

Learning to Communicate and Act in Cooperative Multiagent Systems using Hierarchical Reinforcement Learning

Mohammad Ghavamzadeh & Sridhar Mahadevan
Department of Computer Science
University of Massachusetts Amherst
Amherst, MA 01003-4601
mgh@cs.umass.edu & mahadeva@cs.umass.edu

Abstract

In this paper, we address the issue of rational communication behavior among autonomous agents. We extend our previously reported cooperative hierarchical reinforcement learning (HRL) algorithm to include communication decision and propose a new multiagent HRL algorithm, called COM-Cooperative HRL. In this algorithm, at specific levels of the hierarchy, called cooperation levels, a group of subtasks, in which coordination among agents has significant effect on the performance of the overall task, are defined as cooperative subtasks. Coordination skills among agents are learned faster by sharing information at cooperation levels, rather than the level of primitive actions. We add a communication level to the hierarchical decomposition of the problem, below each cooperation level. A communication action has a certain cost and is used by each agent to obtain the actions selected by the cooperative subtasks of the other agents. Before making a decision at a cooperative subtask, agents decide if it is worthwhile to perform a communication action in order to acquire the actions chosen by the cooperative subtasks of the other agents. Using this algorithm, agents learn a policy to balance the amount of communication needed for proper coordination, and communication cost. We demonstrate the efficacy of the COM-Cooperative HRL algorithm as well as the relation between communication cost and the learned communication policy, using a multiagent taxi domain.¹

1. Introduction

Cooperative multiagent learning studies algorithms for selecting actions for multiple agents coexisting in the same environment and working together to accomplish a task.

The reinforcement learning (RL) framework has been well-studied in cooperative multiagent domains [1, 2, 4, 16]. Multiagent RL has been recognized to be more challenging than single-agent RL for two main reasons: **1) curse of dimensionality**: the number of parameters to be learned increases dramatically with the number of agents, and **2) partial observability**: states and actions of other agents which are required for decision making by an agent are not fully observable and inter-agent communication is usually costly. Prior work in multiagent RL have addressed the *curse of dimensionality* in many different ways. One natural approach is to restrict the amount of information that is available to each agent and hope to maximize global payoff by solving local optimization problems [11, 14]. Another approach is to exploit the structure in the multiagent problem using factored value function architecture [7, 9]. This approach approximates the joint value function as a linear combination of local value functions, each of which relates only to the parts of the system controlled by a small number of agents. Factored value functions allow the agents to find a globally optimal joint-action using a message passing scheme. However, these works do not address the communication cost of their message passing strategy.

Almost all the above methods ignore this fact that an agent might not have free access to other agents' information which are required for its decision making. In general, the world is *partially observable* for agents in distributed multiagent domains. One way to address partial observability in these domains is to use communication to exchange information among agents. However, since communication is usually costly, in addition to its normal actions, each agent needs to make decision about communicating with other agents [12, 17]. The trade-off between the quality of solution and the communication cost is currently a very active area of research in multiagent learning and planning.

In our previous work [10], we introduced a different approach to address *curse of dimensionality* and *partial observability* in cooperative multiagent systems. The key

¹ First author of this paper is a student.

idea underlying the approach is that coordination skills are learned much more efficiently if the agents have a hierarchical representation of the task structure. Agents have only a local view of the overall state space, and learn joint abstract action-values by communicating with each other only the high-level subtasks that they are doing. It reduces the number of parameters to be learned. Furthermore, since high-level subtasks can take a long time to complete, communication is needed only fairly infrequently and this is a significant advantage over flat techniques. Although, the hierarchical RL (HRL) algorithm proposed in that work reduces the amount of communication required for coordination among agents, it does not address the issue of rational communication behavior, which is important when communication is costly. In this paper, we generalize our previous algorithm to include communication decisions and propose a new multiagent HRL algorithm, called *COM-Cooperative HRL*. In this algorithm, at specific levels of the hierarchy, called *cooperation levels*, we define a group of subtasks as *cooperative subtasks*. These are the subtasks, in which coordination among agents has significant effect on the performance of the overall task. Agents learn coordination skills by sharing information at *cooperation levels*, rather than the level of primitive actions. We add a communication level to the hierarchical decomposition of the problem, below each *cooperation level*. A communication action has a certain cost and can be used by each agent to obtain the actions selected by the *cooperative subtasks* of other agents. Using this algorithm, agents learn a policy to balance the amount of communication needed for proper coordination, and communication cost. We demonstrate the efficacy of the *COM-Cooperative HRL* algorithm as well as the relation between communication cost and communication policy, using a multiagent taxi domain.

2. Hierarchical Multiagent RL Framework

In this section, we introduce the hierarchical multiagent RL framework used in the multiagent HRL algorithm introduced in this paper. Our HRL framework builds upon the MAXQ value function decomposition [5], and the options model [15].

2.1. Hierarchical Task Decomposition

To illustrate our hierarchical multiagent RL framework and algorithm, we present a multiagent taxi problem, which will also be used in the experiments of this paper. Consider a 5-by-5 grid world inhabited by two taxis (*T1* and *T2*) shown in Figure 1. There are four specially designated locations in this domain, marked as B(lue), G(reen), R(ed) and Y(ellow). The task is continuing, passengers appear ac-

ording to a fixed passenger arrival rate² at these four locations and wish to be transported to one of the other locations chosen randomly. Taxis must go to the location of a passenger, pick up the passenger, go to its destination location, and put down the passenger there. The throughput of the system is measured in terms of the number of passengers dropped off at their destinations per 5000 steps.

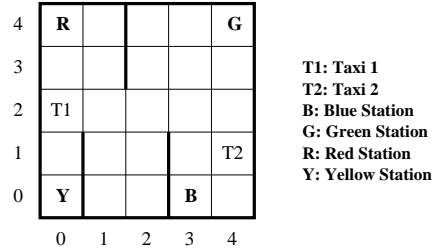


Figure 1. A multiagent taxi domain.

Hierarchical RL methods provide a general framework for scaling RL to problems with large state spaces by using the task structure to restrict the space of policies. In these methods, the overall task is decomposed into a collection of subtasks that are important for solving the problem. Each of these subtasks has a set of termination states, and terminates when one of its termination states is reached. Each primitive action (*North*, *West*, *South*, *East*, *Pickup* and *Put-down*) is a primitive subtask in this decomposition, such that it is always executable and it terminates immediately after execution. On the other hand, non-primitive subtasks such as *root* (the whole taxi problem), *Put*, *Get* B, G, R and Y, *Navigate to* B, G, R and Y, might take more than one time step to complete. After defining subtasks, we must indicate for each subtask, which other primitive or non-primitive subtasks it should employ to reach its goal. For example, navigation subtasks use four primitive actions *North*, *West*, *South* and *East*. *Put* uses four navigation subtasks plus one primitive action *Putdown*, and so on. All of this information is summarized by the task graph shown in Figure 2.

2.2. Temporal Abstraction using SMDPs

Hierarchical RL studies how lower level policies over subtasks or primitive actions can themselves be composed into higher level policies. Policies over primitive actions are semi-Markov when composed at the next level up, because they can take variable, stochastic amount of time to complete. Thus, semi-Markov decision processes (SMDPs) [8] have become the preferred language for modeling temporally extended actions. SMDPs extend the MDP model in

² Passenger arrival rate 10 indicates that on average, one passenger arrives at stations every 10 time steps.

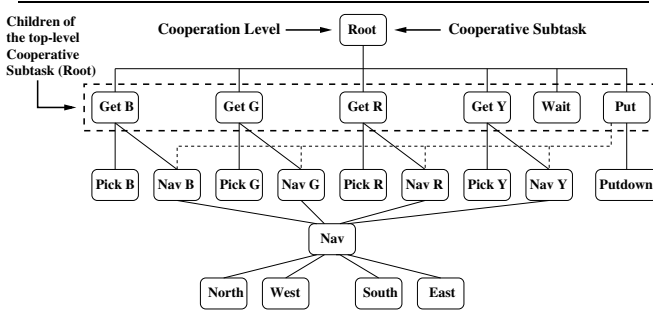


Figure 2. The task graph of the multiagent taxi domain.

several aspects. Decisions are only made at discrete points in time. State of the system may change continually between decisions, unlike MDPs where state changes are only due to the actions. Thus, the time between transitions may be several time units and can depend on the transition that is made. These transitions are at decision epochs only. Basically, the SMDP represents snapshots of the system at decision points, whereas the so-called *natural process* describes the evolution of the system over all times.

In this section, we extend SMDP to multiagent domain, when a team of agents controls the process, and introduce the multiagent SMDP (MSMDP) model. We assume agents are cooperative, i.e., each has the same reward and all maximize the same utility over an extended period of time. The individual actions of agents interact in that the effect of one agent’s action may depend on the actions taken by the others. When a group of agents perform temporally extended actions, these actions may not terminate at the same time. Therefore, unlike the multiagent extension of MDP (MMDP model [1]), the multiagent extension of SMDP is not straight forward.

Definition 1: A multiagent SMDP (MSMDP) is defined as a six tuple $(\alpha, \mathcal{S}, \mathcal{A}, P, R, \tau)$ as follows:

The set α is a finite collection of n agents, with each agent $j \in \alpha$ having a finite set A^j of individual actions. An element $\vec{a} = \langle a^1, \dots, a^n \rangle$ of the joint-action space $\mathcal{A} = \prod_{j=1}^n A^j$, represents the concurrent execution of actions a^j by each agent j . The components \mathcal{S} , R and P are as in an SMDP, the set of states of the system being controlled, the reward function mapping $\mathcal{S} \rightarrow \mathbb{R}$, and the state and action dependent multi-step transition probability function $P : \mathcal{S} \times \mathcal{N} \times \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ (where \mathcal{N} is the set of natural numbers). Since individual actions in a joint-action are temporally extended, they may not terminate at the same time. Therefore, the multi-step transition probability function P depends on how we define decision epochs, and as a result, depends on the termination

scheme τ that is used in the MSMDP model. Three termination strategies τ_{any} , τ_{all} and τ_{cont} for temporally extended joint-actions were investigated in [13]. In τ_{any} termination scheme, the next decision epoch is when the first action within the joint-action currently being executed terminates, where the rest of the actions that did not terminate are interrupted. When an agent reaches a decision point (drops off a passenger), all other agents interrupt their actions, next decision epoch occurs and a new joint-action is selected (taxi $T1$ chooses to pick up passenger at R and taxi $T2$ decides to pick up passenger at B). In τ_{all} termination scheme, the next decision epoch is the earliest time at which all the actions within the joint-action currently being executed have terminated. When an agent completes a subtask, it waits (takes *idle* action) until all other agents complete their current subtask. Then, next decision epoch happens and agents choose next joint-action together. In both these termination strategies, all agents make decision at every decision epoch. τ_{cont} termination scheme is similar to τ_{any} in the sense that the next decision epoch is when the first action within the joint-action currently being executed terminates. However, the other agents are not interrupted and only terminated agents select new actions. In this termination strategy, only a subset of agents choose action at each decision epoch. When an agent terminates a subtask, next decision epoch occurs only for that agent and it selects its next action given the information about the subtasks being performed by other agents. \square

The three termination strategies described above are the most common, but not the only termination schemes in cooperative multiagent activities. A wide range of termination strategies can be defined based on them. Of course, all these termination schemes are not appropriate for every multiagent task. We categorize termination strategies as *synchronous* and *asynchronous*. In *synchronous* schemes, such as τ_{any} and τ_{all} , all agents make decision at every decision epoch and therefore we need a centralized mechanism to synchronize agents at decision epochs. In *asynchronous* strategies, such as τ_{cont} , only a subset of agents make decision at each decision epoch. In this case, there is no need for a centralized mechanism to synchronize agents and decision making can take place in a decentralized fashion.

While SMDP theory provides the theoretical underpinnings of temporal abstraction by allowing for actions that take varying amounts of time, the SMDP model provides little in the way of concrete representational guidance which is critical from a computational point of view. In particular, the SMDP model does not specify how tasks can be broken up into subtasks, how to define policy for subtasks, how to decompose value function etc. We examine these issues in the rest of this section.

Mathematically, a task hierarchy such as the one in Fig-

ure 2 can be modeled by decomposing the overall task MDP M , into a finite set of subtasks $\{M_0, \dots, M_n\}$, where M_0 is the *root* task and solving it solves the MDP M .

Definition 2: Each *non-primitive* subtask i consists of five components $(S_i, I_i, T_i, A_i, R_i)$:

- S_i is the state space for subtask i and is described by those state variables that are relevant to subtask i . The range of the state variables describing S_i might be a subset of their range in S (the state space of the overall task MDP M).
- I_i is the initiation set for subtask i . Subtask i could start only in states belong to I_i .
- T_i is the set of terminal states for subtask i . Subtask i terminates when it reaches a state in T_i .
- A_i is the set of actions that can be performed to achieve subtask i . These actions can either be primitive actions from A (the set of primitive actions for MDP M), or they can be other subtasks.
- R_i is the reward function of subtask i . \square

The goal is to learn a policy for every subtask in the hierarchy. It gives us a policy for the overall task. This collection of policies is called a *hierarchical policy*.

Definition 3: A hierarchical policy π is a set with a policy for each of the subtasks in the hierarchy: $\pi = \{\pi_0, \dots, \pi_n\}$.

The hierarchical policy is executed using a stack discipline, similar to ordinary programming languages. Each subtask policy takes a state and returns the name of a primitive action to execute or a subtask to invoke. When a subtask is invoked, its name is pushed onto the stack and its policy is executed until it enters one of its terminal states. When a subtask terminates, its name is popped off the stack. Under a hierarchical policy π , we define a multi-step transition probability P_i^π for each subtask i in the hierarchy. $P_i^\pi(s', N|s)$ denotes the probability that action $\pi_i(s)$ will cause the system to transition from state s to state s' in N primitive steps.

The action-value function of executing subtask M_a under hierarchical policy π in state s in the context of parent task M_i , $Q^\pi(i, s, a)$, is decomposed into two parts: the value of subtask M_a in state s , $V^\pi(a, s)$, and the value of completing parent task M_i after invoking subtask M_a in state s , which is called the completion function $C^\pi(i, s, a)$ [5, 6]. The value function decomposition is recursively defined as:

$$Q^\pi(i, s, a) = V^\pi(a, s) + C^\pi(i, s, a) \quad (1)$$

$$V^\pi(i, s) = \begin{cases} Q^\pi(i, s, \pi_i(s)) & \text{if } i \text{ is non-primitive} \\ \sum_{s'} P(s'|s, i)R(s'|s, i) & \text{if } i \text{ is primitive} \end{cases}$$

2.3. Multiagent Setup

In our hierarchical multiagent model, we assume that there are n agents in the environment, cooperating with each other to accomplish a task. The task is decomposed by the designer of the system and its task graph is built, as described in Section 2.1. We also assume that agents are *homogeneous*, i.e., all agents are given the same task hierarchy.³ We define a group of subtasks as *cooperative subtasks* at specific levels of the hierarchy, called *cooperation levels*. The set of all *cooperative subtasks* at a *cooperation level* is called the *cooperation set* of that level. Agents actively coordinate only while making decision at *cooperative subtasks* and are ignorant about other agents at non-cooperative subtasks. Subtasks chosen to be cooperative are those in which coordination among agents is necessary and has significant effect on the performance of the overall task. We usually define *cooperative subtasks* at highest level(s) of the hierarchy. Coordination at high-level has two main advantages. First, it increases cooperation skills as agents do not get confused by low level details. Second, since high-level subtasks can take a long time to complete, communication among agents is needed only fairly infrequently. In this model, we specify policies for non-cooperative subtasks as single-agent policies, and policies for *cooperative subtasks* as joint policies.

Definition 4: Under a hierarchical policy π , each *non-cooperative subtask* i can be modeled by a SMDP consists of components (S_i, A_i, P_i^π, R_i) .

Definition 5: Under a hierarchical policy π , each *cooperative subtask* i located at the l th level of the hierarchy can be modeled by a MSMDP as follows:

α is the set of n agents in the team. We assume that agents have only local state information and ignore state of the other agents. Therefore, the state space S_i is defined as single-agent state space S_i (not joint state space). This is certainly an approximation but greatly simplifies the underlying multiagent RL problem. This approximation is based on the fact that an agent can get a rough idea of what state the other agents might be in just by knowing the high-level actions being performed by them. The action space is joint and is defined as $A_i = A_i \times (U_l)^{n-1}$, where $U_l = \bigcup_{k=1}^m A_k$ is the union of the action sets of all the l th level *cooperative subtasks*, and m is the cardinality of the l th level *cooperation set*. In the taxi domain, the set of agents is $\alpha = \{T1, T2\}$, *root* is defined

3 Studying the heterogeneous case where agents are given dissimilar decompositions of the overall task would be more challenging and is not the subject of this paper.

as *cooperative subtask*, and the highest level of the hierarchy as *cooperation level*, see Figure 2. Thus, *root* is the only member of the *cooperation set* at this level, and $U_{root} = A_{root} = \{GetB, GetG, GetR, GetY, Wait, Put\}$. The joint-action space for *root*, A_{root} , is specified as the cross product of the *root* action set, A_{root} , and U_{root} . Finally, since our goal is to design a decentralized multi-agent RL algorithm, we use τ_{cont} termination scheme for joint-action selection. \square

2.4. Incorporating Communication into the Model

Communication is a way for an agent to obtain local information of other agents by paying a certain communication cost. The *Cooperative HRL* algorithm described in our previous paper [10] works under three important assumptions, free, reliable, and instantaneous communication, i.e., communication cost is zero, no message is lost in the environment, and each agent has enough time to receive information of its teammates before taking its next action. Since communication is free, agents do not make decision about communication. Thus, as soon as an agent selects an action at a *cooperative subtask*, it broadcasts it to the team. Using this simple method, and the fact that communication is reliable and instantaneous, whenever an agent is about to choose an action at a *l*th level *cooperative subtask*, it knows the subtasks in U_l being performed by all its teammates.

However, communication is usually costly and unreliable in real-world problems. When communication is not free, it is no longer optimal for a team that agents always broadcast actions taken at their *cooperative subtasks* to their teammates. Therefore, agents must learn to optimally use communication by taking into account its long term return and its immediate cost. In this paper, we examine the case that communication is not free, but still assume that it is reliable and instantaneous. We extend the *Cooperative HRL* algorithm to include communication decision making and propose a new algorithm, called *COM-Cooperative HRL*. In *COM-Cooperative HRL*, we add a communication level to the task graph of the problem, below each *cooperation level*, as shown in Figure 3 for the taxi domain. When an agent is going to select an action at a *l*th level *cooperative subtask*, it first decides whether to communicate (takes *communicate* action) with other agents to acquire their actions in U_l , or takes *not-communicate* action and selects its action without new information about its teammates. The goal of our algorithm is to learn a hierarchical policy (a set of policies for all subtasks including the communication subtasks) to maximize the team utility given the communication cost. We illustrate the algorithm in more details in the next section.

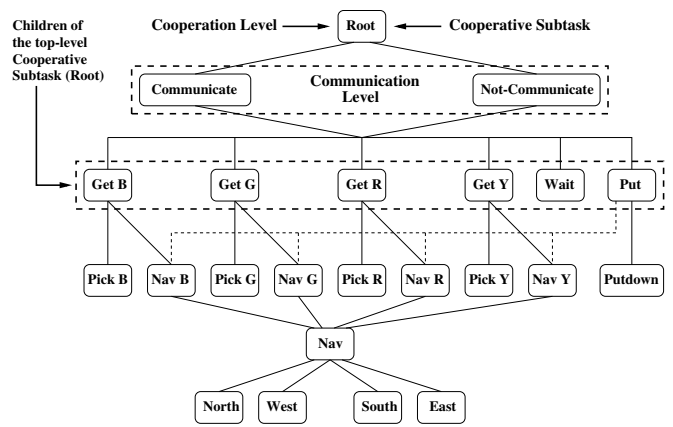


Figure 3. The task graph of the multiagent taxi domain.

3. Cooperative HRL Algorithm with Communication (COM-Cooperative HRL)

In *COM-Cooperative HRL*, agents decide about communication by comparing $Q(\text{Parent}(\text{NotCom}), s, \text{NotCom})$ with $Q(\text{Parent}(\text{Com}), s, \text{Com}) + \text{ComCost}$. If agent j decides not to communicate, it chooses action like a selfish agent by using its action-value function $Q^j(\text{NotCom}, s, a)$, where $a \in \text{Children}(\text{NotCom})$. When agent j decides to communicate, it acquires the actions being executed by all other agents in U_l and then uses its joint-action-value function $Q^j(\text{Com}, s, a^1, \dots, a^{j-1}, a^{j+1}, \dots, a^n, a)$, where $a \in \text{Children}(\text{Com})$ to select its next action in U_l . For instance, in the taxi domain, when taxi $T1$ drops off a passenger and is going to pick a new one, it should first decide whether to communicate with taxi $T2$ in order to acquire its action in U_{root} . To make communication decision, $T1$ compares $Q^1(\text{Root}, s, \text{NotCom})$ with $Q^1(\text{Root}, s, \text{Com}) + \text{ComCost}$. If it chooses not to communicate, it selects its action using $Q^1(\text{NotCom}, s, a)$, where $a \in U_{root}$. If it decides to communicate, after acquiring the $T2$'s action in U_{root} , a^{T2} , it selects its action using $Q^1(\text{Com}, s, a^{T2}, a)$. We can make the model more complicated by making decision for communication with each individual agent. In this case, the number of communication actions would be $C_{n-1}^1 + C_{n-1}^2 + \dots + C_{n-1}^{n-1}$, where C_p^q is the number of distinct combinations selecting q out of p agents. For instance, in a three-agent case, communication actions for agent 1 would be *Com with agent 2*, *Com with agent 3* and *Com with both agents 2 and 3*. It increases the number of communication actions and therefore the number of parameters to be learned. However, there are methods to reduce the number of communication actions in real-world applications. For instance, we can cluster agents in different groups based on their role in the

team and assume each group as a single entity to communicate with. It reduces n from the number of agents to the number of groups.

In *COM-Cooperative HRL* algorithm, *Communicate* subtasks are configured to store joint completion function values. The joint completion function for agent j , $C^j(Com, s, a^1, \dots, a^{j-1}, a^{j+1}, \dots, a^n, a^j)$ is defined as the expected discounted reward of completing subtask a^j by agent j in the context of the parent task i when other agents performing subtasks $a^i, \forall i \in \{1, \dots, n\}, i \neq j$. In the taxi domain, if taxi $T1$ communicates with taxi $T2$, its value function decomposition would be

$$Q^1(Com, s, GetR, GetB) = V^1(GetB, s) + C^1(Com, s, GetR, GetB)$$

which represents the value of $T1$ performing subtask $GetB$, when $T2$ is executing subtask $GetR$. Note that this value is decomposed into the value of the $GetB$ subtask and the completion cost of the remainder of the overall task. If $T1$ does not communicate with $T2$, its value function decomposition would be

$$Q^1(NotCom, s, GetB) = V^1(GetB, s) + C^1(NotCom, s, GetB)$$

which represents the value of $T1$ performing subtask $GetB$, regardless of the action being executed by $T2$.

The V and C values are learned through a standard temporal-difference learning method, based on sample trajectories. Since subtasks are temporally extended in time, the update rules are based on the SMDP model (see [5, 6] for details). Completion function and joint completion function values for an action in U_i are updated when this action is taken under *Not-Communicate* and *Communicate* subtasks respectively. In the later case, the other agents' actions in U_i are known as a result of communication and are used to update the joint completion function values.

4. Experimental Results

In this section, we demonstrate the performance of the *COM-Cooperative HRL* algorithm using the multiagent taxi problem described in Section 2.1. We also investigate the relation between communication policy and communication cost in this domain.

The state variables in this task are locations of taxis $T1$ and $T2$ (25 values each), status of taxis (2 values each, full or empty), status of stations B, G, R, Y (2 values each, full or empty), destination of stations (4 values each, one of the other three stations or without destination, which happens when the station is empty), destination of taxis (5 values each, one of the four stations or without destination, which is when taxi is empty). Thus, in the multiagent flat case, the size of the state space would grow to 256×10^6 . The size of the Q table is this number multiplies by 10, the

number of primitive actions, 256×10^7 . In the hierarchical selfish case (where each agent acts independently without communicating with other agents), using state abstraction and the fact that each agent stores only its own state variables, the number of the C and V values to be learned is reduced to $2 \times 135,895 = 271,790$, which is 135,895 values for each agent. In the hierarchical cooperative case without communication action, the number of values to be learned would be $2 \times 729,815 = 1,459,630$. Finally in the hierarchical cooperative case with communication action, this number would be $2 \times 934,615 = 1,869,230$. All the experiments in this section were conducted five times and the results averaged.

Figure 4 shows the throughput of the system for four algorithms, single-agent HRL, selfish multiagent HRL, *Cooperative HRL* and *COM-Cooperative HRL* when communication cost is zero. The *Cooperative HRL* and *COM-Cooperative HRL* algorithms use the task graphs in Figures 2 and 3 respectively. As seen in Figure 4, *Cooperative HRL* and *COM-Cooperative HRL* with $ComCost = 0$ have better throughput than the selfish multiagent HRL and single-agent HRL. The *COM-Cooperative HRL* learns slower than the *Cooperative HRL*, due to the more parameters to be learned in this model. However, it eventually converges to the same performance as *Cooperative HRL*.

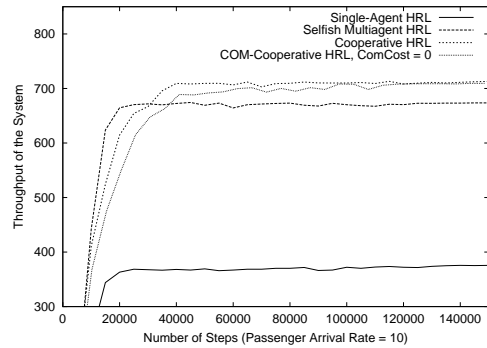


Figure 4. This figure shows the better throughput of *Cooperative HRL* and *COM-Cooperative HRL* with $ComCost = 0$, vs. selfish multiagent HRL and single-agent HRL.

Figure 5 compares the above four algorithms in terms of the average waiting time per passenger. This figure also demonstrates that *Cooperative HRL* and *COM-Cooperative HRL* with $ComCost = 0$, have less average waiting time per passenger than selfish multiagent HRL and single-agent HRL. Like the previous experiment, *COM-Cooperative HRL* eventually converges to the same performance as *Cooperative HRL*, however it is slower due to the more parameters to be learned.

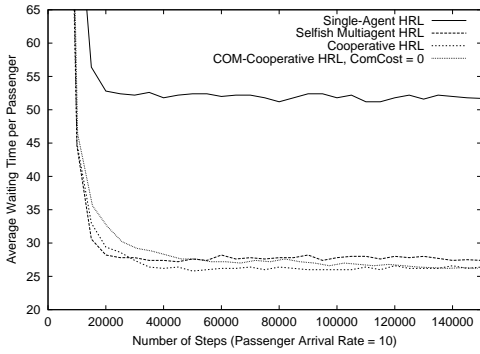


Figure 5. This figure shows that the average waiting time per passenger in *Cooperative HRL* and *COM-Cooperative HRL* with $ComCost = 0$, is less than selfish multiagent HRL and single-agent HRL.

Figure 6 compares the average waiting time per passenger for multiagent selfish HRL and *COM-Cooperative HRL* with $ComCost = 0$, for three different passenger arrival rates (5, 10 and 20). It demonstrates that as the passenger arrival rate becomes smaller, the coordination among taxis becomes more important. When taxis do not coordinate, there is a possibility that both taxis go to the same station. In this case, the first taxi picks up the passenger and the other one returns empty. This case can be avoided by incorporating coordination into the system. However, when the passenger arrival rate is high, there is a chance that a new passenger arrives after the first taxi picked up the previous passenger and before the second taxi reaches the station. This passenger will be picked up by the second taxi. In this case, coordination would not be as crucial as the case when the passenger arrival rate is low.

Figure 7 demonstrates the relationship between the communication policy and communication cost. These two figures show the throughput and the average waiting time per passenger for selfish multiagent HRL and *COM-Cooperative HRL* when communication cost equals 0, 1, 5, 10. In both figures, as the communication cost increases, the performance of the *COM-Cooperative HRL* becomes closer to the selfish multiagent HRL. It indicates that when communication is expensive, agents learn not to communicate and the multiagent system becomes selfish.

5. Conclusion and Future Work

In this paper, we study learning to communicate and act in cooperative multiagent systems using hierarchical reinforcement learning (HRL). The use of hierarchy speeds up learning in multiagent domains by making it possible to

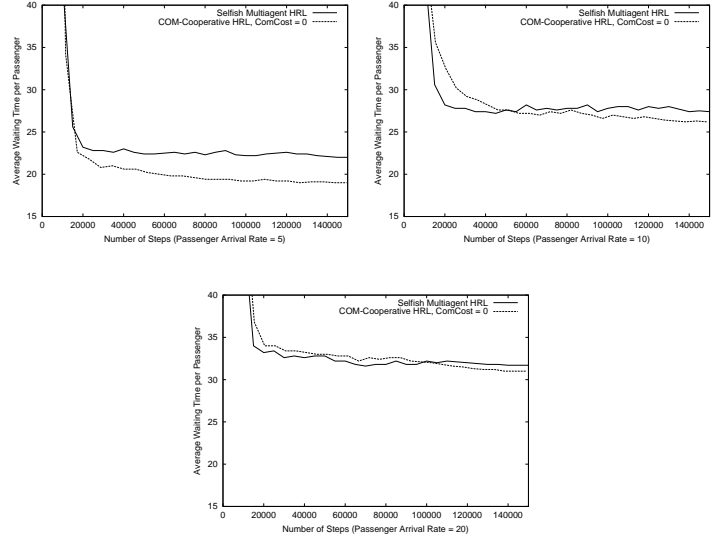


Figure 6. This figure compares the average waiting time per passenger for multiagent selfish HRL and *COM-Cooperative HRL* with $ComCost = 0$, for three different passenger arrival rates (5, 10 and 20). It shows that coordination among taxis becomes more crucial as the passenger arrival rate becomes smaller.

learn coordination skills at the level of subtasks instead of primitive actions. We introduce a new cooperative multiagent HRL algorithm, called *COM-Cooperative HRL*, by extending our previously reported algorithm [10] to include decision making about communication with other agents. In *COM-Cooperative HRL*, we define a group of subtasks as *cooperative subtasks* at specific levels of the hierarchy, called *cooperation levels*. These are subtasks in which coordination among agents is necessary. Each agent learns joint-action-values by communicating with its teammates at *cooperative subtasks*, and is unaware of them at other subtasks. We add a communication level to the task hierarchy, below each *cooperation level*. Before selecting an action at a *cooperative subtask*, agents have to decide if it is worthwhile to communicate with other agents in order to acquire the actions taken by their *cooperative subtasks*. It allows agents to learn a communication policy to balance the amount of communication for proper coordination, and communication cost. We study the empirical performance of the *COM-Cooperative HRL* algorithm in a multiagent taxi problem. We also investigate the relation between communication cost and communication policy in this domain.

A number of extensions would be useful, from studying the scenario where agents are heterogeneous, to recognizing the high-level subtasks being performed by other agents

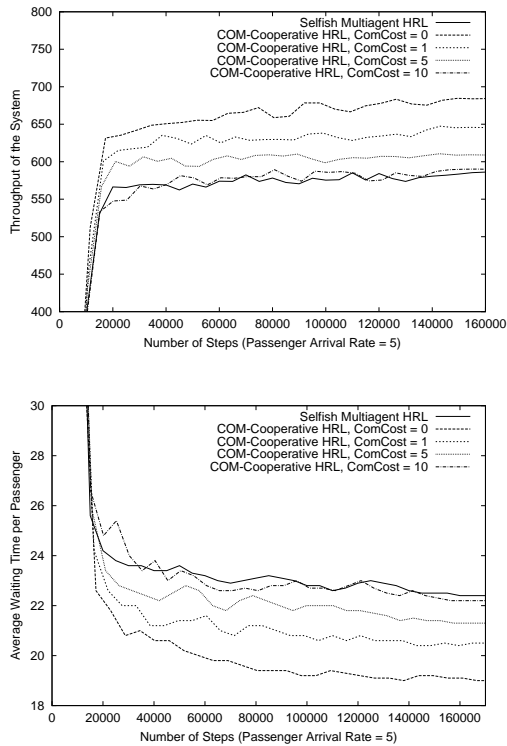


Figure 7. This figure shows the throughputs (top) and the average waiting time per passenger (bottom) for selfish multiagent HRL and COM-Cooperative HRL when communication cost equals 0, 1, 5 and 10.

using a history of observations instead of direct communication. In the later case, we assume that each agent can observe its teammates and uses its observations to extract their high-level subtasks [3]. Good examples for this approach are sport games such as soccer, football or basketball, in which players often extract the strategy being performed by their teammates, using recent observations instead of direct communication. It is obvious that many other manufacturing and robotics problems can benefit from this algorithm. We are currently applying *COM-Cooperative HRL* to the complex four-agent AGV scheduling problem used in experiments of our previous paper [10]. Combining these multiagent algorithms with function approximation and factored action models, which makes them appropriate for continuous state problems, is also an important area of research. The success of the proposed algorithms depends on providing them with a good initial hierarchical task decomposition. Therefore, deriving abstractions automatically is an essential problem to study. Finally, studying those communication features that have not been considered in our model, such as message delay and probability of loss, is another fundamental problem that needs to be addressed.

References

- [1] C. Boutilier. Sequential optimality and coordination in multi-agent systems. In *Proceedings of the Sixteenth International Joint Conference on Artificial Intelligence (IJCAI)*, 1999.
- [2] M. Bowling and M. Veloso. Multiagent learning using a variable learning rate. *Artificial Intelligence*, 136:215–250, 2002.
- [3] H. Bui, S. Venkatesh, and G. West. Policy recognition in the abstract hidden markov models. *Journal of Artificial Intelligence Research*, 17:451–499, 2002.
- [4] R. Crites and A. Barto. Elevator group control using multiple reinforcement learning agents. *Machine Learning*, 33:235–262, 1998.
- [5] T. Dietterich. Hierarchical reinforcement learning with the MAXQ value function decomposition. *Journal of Artificial Intelligence Research*, 13:227–303, 2000.
- [6] M. Ghavamzadeh and S. Mahadevan. Hierarchical multi-agent reinforcement learning. *UMASS Computer Science Technical Report*, 2004.
- [7] C. Guestrin, M. Lagoudakis, and R. Parr. Coordinated reinforcement learning. In *Proceedings of the Nineteenth International Conference on Machine Learning*, 2002.
- [8] R. Howard. *Dynamic Probabilistic Systems: Semi-Markov and Decision Processes*. John Wiley and Sons., 1971.
- [9] M. Lagoudakis and R. Parr. Learning in zero-sum team markov games using factored value functions. In *Neural Information Processing Systems (NIPS)*, 2002.
- [10] R. Makar, S. Mahadevan, and M. Ghavamzadeh. Hierarchical multi-agent reinforcement learning. In *Proceedings of the Fifth International Conference on Autonomous Agents*, 2001.
- [11] L. Peshkin, K. Kim, N. Meuleau, and L. Kaelbling. Learning to cooperate via policy search. In *Proceedings of the Sixteenth International Conference on Uncertainty in Artificial Intelligence (UAI)*, 2000.
- [12] D. Pynadath and M. Tambe. The communicative multiagent team decision problem: Analyzing teamwork theories and models. *Journal of Artificial Intelligence Research (JAIR)*, 16:389–426, 2002.
- [13] K. Rohanimanesh and S. Mahadevan. Learning to take concurrent actions. In *Proceedings of the Sixteenth Annual Conference on Neural Information Processing Systems*, 2002.
- [14] J. Schneider, W. Wong, A. Moore, and M. Riedmiller. Distributed value functions. In *Proceedings of the Sixteenth International Conference on Machine Learning (ICML)*, 1999.
- [15] R. Sutton, D. Precup, and S. Singh. Between MDPs and Semi-MDPs: A framework for temporal abstraction in reinforcement learning. *Artificial Intelligence*, 112:181–211, 1999.
- [16] M. Tan. Multi-agent reinforcement learning: Independent vs. cooperative agents. In *Proceedings of the Tenth International Conference on Machine Learning*, 1993.
- [17] P. Xuan, V. Lesser, and S. Zilberstein. Communication decisions in multi-agent cooperation: Model and experiments. In *Proceedings of the Fifth International Conference on Autonomous Agents*, 2001.