emergence of innovation from social interaction: A

- manifesto for multi-agent intelligence research. arXiv preprint [arXiv:1903.00742](https://arxiv.org/abs/1903.00742)
27. Lillicrap TP, Hunt JJ, Pritzel A, Heess N, Erez T, Tassa Y, Silver D, Wierstra D (2015) Continuous control with deep reinforcement learning. arXiv preprint [arXiv:1509.02971](https://arxiv.org/abs/1509.02971)
  28. Liu M, Deng J, Xu M, Zhang X, Wang W (2017) Cooperative deep reinforcement learning for traffic signal control. In: The 7th International Workshop on Urban Computing (UrbComp 2018)
  29. Lowe R, Wu Y, Tamar A, Harb J, Abbeel P, Mordatch I (2017) Multi-agent actor-critic for mixed cooperative-competitive environments. *Neural Information Processing Systems (NIPS)*
  30. Ma J, Wu F (2020) Feudal multi-agent deep reinforcement learning for traffic signal control. In: Proceedings of the 19th International Conference on Autonomous Agents and Multiagent Systems (AAMAS), pp 816–824
  31. Matignon L, Jeanpierre L, Mouaddib AI (2012) Coordinated multi-robot exploration under communication constraints using decentralized markov decision processes. *AAAI 2012*:p2017-2023
  32. Matignon L, Laurent GJ, Le Fort-Piat N (2012b) Independent reinforcement learners in cooperative markov games: a survey regarding coordination problems
  33. Mnih V, Kavukcuoglu K, Silver D, Graves A, Antonoglou I, Wierstra D, Riedmiller M (2013) Playing atari with deep reinforcement learning. arXiv preprint [arXiv:1312.5602](https://arxiv.org/abs/1312.5602)
  34. Mnih V, Kavukcuoglu K, Silver D, Rusu AA, Veness J, Bellemare MG, Graves A, Riedmiller M, Fidjeland AK, Ostrovski G et al (2015) Human-level control through deep reinforcement learning. *Nature* 518(7540):529–533
  35. Mordatch I, Abbeel P (2018) Emergence of grounded compositional language in multi-agent populations. In: Proceedings of the AAAI Conference on Artificial Intelligence, vol 32
  36. Omidshafiei S, Kim DK, Liu M, Tesauo G, Riemer M, Amato C, Campbell M, How JP (2019) Learning to teach in cooperative multiagent reinforcement learning. *Proc AAAI Conf Artif Intell* 33:6128–6136
  37. Panait L, Luke S (2005) Cooperative multi-agent learning: the state of the art. *Auton Agents Multi-agent Syst* 11(3):387–434
  38. Parker J, Nunes E, Godoy J, Gini M (2016) Exploiting spatial locality and heterogeneity of agents for search and rescue teamwork. *J Field Robot* 33(7):877–900
  39. Peng P, Wen Y, Yang Y, Yuan Q, Tang Z, Long H, Wang J (2017) Multiagent bidirectionally-coordinated nets: Emergence of human-level coordination in learning to play starcraft combat games. arXiv preprint [arXiv:1703.10069](https://arxiv.org/abs/1703.10069)
  40. Pipattanasomporn M, Feroze H, Rahman S (2009) Multi-agent systems in a distributed smart grid: Design and implementation. In: 2009 IEEE/PES Power Systems Conference and Exposition, IEEE, pp 1–8
  41. Van der Pol E, Oliehoek FA (2016) Coordinated deep reinforcement learners for traffic light control. *Proceedings of learning, inference and control of multi-agent systems (at NIPS 2016)* 1
  42. Rashid T, Samvelyan M, De Witt CS, Farquhar G, Foerster J, Whiteson S (2018) Qmix: Monotonic value function factorisation for deep multi-agent reinforcement learning. arXiv preprint [arXiv:1803.11485](https://arxiv.org/abs/1803.11485)
  43. Rogers A, Farinelli A, Stranders R, Jennings NR (2011) Bounded approximate decentralised coordination via the max-sum algorithm. *Artif Intell* 175(2):730–759
  44. Samothrakis S, Lucas S, Runarsson T, Robles D (2012) Coevolving game-playing agents: measuring performance and intransitivities. *IEEE Trans Evolut Comput* 17(2):213–226
  45. Schulman J, Levine S, Abbeel P, Jordan M, Moritz P (2015) Trust region policy optimization. In: International conference on machine learning, pp 1889–1897
  46. Schulman J, Wolski F, Dhariwal P, Radford A, Klimov O (2017) Proximal policy optimization algorithms. arXiv preprint [arXiv:1707.06347](https://arxiv.org/abs/1707.06347)
  47. Seuken S, Zilberstein S (2012) Improved memory-bounded dynamic programming for decentralized pomdps. arXiv preprint [arXiv:1206.5295](https://arxiv.org/abs/1206.5295)
  48. Sheng J, Wang X, Jin B, Yan J, Li W, Chang TH, Wang J, Zha H (2020) Learning structured communication for multi-agent reinforcement learning. arXiv preprint [arXiv:2002.04235](https://arxiv.org/abs/2002.04235)
  49. Shoham Y, Powers R, Grenager T (2007) If multi-agent learning is the answer, what is the question? *Artif Intell* 171(7):365–377
  50. Silva FLD, Warnell G, Costa AHR, Stone P (2020) Agents teaching agents: a survey on inter-agent transfer learning. *Autonomous Agents and Multi-Agent Systems*
  51. Silver D, Huang A, Maddison CJ, Guez A, Sifre L, Van Den Driessche G, Schrittwieser J, Antonoglou I, Panneershelvam V, Lanctot M et al (2016) Mastering the game of go with deep neural networks and tree search. *Nature* 529(7587):484–489
  52. Srivastava N, Hinton G, Krizhevsky A, Sutskever I, Salakhutdinov R (2014) Dropout: a simple way to prevent neural networks from overfitting. *J Mach Learn Res* 15(1):1929–1958
  53. Sukhbaatar S, Fergus R, et al. (2016) Learning multiagent communication with backpropagation. In: Advances in neural information processing systems, pp 2244–2252
  54. Sunehag P, Lever G, Gruslys A, Czarniecki WM, Zambaldi V, Jaderberg M, Lanctot M, Sonnerat N, Leibo JZ, Tuyls K, et al. (2017) Value-decomposition networks for cooperative multi-agent learning. arXiv preprint [arXiv:1706.05296](https://arxiv.org/abs/1706.05296)
  55. Sunehag P, Lever G, Gruslys A, Czarniecki WM, Zambaldi VF, Jaderberg M, Lanctot M, Sonnerat N, Leibo JZ, Tuyls K, et al. (2018) Value-decomposition networks for cooperative multi-agent learning based on team reward. In: AAMAS, pp 2085–2087
  56. Szegedy C, Zaremba W, Sutskever I, Bruna J, Erhan D, Goodfellow I, Fergus R (2013) Intriguing properties of neural networks. arXiv preprint [arXiv:1312.6199](https://arxiv.org/abs/1312.6199)
  57. Tampuu A, Matiisen T, Kodelja D, Kuzovkin I, Korjus K, Aru J, Aru J, Vicente R (2017) Multiagent cooperation and competition with deep reinforcement learning. *PLoS one* 12(4):e0172395
  58. Tan M (1993) Multi-agent reinforcement learning: Independent vs. cooperative agents. In: Proceedings of the tenth international conference on machine learning, pp 330–337
  59. Tang H, Hao J, Lv T, Chen Y, Zhang Z, Jia H, Ren C, Zheng Y, Fan C, Wang L (2018) Hierarchical deep multiagent reinforcement learning. arXiv preprint [arXiv:1809.09332](https://arxiv.org/abs/1809.09332)
  60. Taylor ME, Carboni \* N, Fachantidis A, Vlahavas I, Torrey L, (2014) Reinforcement learning agents providing advice in complex video games. *Connection Science* 26(1):45–63. <https://doi.org/10.1080/09540091.2014.885279>
  61. Tesauo G (1995) Temporal difference learning and td-gammon. *Commun ACM* 38(3):58–68
  62. Tuyls K, Weiss G (2012) Multiagent learning: basics, challenges, and prospects. *Ai Magaz* 33(3):41–41
  63. Vezhnevets AS, Osindero S, Schaul T, Heess N, Jaderberg M, Silver D, Kavukcuoglu K (2017) Feudal networks for hierarchical reinforcement learning. arXiv preprint [arXiv:1703.01161](https://arxiv.org/abs/1703.01161)
  64. Whiteson S, Tanner B, Taylor ME, Stone P (2011) Protecting against evaluation overfitting in empirical reinforcement learning. In: 2011 IEEE symposium on adaptive dynamic programming and reinforcement learning (ADPRL), IEEE, pp 120–127
  65. Wunder M, Littman M, Stone M (2009) Communication, credibility and negotiation using a cognitive hierarchy model. *AAMAS Workshop, Citeseer* 19:73–80
  66. Ying W, Dayong S (2005) Multi-agent framework for third party logistics in e-commerce. *Expert Systems with Applications* 29(2):431–436

67. Zaïem MS, Bennequin E (2019) Learning to communicate in multi-agent reinforcement learning: A review. arXiv preprint [arXiv:1911.05438](https://arxiv.org/abs/1911.05438)
68. Zhou M, Luo J, Villela J, Yang Y, Rusu D, Miao J, Zhang W, Alban M, Fadakar I, Chen Z, et al. (2020) Smarts: Scalable multi-agent reinforcement learning training school for autonomous driving. arXiv preprint [arXiv:2010.09776](https://arxiv.org/abs/2010.09776)

**Publisher's Note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Springer Nature or its licensor (e.g. a society or other partner) holds exclusive rights to this article under a publishing agreement with the author(s) or other rightsholder(s); author self-archiving of the accepted manuscript version of this article is solely governed by the terms of such publishing agreement and applicable law.