



Recent Advances in Hierarchical Reinforcement Learning

ANDREW G. BARTO

*Autonomous Learning Laboratory, Department of Computer Science, University of Massachusetts,
Amherst MA 01003*

SRIDHAR MAHADEVAN

*Autonomous Learning Laboratory Department of Computer Science, University of Massachusetts,
Amherst MA 01003*

Abstract. Reinforcement learning is bedeviled by the curse of dimensionality: the number of parameters to be learned grows exponentially with the size of any compact encoding of a state. Recent attempts to combat the curse of dimensionality have turned to principled ways of exploiting temporal abstraction, where decisions are not required at each step, but rather invoke the execution of temporally-extended activities which follow their own policies until termination. This leads naturally to hierarchical control architectures and associated learning algorithms. We review several approaches to temporal abstraction and hierarchical organization that machine learning researchers have recently developed. Common to these approaches is a reliance on the theory of semi-Markov decision processes, which we emphasize in our review. We then discuss extensions of these ideas to concurrent activities, multiagent coordination, and hierarchical memory for addressing partial observability. Concluding remarks address open challenges facing the further development of reinforcement learning in a hierarchical setting.

1. Introduction

Reinforcement learning (RL) (Bertsekas and Tsitsiklis, 1996; Sutton and Barto, 1998) is an active area of machine learning research that is also receiving attention from the fields of decision theory, operations research, and control engineering. RL algorithms address the problem of how a behaving agent can learn to approximate an optimal behavioral strategy while interacting directly with its environment. In control terms, this involves the on-line approximation of solutions to stochastic optimal control problems, usually under conditions of incomplete knowledge of the system being controlled. Most RL algorithms adapt standard methods of stochastic dynamic programming (DP) so that they can be used on-line for problems with large state spaces. By focusing computational effort along behavioral trajectories and by using function approximation methods for accumulating value function information, RL algorithms have produced good results on problems that pose significant challenges for standard methods (e.g., Crites and Barto, 1998 and Tesauro, 1992). However, current RL methods by no means completely circumvent the curse of dimensionality: the exponential growth of the number of parameters to be learned with the size of any compact encoding of system state. Recent attempts to combat the curse of dimensionality have turned to principled ways of exploiting temporal abstraction, where

decisions are not required at each step, but rather invoke the execution of temporally-extended activities which follow their own policies until termination. This leads naturally to hierarchical control architectures and learning algorithms.

In this article we review several related approaches to temporal abstraction and hierarchical control that have been developed by machine learning researchers, and we then discuss several extensions to these ideas and the open challenges facing the further development of RL in a hierarchical setting. Despite the fact that research in this area began only recently within the machine learning community, we cannot provide a survey of all the related literature, which is already quite extensive. We do attempt, however, to provide enough information so that interested readers can probe the relevant topics more deeply. Another important goal that we do not attempt is to thoroughly relate machine learning research on these topics to the extensive literature in systems and control engineering on hierarchical and multilayer control, hybrid systems, and other closely related topics. The most obvious parallels with some of these approaches have not escaped us, but a careful rapprochement is beyond the scope of this article. We largely adhere to notation and terminology typical of that used in the machine learning community. There is much to be gained from integrating perspectives from these different groups of researchers, and we hope that this article will contribute to the required dialog.

After brief introductions to Markov and semi-Markov decision processes, we introduce the basic ideas of RL, and then we review three approaches to hierarchical RL: the *options* formalism of Sutton et al. (1999), the *hierarchies of abstract machines* (HAMs) approach of Parr (1998) and Parr and Russell (1998), and the MAXQ framework of Dietterich (2000). Although these approaches were developed relatively independently, they have many elements in common. In particular, they all rely on the theory of semi-Markov decision processes to provide a formal basis. Although we take advantage of these commonalities in our exposition, a thorough integration of these approaches awaits a future paper. We also cannot do full justice to these approaches, although we do attempt to provide enough detail to make the approaches concrete; the reader is directed to the original papers for detailed treatments. In the final sections, we briefly present extensions of these ideas that focus on work conducted in our laboratory on concurrent activities, multiagent coordination, and hierarchical memory for addressing partial observability. Concluding remarks address open challenges facing the further development of reinforcement learning in a hierarchical setting.

2. Markov and Semi-Markov Decision Processes

Most RL research is based on the formalism of Markov decision processes (MDPs). Although RL is by no means restricted to MDPs, this discrete-time, countable (in fact, usually finite) state and action formalism provides the simplest framework in which to study basic algorithms and their properties. Here we briefly describe this well-known framework, with a few twists characteristic of how it is used in RL research; additional details can be found in many references (e.g., Bertsekas, 1987, Bertsekas and Tsitsiklis, 1996, Puterman, 1994, Ross, 1983 and Sutton and Barto, 1998). A finite MDP models the following type of problem. At each stage in a sequence of stages, an agent (the controller)

observes a system's state s , contained in a finite set S , and executes an action (control) a selected from a finite, non-empty set, A_s , of admissible actions. The agent receives an immediate reward having expected value $R(s, a)$, and the state at the next stage is s' with probability $P(s'|s, a)$. The expected immediate rewards, $R(s, a)$, and the state transition probabilities, $P(s'|s, a)$, $s, s' \in S$, $a \in A_s$, together comprise what RL researchers often call the *one-step model* of action a .

A (stationary, stochastic) policy $\pi : S \times \bigcup_{s \in S} A_s \rightarrow [0, 1]$, with $\pi(s, a) = 0$ for $a \notin A_s$, specifies that the agent executes action $a \in A_s$ with probability $\pi(s, a)$ whenever it observes state s . For any policy π and $s \in S$, $V^\pi(s)$ denotes the *expected infinite-horizon discounted return* from s given that the agent uses policy π . Letting s_t and r_{t+1} denote the state at stage t and the immediate reward for acting at stage t ,¹ this is defined as:

$$V^\pi(s) = E\{r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \cdots | s_t = s, \pi\}$$

where γ , $0 \leq \gamma < 1$, is a discount factor. $V^\pi(s)$ is the *value* of s under π , and is V^π the *value function* corresponding to π .

This is a finite, infinite-horizon, discounted MDP. The objective is to find an optimal policy, i.e., a policy, π^* , that maximizes the value of each state. The unique *optimal value function*, V^* , is the value function corresponding to any of the optimal policies. Most RL research addresses discounted MDPs since they comprise the simplest class of MDPs, and here we restrict attention to discounted problems. However, RL algorithms have also been developed for MDPs with other definitions of return, such as average reward MDPs (Mahadevan, 1996; Schwartz, 1993).

Playing important roles in many RL algorithms are *action-value functions*, which assign values to admissible state-action pairs. Given a policy π , the value of (s, a) , $a \in A_s$, denoted $Q^\pi(s, a)$, is the expected infinite-horizon discounted return for executing a in state s and thereafter following π :

$$Q^\pi(s, a) = E\{r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \cdots | s_t = s, a_t = a, \pi\} \quad (1)$$

The *optimal action-value function*, Q^* , assigns to each admissible state-action pair (s, a) the expected infinite-horizon discounted return for executing a in s and thereafter following an optimal policy. Action-value functions for other definitions of return are defined analogously.

Dynamic Programming (DP) algorithms exploit the fact that value functions satisfy various *Bellman equations*, such as:

$$V^\pi(s) = \sum_{a \in A_s} \pi(s, a) \left[R(s, a) + \gamma \sum_{s'} P(s'|s, a) V^\pi(s') \right]$$

and

$$V^*(s) = \max_{a \in A_s} \left[R(s, a) + \gamma \sum_{s'} P(s'|s, a) V^*(s') \right] \quad (2)$$

for all $s \in S$. Analogous equations exist for Q^π and Q^* . For example, the Bellman equation for Q^* is:

$$Q^*(s, a) = R(s, a) + \gamma \sum_{s'} P(s'|s, a) \max_{a' \in A_{s'}} Q^*(s', a') \quad (3)$$

for all $s \in S, a \in A_s$.

As an example DP algorithm, consider value iteration, which successively approximates V^* as follows. At each iteration k , it updates an approximation V_k of V^* by applying the following operation for each state s :

$$V_{k+1}(s) = \max_{a \in A_s} \left[R(s, a) + \gamma \sum_{s'} P(s'|s, a) V_k(s') \right] \quad (4)$$

We call this operation a *backup* because it updates a state's value by transferring to it information about the approximate values of its possible successor states. Applying this backup operation once to each state is often called a *sweep*. Starting with an arbitrary initial function V_0 , the sequence $\{V_k\}$ produced by repeated sweeps converges to V^* . A similar algorithm exists for successively approximating Q^* using the following backup:

$$Q_{k+1}(s, a) = R(s, a) + \gamma \sum_{s' \in S} P(s'|s, a) \max_{a' \in A_{s'}} Q_k(s', a') \quad (5)$$

Given V^* , an optimal policy is any policy that for each s assigns non-zero probability only to those actions that realize the maximum on the right-hand side of (4). Similarly, given Q^* , an optimal policy is any policy that for each s assigns non-zero probability only to the actions that maximize $Q^*(s, \cdot)$. These maximizing actions are often called *greedy* actions, and a policy defined in this way is a (stochastic) *greedy policy*. Given sufficiently close approximations of V^* and Q^* obtained by value iteration, then, optimal policies are taken to be the corresponding greedy policies. Note that finding greedy actions via Q^* does not require access to the one-step action models (the $R(s, a)$ and $P(s'|s, a)$) as it does when only V^* is available, where the right-hand side of (4) has to be evaluated. This is one of the reasons that action-value functions play a significant role in RL.

In an MDP, only the sequential nature of the decision process is relevant, not the amount of time that passes between decision stages. A generalization of this is the semi-Markov decision process (SMDP) in which the amount of time between one decision and the next is a random variable, either real- or integer-valued. In the real-valued case, SMDPs model continuous-time discrete-event systems (e.g., Mahadevan et al., 1997; Puterman, 1994). In a *discrete-time SMDP* (Howard, 1971) decisions can be made only at (positive) integer multiples of an underlying time step. In either case, it is usual to treat the system as remaining in each state for a random waiting time (Howard, 1971), at the end of which an instantaneous transition occurs to the next state. Due to its relative simplicity, the discrete-time SMDP formulation underlies most approaches to hierarchical RL, but there are no significant obstacles to extending these approaches to the continuous-time case.

Let the random variable τ denote the (positive) waiting time for state s when action a is executed. The transition probabilities generalize to give the joint probability that a transition from state s to state s' occurs after τ time steps when action a is executed. We write this joint probability as $P(s', \tau | s, a)$. The expected immediate rewards, $R(s, a)$, (which must be bounded) now give the amount of discounted reward expected to accumulate over the waiting time in s given action a . The Bellman equations for V^* and Q^* become

$$V^*(s) = \max_{a \in A_s} \left[R(s, a) + \sum_{s', \tau} \gamma^\tau P(s', \tau | s, a) V^*(s') \right] \quad (6)$$

for all $s \in S$; and

$$Q^*(s, a) = R(s, a) + \sum_{s', \tau} \gamma^\tau P(s', \tau | s, a) \max_{a' \in A_{s'}} Q^*(s', a') \quad (7)$$

for all $s \in S$ and $a \in A_s$. DP algorithms correspondingly extend to SMDPs (e.g., Howard, 1971; Puterman, 1994).

3. Reinforcement Learning

DP algorithms have complexity polynomial in the number of states, but they still require prohibitive amounts of computation for large state sets, such as those that result from discretizing multi-dimensional continuous spaces or representing state sets consisting of all possible configurations of a finite set of structural elements (e.g., possible configurations of a backgammon board Tesauro, 1992). Many methods have been proposed for approximating MDP solutions with less effort than required by conventional DP, but RL methods are novel in their use of Monte Carlo, stochastic approximation, and function approximation techniques. Specifically, RL algorithms combine some, or all, of the following features:

1. Avoid the exhaustive sweeps of DP by restricting computation to states on, or in the neighborhood of, multiple sample trajectories, either real or simulated. Because computation is guided along sample trajectories, this approach can exploit situations in which many states have low probabilities of occurring in actual experience.
2. Simplify the basic DP backup by sampling. Instead of generating and evaluating all of a state's possible immediate successors, estimate a backup's effect by sampling from the appropriate distribution.
3. Represent value functions and/or policies more compactly than lookup-table representations by using function approximation methods, such as linear combinations of basis functions, neural networks, or other methods.

Features 1 and 2 reflect the nature of the approximations usually sought when RL is used. Instead of policies that are close to optimal uniformly over the entire state space, RL methods arrive at non-uniform approximations that reflect the behavior of the agent. The agent's policy does not need high precision in states that are rarely visited. Feature 3 is the least understood aspect of RL, but results exist for the linear case (notably in Tsitsiklis and Van Roy, 1997) and numerous examples illustrate how function approximation schemes that are nonlinear in the adjustable parameters (e.g., multilayer neural networks) can be effective for difficult problems (e.g., Crites and Barto, 1998; Mahadevan et al., 1997; Singh and Bertsekas, 1997; Tesauro, 1992).

Of the many RL algorithms, perhaps the most widely used are Q-learning (Watkins, 1989; Watkins and Dayan, 1992) and Sarsa (Rummery and Niranjan, 1994; Sutton, 1996). Q-learning is based on the DP backup (5) but with the expected immediate reward and the expected maximum action-value of the successor state on the right-hand side of (5) respectively replaced by a sample reward and the maximum action-value for a sample next state. The most common way to obtain these samples is to generate sample trajectories by simulation or by observing the actual decision process over time. Suppose the agent observes a current state s , executes action a , receives immediate reward r , and then observes a next state s' . The Q-learning algorithm updates the current estimate, $Q_k(s, a)$, of $Q^*(s, a)$ using the following update:

$$Q_{k+1}(s, a) = (1 - \alpha_k)Q_k(s, a) + \alpha_k \left[r + \gamma \max_{a' \in A_{s'}} Q_k(s', a') \right] \quad (8)$$

where α_k is a time-varying learning-rate parameter. The values of all the other state-action pairs remain unchanged at this update. If in the limit the action-values of all admissible state-action pairs are updated infinitely often, and α_k decays with increasing k while obeying the usual stochastic approximation conditions, then $\{Q_k\}$ converges to Q^* with probability 1 (Bertsekas and Tsitsiklis, 1996; Jaakkola et al., 1994). As long as these conditions are satisfied, the policy followed by the agent during learning is irrelevant. Of course, when Q-learning is being used, the agent's policy does matter since one is usually interested in the agent's performance throughout the learning process, not just asymptotically. It is usual practice to let the agent select actions using a policy that is greedy with respect to the current estimate of Q^* , while also introducing non-greedy "exploratory actions" in an attempt to widely sample state-action pairs.

Sarsa is similar to Q-learning except that the maximum action-value for the next state on the right-hand side of (8) is replaced by the action-value of the actual next state-action pair:

$$Q_{k+1}(s, a) = (1 - \alpha_k)Q_k(s, a) + \alpha_k [r + \gamma Q_k(s', a')] \quad (9)$$

where a' is the action executed in state s' . (Sutton, 1996, called this algorithm Sarsa due to its dependence on s, a, r, s' , and a' . Equation (9) is actually a special case called Sarsa(0)). Unlike Q-learning, here the agent's policy does matter. Singh et al. (2000) show that if the policy has the property that each action is executed infinitely often in every state that is

visited infinitely often, and it is greedy with respect to the current action-value function in the limit, which Singh et al. (2000) call a GLIE (*Greedy in the limit with infinite exploration*) policy, then with appropriately decaying α_k , the sequence $\{Q_k\}$ generated by Sarsa converges to Q^* with probability 1.

Both the Q-learning and Sarsa learning algorithms also apply to SMDPs, both continuous- and discrete-time, if one interprets the immediate reward, r , as the return accumulated during the waiting time in state s and appropriately adjusts the discounting to reflect the waiting time. For example, in the discrete-time case, if a is executed in state s at time step t and the transition to s' follows after τ time steps, then

$$\begin{aligned} Q_{k+1}(s, a) \\ = (1 - \alpha_k)Q_k(s, a) + \alpha_k \left[r_{t+1} + \gamma r_{t+2} + \cdots + \gamma^{\tau-1} r_{t+\tau} + \gamma^\tau \max_{a' \in A_{s'}} Q_k(s', a') \right] \end{aligned} \quad (10)$$

where r_{t+i} is the immediate reward received at time step $t+i$. The return accumulated during the waiting time must be bounded, and it can be computed recursively during the waiting time. Bradtke and Duff (1995) showed how to do this for continuous-time SMDPs, Parr (1998) proved that it converges under essentially the same conditions required for Q-learning convergence, and Das et al. (1999) developed the average reward case.

Crites (1996) and Crites and Brato (1998) used SMDP Q-learning in a continuous-time discrete-event formulation of an elevator dispatching problem, an application that illustrates two useful features of RL methods for discrete-event systems. First, Q-learning and Sarsa do not require explicit knowledge of the expected immediate rewards or the state-transition probabilities. Instead, they use samples from the respective distributions, which can come from a stochastic simulation or from the real world itself. This can be a significant advantage since in many problems it is often much easier to produce a simulation than it is to make the expected rewards and transition probabilities explicit. In Crites' elevator dispatcher, for example, SMDP Q-learning was applied along trajectories generated by a simulation of a 10-story, 4-elevator system. An explicit SMDP model of this system would have been difficult to make explicit. This illustrates advantages of so-called *model-free* RL algorithms such as Q-learning and Sarsa, meaning that they do not need access to explicit representations of the expected immediate reward function and the state-transition probabilities. Importantly, however, RL is not restricted to such model-free methods (Bertsekas and Tsitsiklis, 1996; Sutton and Barto, 1998).

RL algorithms that estimate action-values, such as Q-learning and Sarsa, have a second advantage when applied to discrete-event systems. As mentioned in Section 2, finding optimal actions via Q^* does not require access to the one-step action models (the $R(s, a)$ and $P(s'|s, a)$) as it does when only V^* is available. That is, a one-step ahead search is not needed to determine optimal actions. In many problems involving discrete-event systems, such as the elevator dispatching problem, it is not clear how to conduct a one-step ahead search since the next event can occur at any of an infinite number of times in the future. The use of action-values eliminates this difficulty.

Finally, we point out that in our brief presentation of RL algorithms we assumed that it

was possible to explicitly store values for every state, or action-values for every state-action pair. This is obviously not feasible for large-scale problems, and extensions of these algorithms need to be considered that adjust the parameters of parametric representations of value functions. It is relatively easy to produce such extensions, although, as we mentioned above, the theory of their behavior still contains many open questions which are beyond the scope of this article.

The view of DP-based RL just outlined, which has been assembled by many researchers over roughly the last ten years, represents our current state of understanding rather than the intuition underlying the origination of these methods. Indeed, DP-based learning originated at least as far back as Samuel's famous checkers player of the 1950s (Samuel, 1963; 1967), which, however, made no reference to the DP literature existing at that time. Other early RL research was explicitly motivated by animal behavior and its neural basis (Klopf, 1972; Klopf, 1982; Minsky, 1954; Sutton and Barto, 1981). Much of the current interest is attributable to Werbos (1977, 1987, 1992), Watkins (1989), and Tesauro's backgammon-playing system TD-Gammon (Tesauro, 1992, 1994). Additional information about RL can be found in several references (e.g., Barto et al., 1995; Bertsekas and Tsitsiklis, 1996; Kaelbling et al., 1996; Sutton and Barto, 1998). Despite the utility of RL methods in many applications, the amount of time they can take to form acceptable approximate solutions can still be unacceptable. As a result, RL researchers are investigating various means for introducing abstraction and hierarchical structure into RL algorithms. In the following sections we review several of the proposed approaches.

4. Approaches to Hierarchical Reinforcement Learning

Artificial intelligence researchers have addressed the need for large-scale planning and problem solving by introducing various forms of abstraction into problem solving and planning systems, e.g., Fikes et al. (1972) and Korf (1985). Abstraction allows a system to ignore details that are irrelevant for the task at hand. One of the simplest types of abstraction is the idea of a "macro-operator," or just a "macro," which is a sequence of operators or actions that can be invoked by name as if it were a primitive operator or action. Macros form the basis of hierarchical specifications of operator or action sequences because macros can include other macros in their definitions: a macro can "call" other macros. Also familiar is the idea of a subroutine that can call other subroutines as well as execute primitive commands. Most of the current research on hierarchical RL focuses on action hierarchies that follow roughly the same semantics as hierarchies of macros or subroutines.

From a control perspective, a macro is an open-loop control policy and, as such, is inappropriate for most interesting control purposes, especially the control of stochastic systems. Hierarchical approaches to RL generalize the macro idea to closed-loop policies, or more precisely, closed-loop *partial* policies because they are generally defined for a subset of the state set. The partial policies must also have well-defined termination conditions. These partial policies are sometimes called *temporally-extended actions*, *options* (Sutton et al., 1999), *skills* (Thrun and Schwartz, 1995), *behaviors* (Brooks, 1986;

Huber and Grupe, 1997), or the more control-theoretic *modes* (Grudic and Ungar, 2000). When not discussing a specific formalism, we will use the term *activity*, as suggested by Harel (1987).

For MDPs, this extension adds to the sets of admissible actions, A_s , $s \in S$, sets of activities, each of which can itself invoke other activities, thus allowing a hierarchical specification of an overall policy. The original one-step actions, now called the “primitive actions,” may or may not remain admissible. Extensions along these general lines result in decision processes modeled as SMDPs, where the waiting time in a state now corresponds to the duration of the selected activity. If τ is the waiting time in state s upon execution of activity a , then a takes τ steps to complete when initiated in s , where the distribution of the random variable τ now depends on the policies and termination conditions of all of the lower-level activities that comprise a . To the best of our knowledge, (and as pointed out by Parr, 1998), the approach most closely related to this in the control literature is that of Forestier and Varaiya (1978), which we discuss briefly in Section 4.2.

4.1. Options

Sutton et al. (1999) formalize this approach to including activities in RL with their notion of an *option*. Starting from a finite MDP, which we call the *core* MDP, the simplest kind of option consists of a (stationary, stochastic) policy $\pi : S \times \cup_{s \in S} A_s \rightarrow [0, 1]$, a termination condition $\beta : S \rightarrow [0, 1]$, and an input set $I \subseteq S$. The option $\langle I, \pi, \beta \rangle$ is available in state s if and only if $s \in I$. If the option is executed, then actions are selected according to π until the option terminates stochastically according to β . For example, if the current state is s , the next action is a with probability $\pi(s, a)$, the environment makes a transition to state s' , where the option either terminates with probability $\beta(s')$ or else continues, determining the next action a' with probability $\pi(s', a')$, and so on. When the option terminates, the agent can select another option.

It is usual to assume that for any state in which an option can continue, it can also be initiated, that is, $\{s : \beta(s) < 1\} \subseteq I$. This implies that an option’s policy only needs to be defined over its input set I . Note that any action of the core MDP, a *primitive action* $a \in \cup_{s \in S} A_s$, is also an option, called a *one-step option*, with $I = \{s : a \in A_s\}$ and $\beta(s) = 1$ for all $s \in S$. Sutton et al. (1999) give the example of an option named `open-the-door` for a hypothetical robot control system. This option consists of a policy for reaching, grasping and turning the door knob, a termination condition for recognizing that the door has been opened, and an input set restricting execution of `open-the-door` to states in which a door is within reach.

An option of the type just defined is called a *Markov option* because its policy is Markov, that is, it sets action probabilities based solely on the current state of the core MDP. To allow more flexibility, especially with respect to hierarchical architectures, one must include *semi-Markov options* whose policies can set action probabilities based on the entire history of states, actions, and rewards since the option was initiated (Sutton et al., 1999). Semi-Markov options include options that terminate after a pre-specified number of time steps, and most importantly, they are needed when policies over options are considered, i.e., policies $\mu : S \times \cup_{s \in S} O_s \rightarrow [0, 1]$, where O_s is the non-empty set of

admissible options for state s (which can include one-step options corresponding to the admissible primitive actions in A_s).

A policy μ over options selects option o in state s with probability $\mu(s, o)$; o 's policy in turn selects other options until o terminates. The policy of each of these selected options selects other options, and so on, until one-step options are selected that correspond to actions of the core MDP. We see, then, that a policy over options, μ , determines a policy of the core MDP, which Sutton et al. (1999) call the *flat* (i.e., non-hierarchical) policy corresponding to μ , denoted $flat(\mu)$. Flat policies corresponding to policies over options are generally not Markov even if all the options are Markov. The probability of a primitive action at any time step depends on the current core state plus the policies of all the options currently involved in the hierarchical specification. This dependence is made more explicit in the work of Parr (1998) and Dietterich (2000), which we discuss below. Using this machinery (made precise by Precup, 2000), one can define *hierarchical options* as triples $\langle I, \mu, \beta \rangle$, where I and β are the same as for Markov options but μ is a semi-Markov policy over options.

Value functions for option policies can be defined in terms of value functions of semi-Markov flat policies. For a semi-Markov flat policy π :

$$V^\pi(s) = E\{r_{t+1} + \gamma r_{t+2} + \cdots + \gamma^{\tau-1} r_{t+\tau} + \cdots | E(\pi, s, t)\}$$

where $E(\pi, s, t)$ is the event of π being initiated at time t in s . Note that this value can depend on the complete history from t onwards, but not on events earlier than t since π is semi-Markov. Given this definition for flat policies, $V^\mu(s)$, the value of s for a policy μ over options, is defined to be $V^{flat(\mu)}(s)$. Similarly, one can define the option-value function for μ as follows:

$$Q^\mu(s, o) = E\{r_{t+1} + \gamma r_{t+2} + \cdots + \gamma^{\tau-1} r_{t+\tau} + \cdots | E(o\mu, s, t)\}$$

where $o\mu$ is the semi-Markov policy that follows o until it terminates after τ time steps and then continues according to μ .

Adding any set of semi-Markov options to a core finite MDP yields a well-defined discrete-time SMDP whose actions are the options and whose rewards are the returns delivered over the course of an option's execution. Since the policy of each option is semi-Markov, the distributions defining the next state (the state at an option's termination), waiting time, and rewards depend only on the option executed and the state in which its execution was initiated. Thus, all of the theory and algorithms applicable to SMDPs can be appropriated for decision making with options.

In their effort to treat options as much as possible as if they were conventional single-step actions, Sutton et al. (1999) introduced the interesting concept of a *multi-time model* of an option that generalizes the single-step model consisting of $R(s, a)$ and $P(s'|s, a)$, $s, s' \in S$, of a conventional action a . For any option o , let $E(o, s, t)$ denote the event of o being initiated in state s at time t . Then the reward part of the multi-time model of o for any $s \in S$ is:

$$R(s, o) = E\{r_{t+1} + \gamma r_{t+2} + \cdots + \gamma^{\tau-1} r_{t+\tau} | E(o, s, t)\}$$

where $t + \tau$ is the random time at which o terminates. The state-prediction part of the model of o for s is:

$$P(s'|s, o) = \sum_{\tau=1}^{\infty} p(s', \tau) \gamma^{\tau},$$

for all $s' \in S$, where $p(s', \tau)$ is the probability that the o terminates in s' after τ steps when initiated in state s . Though not itself a probability, $P(s'|s, o)$ is a combination of the probability that s' is the state in which o terminates together with a measure of how delayed that outcome is in terms of γ .

The quantities $R(s, o)$ and $P(s'|s, o)$ respectively generalize the reward and transition probabilities, $R(s, a)$ and $P(s'|s, a)$, of the usual MDP in such a way that one can write a generalized form of the Bellman optimality equation. If V_o^* denotes the optimal value function over an option set O , then

$$V_o^*(s) = \max_{o \in O_s} \left[R(s, o) + \sum_{s'} P(s'|s, o) V_o^*(s') \right] \quad (11)$$

which reduces to the usual Bellman optimality equation (2) if all the options are one-step options ($\beta(s) = 1, s \in S$). A Bellman equation analogous to (3) can be written for Q_o^* :

$$Q_o^*(s, o) = R(s, o) + \sum_{s'} P(s'|s, o) \max_{o' \in O_{s'}} Q_o^*(s', o'), \quad (12)$$

for all $s \in S$ and $o \in O_s$.

The system of equations (11) and (12) can be solved respectively for V_o^* and Q_o^* , exactly or approximately, using methods that generalize the usual DP and RL algorithms (Precup et al., 1998). For example, the DP backup analogous to (5) for computing option-values is:

$$Q_{k+1}(s, o) = R(s, o) + \sum_{s' \in S} P(s'|s, o) \max_{o' \in O_{s'}} Q_k(s', o')$$

and the corresponding Q-learning update analogous to (8) is:

$$Q_{k+1}(s, o) = (1 - \alpha_k) Q_k(s, o) + \alpha_k \left[r + \gamma^{\tau} \max_{o' \in O_{s'}} Q_k(s', o') \right] \quad (13)$$

This update is applied upon the termination of o at state s' after executing for τ time steps, and r is the return accumulated during o 's execution. This is the specialization of SMDP Q-learning given by (10) to the SMDP that results from adding options to an MDP. It reduces to conventional Q-learning when all the options are one-step options. In addition to Precup (2000) and Sutton et al. (1999), see Parr (1998), Precup and Sutton (1998) and

Precup et al. (1998) for discussion of this formulation and its relation to conventional SMDP theory.

As in the case of conventional MDPs, given V_O^* or Q_O^* , optimal policies over options can be determined as (stochastic) greedy policies. If each set of admissible options, O_s , $s \in S$, contains one-step options corresponding to all of the primitive actions $a \in A_s$ of the core MDP, then it is clear that the optimal policies over options are the same as the optimal policies for the core MDP (since primitive actions give the most refined degree of control). On the other hand, if some of the primitive actions are not available as one-step options, then optimal policies over the set of available options are in general suboptimal policies of the core MDP.

The theory described above all falls within conventional SMDP theory, which is specialized by the focus on semi-Markov options. A shortcoming of this is that the internal structure of an option is not readily exploited. Precup (2000) writes:

SMDP methods apply to options, but only when they are treated as opaque indivisible units. Once an option has been selected, such methods require that its policy be followed until the option terminates. More interesting and potentially more powerful methods are possible by looking inside options and by altering their internal structure.

Precup (2000), p. 58.

For example, the Q-learning update (13) for options does nothing until an option terminates (and so does not apply to non-terminating options), and it only applies to one option at a time. This is the motivation for *intra-option learning methods* which allow learning useful information before an option terminates and can be used for multiple options simultaneously. For example, a variant of Q-learning, called *one-step intra-option Q-learning* (Precup, 2000), works as follows. Suppose (primitive) action a_t is taken in s_t , and the next state and immediate reward are respectively s_{t+1} and r_{t+1} . Then for every Markov option $o = \langle I, \pi, \beta \rangle$ whose policy could have selected a_t according to the same distribution $\pi(s_t, \cdot)$, this update is applied:

$$Q_{k+1}(s_t, o) = (1 - \alpha_k)Q_k(s_t, o) + \alpha_k[r_{t+1} + \gamma U_k(s_t, o)] \quad (14)$$

where

$$U_k(s, o) = (1 - \beta(s))Q_k(s, o) + \beta(s) \max_{o' \in O} Q_k(s, o')$$

which is an estimate of the value of state-option pair (s, o) upon arrival in state s . In the case of deterministic option policies, for example, this update is applied to all options whose policies select a_t in s_t *whether that option was executing or not*. If all the options in O are deterministic and Markov, then for every option in O this converges to Q_O^* with probability 1, provided the α_k decay appropriately and that in the limit every primitive action is executed infinitely often in every state (Precup, 2000).

Although its utility for hierarchical RL is limited due to its restriction to Markov options, one-step intra-option Q-learning is one example of a class of methods that take

advantage of the structure of the core MDP. Related methods have been developed for estimating multi-time models of many options simultaneously by exploiting Bellman-like equations relating the components $R(s, a)$ and $P(s'|s, a)$ of multi-time models for successive states. Results also exist on interrupting an option's execution in favor of an option more highly valued for the current state, and for adjusting an option's termination condition to allow the longest expected execution without sacrificing performance. See Sutton et al. (1999) and Precup (2000) for details on these and other option-related algorithms and illustrations of their performance.

The primary motivation for the options framework is to permit one to add temporally-extended activities to the repertoire of choices available to an RL agent, while at the same time not precluding planning and learning at the finer grain of the core MDP. The emphasis is therefore on augmentation rather than simplification of the core MDP. If all the primitive actions remain in the option set as one-step options, then clearly the space of realizable policies is unrestricted so that the optimal policies over options are the same as the optimal policies for the core MDP. But since finding optimal policies in this case takes more computation via conventional DP than does just solving the core MDP, one is tempted to ask what one gains from this augmentation of the core MDP. One answer is to be found in the use of RL methods. For RL, the availability of temporally-extended activities can dramatically improve the agent's performance while it is learning, especially in the initial stages of learning. Invoking multi-step options provides one way to prevent the prolonged period of "flailing" that one often sees in RL systems. Options also can facilitate transfer of learning to related tasks. Of course, only some options can facilitate learning in this way, and a key question is how does a system designer decide on what options to provide. On the other hand, if the set of options does not include the one-step options corresponding to all of the primitive actions, then the space of policies over options is a proper subset of the set of all policies of the core MDP. In this case, though, the resulting SMDP can be much easier to solve than the core MDP: the options simplify rather than augment. This is the primary motivation of two other approaches to abstraction in RL that we consider in the next two sections.

In the current state-of-the-art, the designer of an RL system typically uses prior knowledge about the task to add a specific set of options to the set of primitive actions available to the agent. In some cases, complete option policies can be provided; in other cases, option policies can be learned using, for example, intra-option learning methods together with option-specific reward functions that are provided by the designer. Providing options and their policies a priori is an opportunity to use background knowledge about the task to try to accelerate learning and/or provide guarantees about system performance during learning. Perkins and Barto (to appear); (2001), for example, consider collections of options each of which descends on a Lyapunov function. Not only is learning accelerated, but the goal state is reached on every learning trial while the agent learns to reach the goal more quickly by approximating a minimum-time policy over these options.

When option policies are learned, they usually are policies for efficiently achieving *subgoals*, where a subgoal is often a state, or a region of the state space, such that reaching that state or region is assumed to facilitate achieving the overall goal of the task. The canonical example of a useful subgoal is a doorway in a robot navigation scenario: the doorway has to be passed through to reach any goal outside the room. Given a collection of

subgoals, one can define subgoal-specific reward functions that positively reward the agent for achieving the subgoal (while possibly penalizing it until the subgoal is achieved). Options are then defined which terminate upon achieving a subgoal, and their policies can be learned using the subgoal-specific reward function and standard RL methods. Precup (2000) discusses one way to do this by introducing *subgoal values*, and Dietterich (2000), whose approach we discuss in Section 4.3, proposes a similar scheme using *pseudo-reward* functions.

A natural question, then, is how are useful subgoals determined? McGovern (2002) and McGovern and Brato (2001) developed a method for automatically identifying potentially useful subgoals by detecting regions that the agent visits frequently on successful trajectories but not on unsuccessful trajectories. An agent using this method selects such regions that appear early in learning and persist throughout learning, creates options for achieving them and learns their policies, and at the same time learns a higher-level policy that invokes these options appropriately to solve the overall task. Experiments with this method suggest that it can be useful for accelerating learning on single tasks, and that it can facilitate knowledge transfer as previously-discovered options are reused in related tasks. This approach builds on previous work in artificial intelligence that addresses abstraction, particularly that of Iba (1989), who proposed a method for discovering macro-operators in problem solving. Related ideas have been studied by Digney (1996, 1998).

4.2 Hierarchies of Abstract Machines

Parr (1998) and Parr and Russell (1988) developed an approach to hierarchically structuring MDP policies called *Hierarchies of Abstract Machines* or HAMs. Like the options formalism, HAMs exploit the theory of SMDPs, but the emphasis is on simplifying complex MDPs by restricting the class of realizable policies rather than expanding the action choices. In this respect, as pointed out by Parr (1998), it has much in common with the multilayer approach for controlling large Markov chains described by Forestier and Varaiya (1978) who considered a two-layer structure in which the lower level controls the plant via one of a set of pre-defined regulators. The higher level, the supervisor, monitors the behavior of the plant and intervenes when its state enters a set of boundary states. Intervention takes the form of switching to a new low-level regulator. This is not unlike many hybrid control methods (Bramichy et al., 1998) except that the low-level process is formalized as a finite MDP and the supervisor's task as a finite SMDP. The supervisor's decisions occur whenever the plant reaches a boundary state, which effectively "erases" the intervening states from the supervisor's decision problem, thereby reducing its complexity (Forestier and Varaiya, 1978). In the options framework, each option corresponds to a low-level regulator, and when the option set does not contain the one-step options corresponding to all primitive actions, the same simplification results. HAMs extend this idea by allowing policies to be specified as hierarchies of stochastic finite-state machines.

The idea of the HAM approach is that policies of a core MDP are defined as programs which execute based on their own states as well as the current states of the core MDP. Departing somewhat from Parr's (1998) notation, let M be a finite MDP with state set S and

action sets A_s , $s \in S$. A HAM policy is defined by a collection of stochastic finite-state machines, $\{H_i\}$, with state sets S_i , stochastic transition functions δ_i , and input sets all equal to M 's state set, S . Each machine i also has a stochastic function $I_i : S \rightarrow S_i$ that sets the initial state of M in the manner described below. Each H_i has four types of states: *action*, *call*, *choice*, and *stop*. An action state generates an action of the core MDP, M , based on the current state of M and the current state of the currently executing machine, say H_i . That is, at time step t , the action $a_t = \pi(m_t^i, s_t) \in A_{s_t}$, where m_t^i is the current state of H_i and s_t is the current state of M . A call state suspends execution of the currently executing H_i and initiates execution of another machine, say H_j , where j is a function of H_i 's state m_t^i . Upon being called, the state of H_j is set to $I_j(s_t)$. A choice state nondeterministically selects a next state of H_i . Finally, a stop state terminates execution of H_i and returns control to the machine that called it (whose execution commences where it was suspended). Meanwhile, the core MDP, upon receiving an action, makes a transition to a next state according to its transition probabilities and generates an immediate reward.² If no action is generated at step t , then M remains in its current state. Parr defines a HAM H to be the initial machine together with the closure of all machine states in all machines reachable from the possible initial states of the initial machine. Let us call this state set S_H . For convenience, he also assumes that the initial machine does not have a stop state and that there are no infinite, probability 1, loops that do not contain action states. This ensures that the core MDP continues to receive primitive actions.

Figure 1 shows a simple HAM state-transition diagram similar to an example given by Parr and Russell (1998) for control simple simulated mobile robot. This HAM runs until the robot reaches an intersection. Whenever an obstacle is encountered, a choice state is entered that allows the robot to decide to back away from the obstacle by calling the machine *back-off* or to try to get around the obstacle by calling the machine *follow-wall*. Each of these machines has its own state-transition structure, possibly containing additional choice and call states. When this HAM is selected, it deterministically starts by calling the *follow-wall* machine.

The composition of a HAM H and an MDP M , as described above, yields a discrete-time SMDP denoted $H \circ M$. The state set of $H \circ M$ is $S \times S_H$, and its transitions are determined by the parallel action of the transition functions of H and M . The only actions of $H \circ M$ are the choices allowed in the *choice points* of $H \circ M$, which are the states whose H components are choice states. These actions change only the HAM component of each state. This is an SMDP because after a choice is made, the system—the composition of H and M —runs autonomously until another choice point is reached. All the primitive actions to M during this period are fully determined by the action states of H . The expected immediate rewards of $H \circ M$ are the expected returns accumulated during these periods between choice points, and they are determined by the immediate rewards of M together with rewards of zero for the time steps in which M 's state does not change. Thus, one can think of a HAM as a method for delineating a possibly drastically restricted set of policies for M . This restriction is determined by the prior knowledge that the HAM's designer, or programmer, has about what might be good ways to control M .

The next step, which corresponds to the main observation of Forestier and Varaiya (1978), is to note that in determining an optimal policy for $H \circ M$, the only relevant states are the choice points; the rest can be “erased.” Therefore, there is an SMDP, called

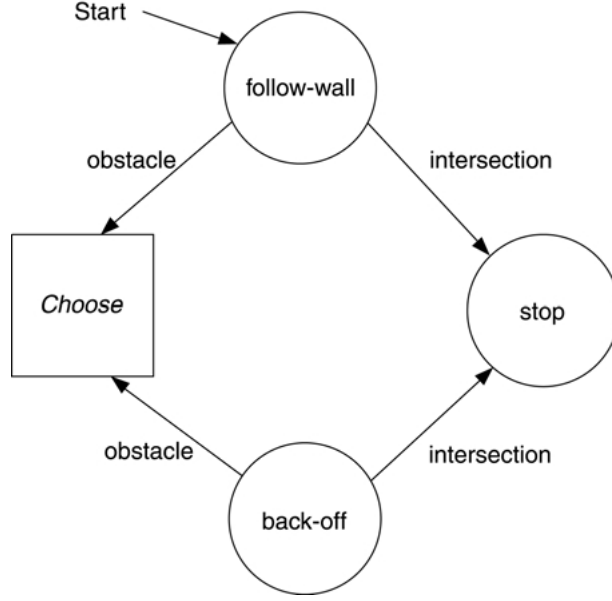


Figure 1. State-transition structure of a simple HAM (after Parr, 1998).

$reduce(H \circ M)$, that is equivalent to $H \circ M$ but whose states are just the choice points of $H \circ M$. Optimal policies of $reduce(H \circ M)$ are the same as the optimal policies of $H \circ M$. Of course, how close these policies will be to optimal policies of M will depend on the programmer's knowledge and skill.

How does RL enter into the HAM framework? It is easy to see that SMDP Q-learning given by (10) can be applied to $reduce(H \circ M)$ to approximate optimal policies for $H \circ M$. The important strength of an RL method like SMDP Q-learning in this context is that it can be applied to $reduce(H \circ M)$ without performing any explicit reduction of $H \circ M$. Trajectories of M under control of HAM H , i.e., trajectories of $H \circ M$, are generated, either through simulation or observed from the real system. The update (10) is applied from choice point to choice point. In more detail, the learning system maintains a Q-table with entries $Q([s, m], a)$ for each state, s , of M , each choice state, m , of H , and each action, a , that can be taken from the corresponding choice point. It also has to store the previous choice point, $[s_c, m_c]$, and the action selected at that choice point, a_c . Suppose choice point $[s_c, m_c]$ is encountered at time step t , and the next choice point, $[s'_c, m'_c]$, is encountered at $t + \tau$. Then the SMDP Q-learning update (10) appears as follows:

$$Q_{k+1}([s_c, m_c], a_c) = (1 - \alpha_k)Q_k([s_c, m_c], a_c) + \alpha_k[r_{t+1} + \gamma r_{t+2} + \dots + \gamma^{\tau-1}r_{t+\tau} + \gamma^\tau \max_{a'} Q_k([s'_c, m'_c], a')]$$

Here, the max is taken over all actions available in choice point $[s'_c, m'_c]$. Clearly the between-choice return, $r_{t+1} + \gamma r_{t+2} + \dots + \gamma^{\tau-1}r_{t+\tau}$, can be accumulated iteratively

between the visits to these choice points. If the SMDP Q-learning convergence conditions are in force (Section 3), then the sequence of action-values functions generated by this algorithm converges to Q^* of $\text{reduce}(H \circ M)$ with probability 1. Consequently, any sequence of policies that are greedy with respect to these successive action-value functions converges with probability 1 to an optimal policy.

We know of no large-scale applications of HAMs at the present time, but Parr (1998) and Parr and Russell (1998) illustrate the advantages that HAMs can provide in several simulated robot navigation tasks. Andre and Russell (2001) increased the expressive power of HAMs by introducing Programmable HAMs (PHAMs), which add interrupts, aborts, local state variables, and the ability to pass parameters. PHAMs, and other extensions of HAMs that may occur in the future, point the way toward methods theoretically grounded in stochastic optimal control that use expressive programming languages to provide a knowledge-rich context for RL.

4.3 MAXQ Value Function Decomposition

Dietterich (2000) developed another approach to hierarchical RL called the MAXQ Value Function Decomposition, which we call simply MAXQ. Like options and HAMs, this approach also relies on the theory of SMDPs. Unlike options and HAMs, however, the MAXQ approach does not rely directly on reducing the entire problem to a single SMDP. Instead, a hierarchy of SMDPs is created whose solutions can be learned simultaneously. The MAXQ approach starts with a decomposition of a core MDP M into a set of subtasks $\{M_0, M_1, \dots, M_n\}$. The subtasks form a hierarchy with M_0 being the root subtask, which means that solving M_0 solves M . Actions taken in solving M_0 consist of either executing primitive actions or policies that solve other subtasks, which can in turn invoke primitive actions or policies of other subtasks, etc.

The structure of the hierarchy is summarized in a *task graph*, an example of which is given in Figure 2 for a Taxi problem that Dietterich used as an illustration. Each episode of the overall task consists of picking up, transporting, and dropping off a passenger. The overall problem, corresponding to the root node of the graph, is decomposed into the subtask *Get*, which is the subtask of going to the passenger's location and picking them up, and the subtask *Put*, which is the subtask of going to the passenger's destination and dropping them off. These subtasks, in turn, are respectively decomposed into the primitive actions *Pickup* or *Dropoff*, which respectively pick up and drop off a passenger, and the subtask *Navigate(t)*, which consists of navigating to one of the locations indicated by the parameter t . (A subtask parameterized like this is shorthand for multiple copies of the subtask, one for each value of the parameter.) This subtask *Navigate(t)* is decomposed into the primitive actions that are moves *North*, *South*, *East*, or *West*. The subtasks and primitive actions into which a subtask M_i is decomposed are called the “children” of M_i . An important aspect of a task graph is that the order in which a subtask's children are shown is arbitrary. Which choice the higher level controller makes depends on its policy. The graph just restricts the action choices that can be made at each level. See Dietterich (2000) for details.

Each subtask, M_i , consists of three components. First, it has a subtask policy, π_i , that

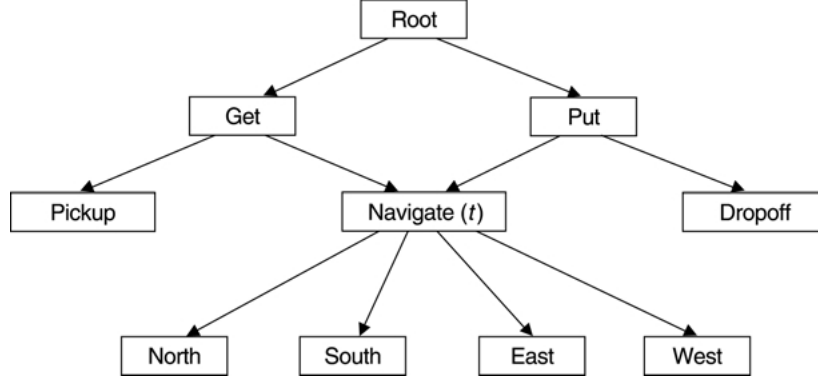


Figure 2. A task graph for the taxi problem (after Dietterich, 2000).

can select other subtasks from the set of M_i 's children. Here, as with options, primitive actions are special cases of subtasks. We also assume the subtask policies are deterministic. Second, each subtask has a termination predicate that partitions the state set, S , of the core MDP into S_i , the set of *active states*, in which M_i 's policy can execute, and T_i , the set of *termination states*, which when entered causes the policy to terminate. Third, each subtask M_i has a *pseudo-reward function* that assigns reward values to the states in T_i . The pseudo-reward function is only used during learning, which we discuss after first describing how the task graph hierarchy allows value functions to be decomposed.

A subtask is very much like a hierarchical option, $\langle I_i, \mu_i, \beta_i \rangle$, as defined in Section 4.1, with the addition of a pseudo-reward function. The policy over options, μ_i , corresponds to the subtask's π_i ; the termination condition, β_i , in this case assigns to states termination probabilities of only 1 or 0; and the option's input set I_i corresponds to S_i . Unlike the option formalism, however, which treats semi-Markov options, MAXQ explicitly adds a component to each state that gives the current contents, K , of a pushdown stack containing the names and parameter values of the hierarchy of calling subtasks, as in subroutine handling of ordinary programming languages. At any time step, the top of the stack contains the name of the subtask currently being executed. Thus, while a subtask's policy is non-Markov with respect to the state set of the core MDP, it is Markov with respect to this augmented state set. As a consequence, each subtask policy has to assign actions to every combination, $[s, K]$, of core state s and stack contents K .

Given a hierarchical decomposition of M into n subtasks as given by a task graph, a *hierarchical policy* is defined to be $\pi = \{\pi_0, \dots, \pi_n\}$, where π_i is the policy of M_i . The *hierarchical value function* for π gives the value, i.e., the expected return, for each state-stack pair, $[s, K]$, given that π is followed from s when the stack contents are K . This value is denoted $V^\pi([s, K])$. The "top level" value of a state s is $V^\pi([s, nil])$, indicating that the stack is empty. This is the value of s under the flat policy induced by the hierarchy of subtask calls starting from the root of the hierarchy. The *projected value function* of hierarchical policy π on subtask M_i gives the expected return of each state s under the assumption that π_i is executed until it terminates. This projected value is denoted $V^\pi(i, s)$.

This is the same as the value function of the hierarchical option corresponding to M_i as defined in Section 4.1.

Given a hierarchical decomposition of M and a hierarchical policy, π , each subtask M_i defines a discrete-time SMDP with state set is S_i . Its actions are its child subtasks, M_a , and its transition probabilities, $P_i(s', \tau | s, a)$ (cf. equations (6) and (7)), are well-defined given the policies of the lower-level subtasks. A key observation, which follows that of Singh (1992, 1992), is that this SMDP's expected immediate reward, $R_i(s, a)$, for executing action (subtask) a is the projected value of π on subtask M_a . That is, for all $s \in S_i$ and all child subtasks M_a of M_i , $R_i(s, a) = V^\pi(a, s)$. To see why this is true, suppose the core MDP is in state s when subtask M_i selects one of its child subtasks, M_a , for execution. Expanding M_a 's policy down to primitive actions results in a flat policy that accumulates rewards according to the core MDP until M_a terminates at a state $s' \in T_a$. The waiting time for s when this action is chosen is the execution time of M_a 's policy. The reward accumulated during this waiting time is $V^\pi(a, s)$. Given this, one can write a Bellman equation for the SMDP corresponding to subtask M_i :

$$V^\pi(i, s) = V^\pi(\pi_i(s), s) + \sum_{s', \tau} P_i^\pi(s', \tau | s, \pi_i(s)) \gamma^\tau V^\pi(i, s') \quad (15)$$

where $V^\pi(i, s')$ is the expected return for completing subtask M_i starting in state s' (cf. equation (6)).

The action-value function, Q , defined by (1) can be extended to apply to subtasks: for hierarchical policy π , $Q^\pi(i, s, a)$ is the expected return for action a (a primitive action or a child subtask) being executed in subtask M_i and then π being followed until M_i terminates. In terms of this subtask action-value function, the observation expressed in (15) takes the form:

$$Q^\pi(i, s, a) = V^\pi(a, s) + \sum_{s', \tau} P_i^\pi(s', \tau | s, a) \gamma^\tau Q^\pi(i, s', \pi(s')) \quad (16)$$

Dietterich (2000) calls the second term on the right of this equation the *completion function*:

$$C^\pi(i, s, a) = \sum_{s', \tau} P_i^\pi(s', \tau | s, a) \gamma^\tau Q^\pi(i, s', \pi(s'))$$

This gives expected return for completing subtask M_i after subtask M_a terminates. Rewriting (16) using the completion function, we have

$$Q^\pi(i, s, a) = V^\pi(a, s) + C^\pi(i, s, a) \quad (17)$$

Now it is possible to describe the MAXQ hierarchical value function decomposition. Given a hierarchical policy π and a state s of the core MDP, suppose the policy of the top-level subtask, M_0 , selects subtask M_{a_1} , and that this subtask's policy selects subtask M_{a_2} ,

whose policy in turn selects subtask M_{a_3} , etc. until finally subtask $M_{a_{n-1}}$'s policy selects a_n , a primitive action that is executed in the core MDP. Then the projected value of s for the root subtask, $V^\pi(0, s)$, which is the value of s in the core MDP, can be written as follows:

$$V^\pi(0, s) = V^\pi(a_n, s) + C^\pi(a_{n-1}, s, a_n) + \cdots + C^\pi(a_1, s, a_2) + C^\pi(0, s, a_1) \quad (18)$$

where $V^\pi(a_n, s) = \sum_{s'} P(s'|s, a_n) R(s'|s, a_n)$.

The decomposition (18) is the basis of a learning algorithm that is able to learn hierarchical policies from sample trajectories. The details of this algorithm are somewhat complex and beyond the scope of this review, but we describe the basic idea. This algorithm is a recursively applied form of SMDP Q-learning that takes advantage of the MAXQ value function decomposition to update estimates of subtask completion functions. The simplest version of this algorithm applies when the pseudo-reward functions of all the subtasks are identically zero. The key component of the algorithm is a function—called MAXQ-0 for this special case—that calls itself recursively to descend through the hierarchy to finally execute primitive actions. When it returns from each call, it updates the completion function corresponding to the appropriate subtask, with discounting determined by the returned count of the number of primitive actions executed. For details see Dietterich (2000).

If the agent follows a policy that is GLIE (Section 3) and also further constrained to always break ties in the same order, and the step-size parameter converges to zero according to the usual stochastic approximation conditions, then the algorithm sketched above converges with probability 1 to the unique *recursively optimal* policy for M that is consistent with the task graph. A recursively optimal policy is a hierarchical policy, $\pi = \{\pi_0, \dots, \pi_n\}$, such that for each subtask M_i , $i = 0, \dots, n$, π_i is optimal for the SMDP corresponding to M_i , given the policies of M_i 's children subtasks. This form of optimal policy stands in contrast to a *hierarchically optimal* policy, which is a hierarchical policy that is optimal among all the policies that can be expressed within the constraints imposed by the given hierarchical structure.

Examples of policies that are recursively optimal but not hierarchically optimal are easy to construct. This is because the hierarchical optimality of a subtask generally depends not only on that subtask's children, but also on how the subtask participates in higher-level subtasks. For example, the hierarchically optimal way to travel to a given destination may depend on what you intend to do after arriving there in addition to properties of the trip itself. SMDP DP or RL methods applied to HAMs, or to MDPs with options whose policies are fixed a priori, yield hierarchically optimal policies, and it is relatively easy to define a learning algorithm for MAXQ hierarchies that also produces hierarchically optimal policies. Dietterich's interest in the weaker recursive optimality stems from the fact that this form of optimality can be determined without considering the context of a subtask. This facilitates using subtask policies as building blocks for other tasks and has implications for state abstraction, although we do not touch on the latter aspect in this review.

What about the pseudo-reward functions? These functions allow a system designer within the MAXQ framework to specify subtasks by defining subgoals that they must

achieve but without specifying policies for achieving them. In this respect, they play the same role as do the auxiliary reward functions that can be used in the options framework for learning option policies. The MAXQ-0 learning algorithm sketched above can be extended to an algorithm, called MAXQ-Q, that learns a policy that is recursively optimal hierarchical with respect to the sum of the original reward function and the pseudo-reward functions.

We have not been able to do justice to the MAXQ approach and associated algorithms in this short review. However, we have presented enough to provide a basis for seeing how SMDP theory again plays a central role in an approach to hierarchical RL. Like HAMs, MAXQ is an excellent example of how concepts from programming languages can be fruitfully integrated with a stochastic optimal control framework.

5. Recent Advances in Hierarchical RL

Thus far we have surveyed different approaches to hierarchical RL, all of which are based on the underlying framework of SMDPs. These approaches collectively suffer from some key limitations: policies are restricted to *sequential* combinations of activities; agents are assumed to act alone in the environment; and finally, states are considered to be fully accessible. In this section, we address these assumptions and describe how the SMDP framework can be extended to *concurrent activities*, *multiagent domains*, and *partially observable* states.

5.1 Concurrent Activities

Here we summarize recent work by Rohanimanesh and Mahadevan, 2001 towards a general framework for modeling concurrent activities. This framework is motivated by situations in which a single agent can execute multiple parallel processes, as well as by the multiagent case in which many agents act in parallel (addressed in Section 5.2). Managing concurrent activities clearly is a central component of our everyday behavior: in making breakfast, we interleave making toast and coffee with other activities such as getting milk; in driving, we search for road signs while controlling the wheel, accelerator and brakes.

Building on the SMDP framework described above, this approach focuses on modeling concurrent activities, where each component activity is temporally extended. One immediate question that arises in modeling concurrent activities is that of termination: unlike the purely sequential case, when multiple activities are executed simultaneously, concurrently executing activities do not generally terminate at the same time. How does one define the termination of a set of concurrent activities? Among the termination schemes studied by Rohanimanesh and Mahadevan (2001) are the following two. The *any* scheme terminates all other activities when the first activity terminates, whereas the *all* scheme waits until all existing activities terminate before choosing a new concurrent set of activities. These researchers have also studied other schemes, such as *continue*, which replaces the terminated activity with a set of new activities that includes

the activities already executing. For simplicity, we restrict our discussion to the first two schemes.

For concreteness, we describe the concurrent activity model using the options formalism described in Section 4.1. The treatment here is restricted to options over discrete-time SMDPs and having deterministic policies, but the main ideas extend readily to the other variants (HAMs, MAXQ), as well as to continuous-time SMDPs. (See Ghavamzadeh and Mahadevan, 2000 for a treatment of hierarchical RL for continuous-time SMDPs.) The sequential option model is generalised to a *multi-option*, which is a set of options that can be executed in parallel. Here we discuss the simple case in which the options comprising a multi-option influence different sets of state variables that do not interact. (This assumption can be generalized, but for simplicity, we restrict our treatment here.) For example, turning the radio off and pressing the brake can always be executed in parallel since they affect different non-interacting state variables.

We now define a multi-option (denoted by \vec{o}) more precisely. Let $o = \langle I, \pi, \beta \rangle$ be an option as defined in Section 4.1 on a core MDP with state set $S \subseteq \prod_1^n S_i$, where S_i is the range of state variable s_i , $i = 1, \dots, n$. Suppose that while option o is executing, it influences only a subset of the state variables $S_o \subseteq \{s_1, s_2, \dots, s_n\}$. Assume also that while o is executing, these variables do not influence, and are not influenced by, any variables not in S_o . A set of options $\{o_1, \dots, o_m\}$, each defined on the same core MDP, is *coherent* if $\cap_1^m S_{o_i} = \emptyset$. This ensures that each option will affect different components of the state set so that any subset of them can be executed in parallel without interfering with one another. In other words, while a coherent set of options is executing, the system admits a parallel decomposition with each component corresponding to the state variables influenced by one of the options. In the driving example, the `turn-radio-on` and `brake` options comprise a coherent set, but the `turn-right` and `accelerate` options do not since the state variable `position` is influenced by both. A multi-option over an MDP M is a coherent set of options over M .

When multi-option \vec{o} is executed in state s , several options $o_i \in \vec{o}$ are initiated. Each option o_i will terminate at some random time t_{o_i} . We can define the termination of a multi-option based on either of the following events: (1) $T_{\text{all}} = \max_i(t_{o_i})$: when all the options $o_i \in \vec{o}$ terminate according to $\beta_i(s)$, multi-option \vec{o} is declared terminated, or (2) $T_{\text{any}} = \min_i(t_{o_i})$: when any (i.e., the first) of the options terminate, the options that are still executing at that time are interrupted.

We then have the following result. Given a finite MDP and a collection of multi-options defined on it, where the underlying options are Markov, the decision process that selects only among multi-options, and executes each one until its termination according to the T_{all} or T_{any} termination conditions, is a discrete-time SMDP. The proof requires showing that the state transition probabilities and the rewards corresponding to any concurrent option \vec{o} defines an SMDP (Rohanimanesh and Mahadevan, 2001). The significance of this result is that SMDP Q-learning methods can be extended to learn policies over concurrent options under this model.

The extended SMDP Q-learning algorithm for learning policies over multi-options updates the multi-option value function $Q(s, \vec{o})$ after each decision epoch in which the multi-option \vec{o} is taken in state s and terminates in state s' (under either termination condition):

$$Q_{k+1}(s, \vec{o}) = (1 - \alpha_k)Q_k(s, \vec{o}) + \alpha \left[r + \gamma^\tau \max_{\vec{o}' \in \mathcal{O}_{s'}} Q_k(s', \vec{o}') \right] \quad (19)$$

where τ denotes the number of time steps between initiation of the multi-option \vec{o} in state s and its termination in state s' , and r denotes the cumulative discounted reward over this period. This learning rule generalizes (13) to multi-options. It is straightforward to also generalize other learning algorithms, such as the intra-option Q-learning algorithm (14) to multi-options.

As a simple illustrative example of using this concurrent SMDP Q-learning algorithm, Rohanimanesh and Mahadevan (2001) considered an environment consisting of four “rooms”, each of which is divided into a discrete set of cells (Figure 3). Each room contains two doors, each of which can be “opened” by two keys. The agent is given options for getting to each door from any interior room state, and for opening a locked door. It has to learn the shortest path to the goal by concurrently combining these options. The agent can reach the goal more quickly if it learns to parallelize the option for retrieving the key before it reaches a locked door. However, retrieving the key too early is counterproductive since it can be dropped with some probability. The process of retrieving a key is modeled as an SMDP consisting of a linear chain of states to model the waiting process before the agent is holding the key.

Figure 4 compares the policies learned with strictly sequential options, as described in Section 4.1, using the Q-learning algorithm (13), with policies over multi-options learned using the extended Q-learning rule (19). The vertical axis plots the median steps to the goal in each trial. Note that the sequential solution is the slowest, taking the agent about 50 steps to reach the goal. The termination condition T_{all} is faster than the sequential case,

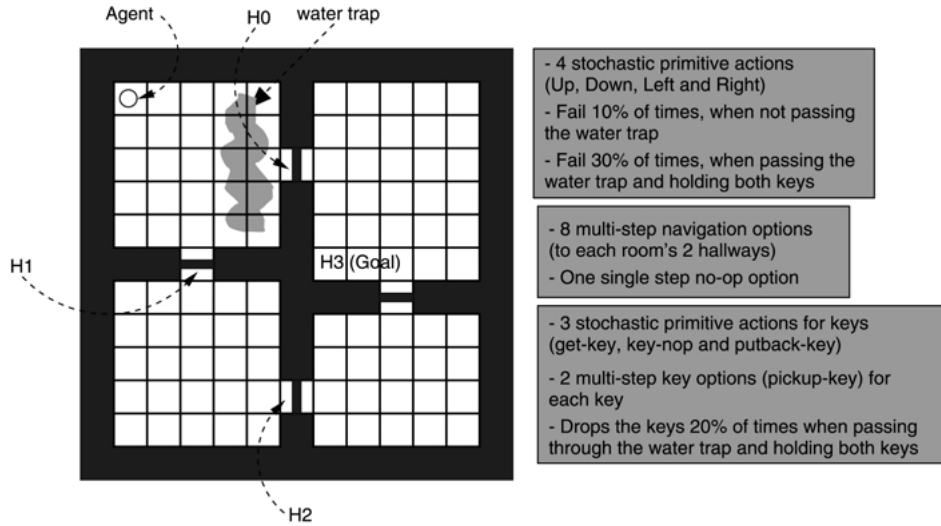


Figure 3. An Illustrative Problem for Concurrent Options.

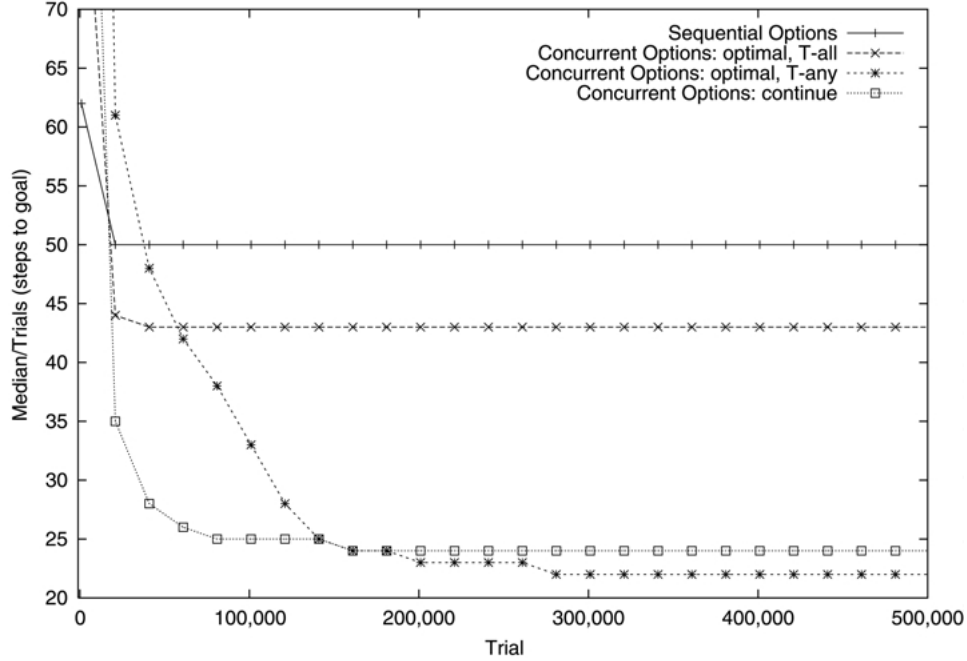


Figure 4. Multi-Option Comparison. This graph compares a multi-option SMDP Q-learning system (under different termination schemes) with one that learns sequential policies. Policies over multi-options easily outperform sequential policies, and termination makes a large difference in of convergence rate and quality of the learned policy.

taking about 45 steps to reach the goal. Both of these, however, are significantly slower than the T_{any} and $T_{continue}$ policies, which reach the goal in half the time. T_{any} converges the slowest, since it terminates multi-options fairly aggressively resulting in more decision epochs. $T_{continue}$ provides the best tradeoff between speed of convergence and optimality of the learned policy.

The multi-option formalism can be extended to allow cases in which options executing in parallel modify the same shared variables at the same time. Details of this extension is beyond the scope of the present article.

5.2. Multiagent Coordination

Concurrency is the basis for modeling coordination among multiple simultaneously behaving agents. From a theoretical point of view, it matters little if concurrent activities are being executed by a single agent or by multiple cooperating agents. However, the multiagent problem usually involves other complexities as well, such as the fact that one agent cannot usually observe the actions or the states of other agents in the environment. We address the issue of hidden state in the next section; here we focus on the problem of learning a policy across *joint states* and *joint actions*.

In general, problems involving multiagent coordination can be modeled by assuming that states now represent the *joint state* of n agents, where each agent i may only have access to a partial view s_i . Also, the joint action is represented as (a_1, \dots, a_n) , where again, each agent may not have knowledge of the other agents' actions. Typically, the assumption made in multiagent studies is that the set of joint actions defines an SMDP (or MDP) over the joint state set. Of course, in most practical problems, the joint state and action sets are exponential in the number of agents, and the aim is to find a *distributed* solution that does not require combinatorial enumeration over joint states and actions.

One general approach to learning task-level coordination is to extend the above concurrency model to the joint state and action spaces, where each action is a fixed option with a pre-specified policy. This approach requires very minimal modification of the approach described in the previous section. In contrast, we now describe an extension of this approach due to Makar and Mahadevan (2001) in which agents learn both coordination skills and the base-level policies using a multiagent MAXQ-like task graph. However, convergence to (hierarchically) optimal policies is no longer assured since lower-level subtask policies are varying at the same time when learning higher-level policies. The ideas can be extended to other formalisms also, but for the sake of clarity, we focus on the MAXQ value function decomposition approach described in Section 4.3.

It is necessary to generalize the MAXQ decomposition from its original sequential single-agent setting to the multiagent coordination problem. Let $\vec{o} = (o_1, \dots, o_n)$ denote a multi-option, where o_i is the option executed by agent i .³ Let $s = (s_1, \dots, s_n)$ denote a joint state. The joint action-value of a multi-option \vec{o} in a joint state s , and in the context of doing parent task p , is denoted $Q(p, s, \vec{o})$. The MAXQ decomposition of the Q-function given by (17) can be extended to joint action-values as follows. The *joint completion function* for agent j assigns values $C^j(p, s_j, \vec{o})$ giving the discounted return for agent j completing a multi-option in the context of doing parent task p , when the other agents are performing multi-options o_k , for all $k \in \{1, \dots, n\}, k \neq j$. The joint multi-option value $Q(p, \vec{s}, \vec{o})$ is now approximated by each agent j (given only its local state s_j) as:

$$Q^j(p, s_j, \vec{o}) \approx V^j(\vec{o}_j, s_j) + C^j(p, s_j, \vec{o})$$

where

$$V^j(p, s_j) = \begin{cases} \max_{\vec{o}_k} Q^j(p, s_j, \vec{o}_k) & \text{if parent task } p \text{ is non- primitive} \\ \sum_{s'_j} \bar{P}(s'_j | s_j, p) R(s'_j | s_j, p) & \text{if } p \text{ is a primitive action} \end{cases}$$

The first term in the $Q(p, s_j, \vec{o})$ expansion above refers to the discounted sum of rewards received by agent j for doing concurrent action \vec{o}_j in state s_j . The second term ‘‘completes’’ the sum by accounting for rewards earned for completing the parent task p after finishing \vec{o}_j . The completion function is updated from sample values using an SMDP learning rule. Note that the correct action value is approximated by only considering local state s_j and also by ignoring the effect of concurrent actions $\vec{o}_k, k \neq j$ by other agents when agent j is performing \vec{o}_j . In practice, a human designer can configure the task graph to store joint concurrent action values at the highest level(s) of the hierarchy as needed.

Figure 5 illustrates a robot trash collection task, where the two agents, A1 and A2, will maximize their performance at the task if they learn to coordinate with each other. Here, we want to design learning algorithms for *cooperative* multiagent tasks (Weiss, 1999), where the agents learn the coordination skills by trial and error. The key idea here is that coordination skills are learned more efficiently if agents learn to synchronize using a hierarchical representation of the task structure (Sugawara and Lesser, 1998). In particular, rather than each robot learning its response to low-level primitive actions of the other robots (for instance, if A1 goes forward, what should A2 do), they learn high-level coordination knowledge (what is the utility of A2 picking up trash from T1 if A1 is also picking up from the same bin, and so on). The proposed approach differs significantly from previous work in multiagent reinforcement learning (Littman, 1994; Tan, 1993) in using hierarchical task structure to accelerate learning, and as well in its use of concurrent activities.

To illustrate the use of this decomposition in learning multiagent coordination, for the two-robot trash collection task, if the joint action-values are restricted to only the highest level of the task graph under the root, we get the following value function decomposition for agent A1:

$$Q^1(\text{Root}, s_1, (\text{NavT1}, \text{NavT2})) \approx V_t^1((\text{NavT1}), s_1) + C_t^1(\text{Root}, s_1, (\text{NavT1}, \text{NavT2}))$$

which represents the value of agent A1 doing task *NavT1* in the context of the overall *Root* task, when agent A2 is doing task *NavT2*. Note that this value is decomposed into the value

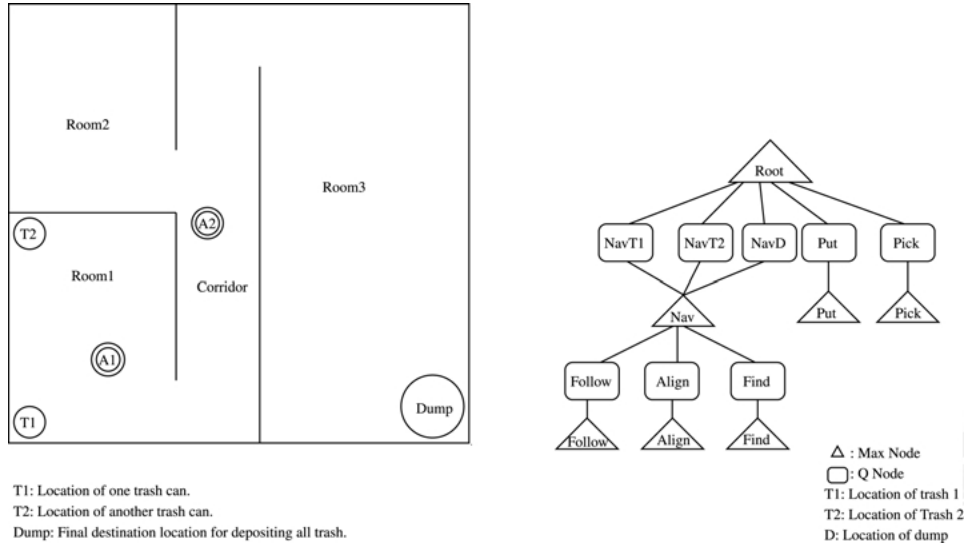


Figure 5. A two-robot (*a1* and *a2*) trash collection task. The robots can learn to coordinate much more rapidly using the task structure than if they attempted to coordinate at the level of primitive movements. On the right is shown a MAXQ graph, which is MAXQ task graph with two kinds of nodes representing subtasks and the actions those subtasks can select Dietterich (2000).

of agent A1 doing *NavT1* subtask itself and the completion sum of the remainder of the overall task done by both agents. In this example, the multiagent MAXQ decomposition embodies the constraint that the value of A1 navigating to trash bin T1 is independent of whatever A2 is doing.

5.3 Hierarchical Memory

In multiagent environments, agents cannot observe joint states and joint actions, but must act based on estimates of these hidden variables (the problem of hidden state occurs in many single-agent tasks as well, such as a robot navigating in an indoor office environment). One approach is formalized in terms of Partially observable Markov decision processes (POMDPs), where agents learn policies over *belief states*, i.e., probability distributions over the underlying state set (Kaelbling et al., 1998). It can be shown that belief states satisfy the Markov property and consequently yield a new (and more complex) MDP over information states. Belief states can be recursively updated using the transition model, and an observation model $O(y | s, a)$ specifying the likelihood of observing y if action a was performed and resulted in state s . However, mapping belief states to optimal actions is known to be intractable, particularly in the decentralized multiagent formulation (Bernstein et al., 2000). Also, learning a perfect model of the underlying POMDP is a challenging task. An empirically more effective (but theoretically less powerful) approach is to use finite memory models as linear chains or nonlinear trees over histories (McCallum, 1996). However, such finite memory structures can be defeated by long sequences of mostly irrelevant observations and actions that conceal a critical past observation.

We briefly summarize three multiscale memory models that have been explored recently by Hernandez and Mahadevan (2001), Theodorou and Mahadevan (2001), and Jonsson and Barto (2001). These models combine temporal abstraction with previous methods for dealing with hidden state. Hierarchical suffix memory (HSM) (Hernandez and Mahadevan, 2001) generalizes the suffix tree model (McCallum, 1996), to SMDP-based temporally-extended activities. Suffix memory constructs state estimators from finite chains of observation-action-reward triples. In addition to extending suffix models to temporally-extended activities, HSM also uses multiple layers of temporal abstraction to form longer-term memories at more abstract levels. Figure 6 illustrates this idea for robot navigation for the simpler case of a linear chain, although the tree-based model has also been investigated. An important side-effect is that the agent *can look back many steps back in time* while ignoring the exact sequence of low-level observations and actions that transpired. Tests in a robot navigation domain showed that HSM outperformed “flat” suffix tree methods, as well as hierarchical methods that used no memory (Hernandez and Mahadevan, 2001).

POMDPs are theoretically more powerful than finite memory models, but past work on POMDPs has mostly studied “flat” models for which learning and planning algorithms scale poorly with model size. Theodorou et al. (2001) developed a *hierarchical POMDP* formalism, termed H-POMDPs (Figure 7), by extending the hierarchical hidden Markov model (HHMM) (Fine et al., 1998) to include rewards and temporally-extended activities.

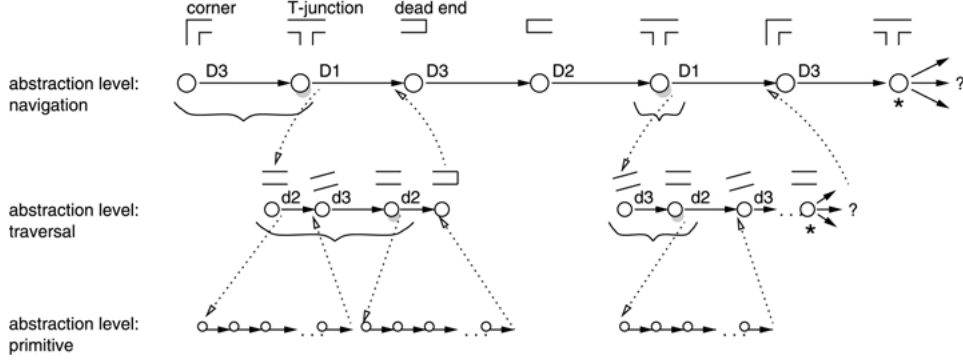


Figure 6. Hierarchical suffix memory state estimator for a robot navigation task. At the abstract (navigation) level, observations and decisions occur at intersections. At the lower (corridor-traversal) level, observations and decisions occur within the corridor. At each level, each agent constructs state representations from its past experience with similar history (shown with shadows).

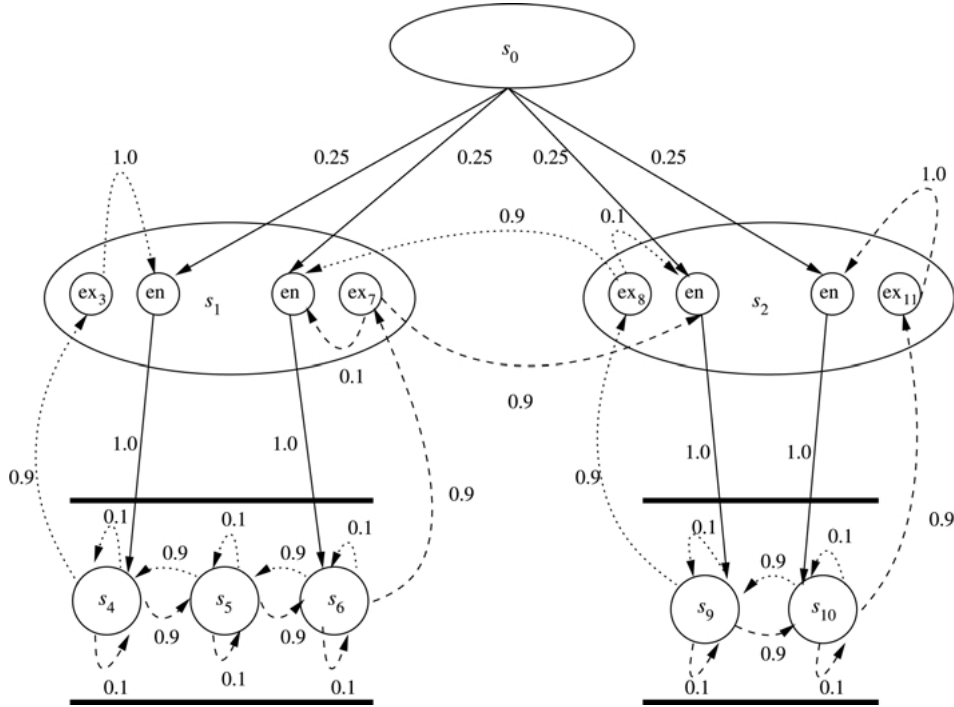


Figure 7. Hierarchical POMDP Model. This hierarchical POMDP model represents two adjacent corridors in a robot navigation task. The model has two primitive actions, “go-left” indicated with the dotted arrows and “go-right” indicated with the dashed arrows. This HPOMDP has two (unobservable) abstract states s_1 and s_2 , and each abstract state has two entry and two exit states. The (hidden) product states s_4, s_5, s_6, s_9 , and s_{10} have associated observation models.

They also developed a hierarchical EM algorithm for learning the parameters of an H-POMDP model from sequences of observations and actions. Extensive tests on a robot navigation domain show learning and planning performance is much improved over flat POMDP models (Theodorou and Mahadevan, 2002; Theodorou et al., 2001). The hierarchical EM-based parameter estimation algorithm scales more gracefully to large models because previously learned sub-models can be reused when learning at higher levels. In addition, the effect of temporally-extended activities in H-POMDPs (e.g., exit the corridor) can be modeled at abstract and product-level states, which supports planning at multiple levels of abstraction.

H-POMDPs have an inherent advantage in allowing belief states to be computed at different levels of the tree. In addition, there is often less uncertainty at higher levels (e.g., a robot is more sure of which corridor it is in than exactly where it is in the corridor). A number of heuristics for mapping belief states to actions provide good performance in robot navigation (e.g., the most-likely-state (MLS) heuristic assumes the agent is in the state corresponding to the “peak” of the belief state distribution) (Koenig and Simmons, 1997; Nourbakhsh et al., 1995; Shatkay and Kaelbling, 1997). Such heuristics work much better in H-POMDPs because they can be applied at multiple levels, and belief states over abstract states usually have lower entropy (Figure 8). For a detailed study of the H-POMDP model, as well as its application to robot navigation, see Theodorou (2002).

Jonsson and Barto (2001) also addressed partial observability by adapting suffix tree methods to hierarchical RL systems. Their approach focused on automating the process of constructing activity-specific state representations by applying McCallum’s U-Tree algorithm (McCallum, 1996) to individual options. The U-Tree algorithm employs the concept of a suffix tree to automatically construct a state representation starting from one that makes no distinctions between different observation vectors. Thus, no specification of

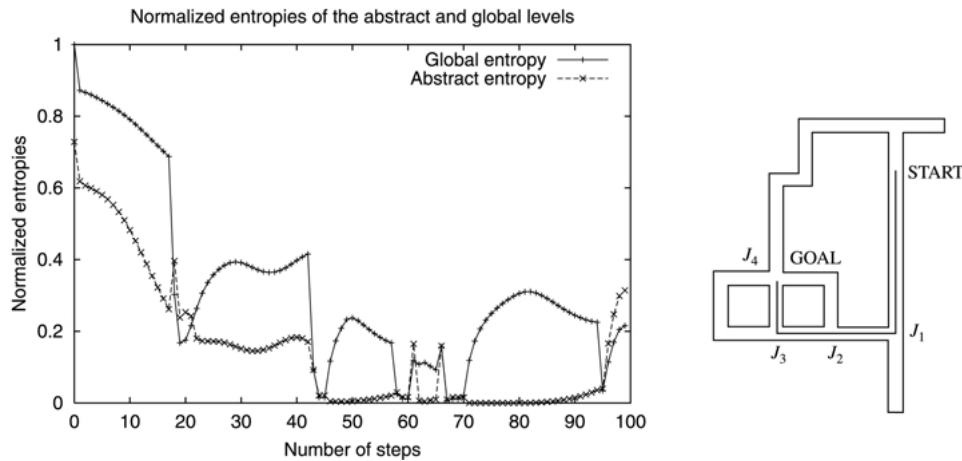


Figure 8. Entropy of a sample robot navigation run. This graph shows a sample robot navigation run whose trace is on the right, where positional uncertainty (measured by belief state entropy) at the abstract (corridor) level is less than at the product state level. Spatiotemporal abstraction reduces the uncertainty and requires less frequent decision-making, allowing the robot to get to goals without initial positional information.

state-feature dependencies is necessary prior to learning. With a separate U-Tree assigned to each option, it is possible to perform state abstraction separately for each option.

The U-Tree algorithm retains a history of transition instances $T_t = \langle T_{t-1}, a_{t-1}, r_t, s_t \rangle$ composed of the observation vector, s_t , at time step t , the previous action, a_{t-1} , the reward, r_t , received during the transition into s_t , and the previous instance, T_{t-1} . A decision tree—the U-Tree—sorts a new instance T_t based on its components and assigns it to a unique leaf $L(T_t)$ of the tree. The distinctions associated with a leaf are determined by the root-to-leaf path. For each leaf-action pair (L_j, a) , the algorithm keeps an action-value $Q(L_j, a)$ estimating the future discounted reward associated with being in L_j and executing a . These action-values can be updated in a variety of ways, such as via a DP algorithm if a system model is available (or can be learned) or RL algorithms such as Q-learning.

The U-Tree algorithm periodically adds new distinctions to the tree in the form of temporary nodes, called fringe nodes, and performs statistical tests to see whether the added distinctions increase the predictive power of the U-Tree. The tree is extended with new distinctions when these are estimated to increase the tree’s predictive power. Whenever the tree is extended, the action-values of the previous leaf node are passed on to the new leaf nodes. Each distinction is based on (1) a perceptual dimension, which is either an observation or a previous action, and (2) a history index, indicating how far back in the current history the dimension will be examined.

Jonsson and Barto (2001) adapted the U-Tree algorithm for use with options and hierarchical learning architectures. Given a finite MDP and a set of options whose policies were not yet defined, they assigned each option a separate U-Tree, which was updated using the U-Tree algorithm (with some modifications) based on each option’s local history. What makes the U-Tree algorithm suitable for performing option-specific state abstraction is that a U-Tree simultaneously defines a state representation and a policy over this representation. By assigning one U-Tree to each option, the algorithm is able to perform state abstraction separately for each option while also modifying its policy. It is also possible to use intra-option learning methods (Section 4.1) so that something about an option can be learned from the behavior of other options.

This version of the U-Tree algorithm was illustrated using Dietterich’s (2000) Taxi task that we briefly described in Section 4.3. Results show that starting with no distinctions being made in the state representation (i.e., all states were represented as a single block), the algorithm was able to solve this task by introducing distinctions that partitioned the state-action set into subsets requiring different action-values, at a significant savings over an approach that initially distinguished between all possible states. Although this example is merely illustrative, it suggests that automated option-specific state abstraction is an attractive approach to making hierarchical learning systems more powerful.

6. Topics for Future Research

6.1. Compact Representations

In most interesting real-world tasks, states have significant internal structure. For example, states are very often represented as vectors of state variables (usually called *factored states*

by machine learning researchers), or even possess richer relational structure (Driessens and Dzeroski, 2002). Much work in artificial intelligence has focused on exploiting this structure to develop compact representations of single-step actions (e.g., the Dynamic Bayes Net representation, Dean and Kanazawa, 1989). A natural question to consider is how to extend these single-step compact models into compact models of temporally-extended activities, such as options. The problem is a bit subtle since even if all actions have limited single-step influence on state variables, this property generally does not hold over an extended activity. One approach that Rohanimanesh and Mahadevan (in preparation) have been studying is how to exploit results from approximation of structured stochastic processes (Boyen and Koller, 1998) to develop structured ways of approximating the next-state predictions of temporally-extended activities. The key idea is that by clustering the state variables into disjoint subsets, and keeping track of a next-state distribution for each local cluster, it is possible to efficiently approximate the underlying next-state distribution for a temporally-extended activity. Preliminary analysis of this approach appears promising, and further theoretical and experimental study is under way.

6.2. *Learning Task Hierarchies*

In the approaches discussed above, the components of the hierarchy, their places in the hierarchy, and the abstractions that are used are decided upon in advance. A key open question is how to form task hierarchies automatically, such as those used in the MAXQ framework. We briefly discussed in Section 4.1 automated methods for identifying useful subgoals (Digney, 1996; Digney, 1998; McGovern, 2002; McGovern and Barto, 2001) which address some aspects of this problem. Another approach called HEXQ was recently proposed by Hengst (2002). It exploits a factored state representation and sorts state variables into an ordered list, beginning with the variable that changes most rapidly. HEXQ builds a task hierarchy, consisting of one level for each state variable, where each level contains a simpler MDP that is connected to other smaller MDPs through a set of “bottleneck” states (for example, if the environment is a set of rooms, as in the navigation example shown in Figure 3 above, each room would correspond to such a smaller MDP). The obvious limitation of HEXQ is that it is limited to considering each state variable in isolation, an approach that fails for more complex problems. Further work is required for understanding how to build task hierarchies in such cases, and how to integrate this approach to related systems approaches such as singular perturbation methods (Kokotovic et al., 1986; Naidu, 1988).

6.3. *Dynamic Abstraction*

Systems such as those outlined in this article naturally provide opportunities for using different state representations depending on the activity that is currently executing. There is a crucial distinction between static abstractions, which remain fixed throughout all phases of a sequential decision task, and what we call *dynamic abstractions* that are conditional on the execution of particular temporally-extended activities. In other words,

the variables that a dynamic abstraction renders relevant or irrelevant are afforded that status only for a temporally-confined segment of time. For example, during the course of driving, the steering wheel angle, the gear position, or the radio status may all be relevant, but become irrelevant during other activities. Dietterich (2000) introduced activity-specific state abstraction in the MAXQ framework, and Jonsson and Barto (2001) explored automatic methods for constructing such representation from experience, as described in Section 5.3. The ability to use dynamic abstraction in RL is one of the key reasons that the hierarchical RL approaches discussed in this article appear so attractive to machine learning researchers. This is an area in which future research can have a significant impact.

6.4. *Large Applications*

Although the proposed ideas for hierarchical RL described above appear promising, to date there has been insufficient experience in experimentally testing the effectiveness of these ideas on large applications. Makar et al. (2001) extended the MAXQ approach to multiagent domains, and applied it to a large multi-vehicle autonomous guided vehicle (AGV) routing problem. They demonstrated that the policies learned for this problem were better than standard heuristics used in industry, such as the “go to the nearest free machine” heuristic. Stone and Sutton (2001) applied the framework of options to a “keep away” task in Robot soccer. This task involves a set of players from one team passing the ball between them and keeping the ball in their possession against the defending opponents. While these initial studies are promising, much further work is necessary to establish the effectiveness of hierarchical RL, particularly on large complex continuous control tasks.

7. Conclusion

It has been our goal in the article to review several closely related approaches to temporal abstraction and hierarchical control that have been developed by machine learning researchers: the options formalism of Sutton, Precup, and Singh, the *hierarchies of abstract machines* (HAMs) approach of Parr and Russell, and Dietterich’s MAXQ framework. We also discussed extensions of these ideas addressing on concurrent activities, multiagent coordination, and hierarchical memory for partial observability. Although many of these ideas are closely related to similar concepts in systems and control engineering on hierarchical, hybrid, and multilayer control, we have not attempted to provide a careful rapprochement between these areas and what machine learning, and other researchers associated with artificial intelligence, have been developing. However, we strongly believe that there is much to be gained—on both sides—from such a rapprochement, and it is our hope that this article will prove to be useful in stimulating the needed dialog.

Acknowledgments

The authors would like to thank Anders Jonsson, Khashayar Rohanimanesh, and Mohammad Ghavamzadeh for their critical reading of an early draft of this article. This work was funded by National Science Foundation under Grant No. ECS-9980062. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the National Science Foundation.

Notes

1. We follow Sutton and Barto (1998) in denoting the reward for the action at stage t by r_{t+1} instead of the more usual r_t .
2. Parr (1998) restricts the HAM call graph to be a tree so that call stack contents do not need to be treated as part of the program state, a point we gloss over in our discussion. This kind of machine hierarchy is an instance of a recursive transition network as discussed by Woods (1970).
3. For multiagent problems, we treat a multi-option as a tuple rather than a set since its elements are associated with specific agents. It is also possible to generalize this so that each o_i is itself a multi-option.

References

- Andre, D., and Russell, S. J. 2001. Programmable reinforcement learning agents. In *Advances in Neural Information Processing Systems: Proceedings of the 2000 Conference*. Cambridge, MA: MIT Press, pp. 1019–1025.
- Barto, A. G., Bradtke, S. J., and Singh, S. P. 1995. Learning to act using real-time dynamic programming. *Artificial Intelligence* 72: 81–138.
- Bernstein, D., Zilberstein, S., and Immerman, N. 2000. The complexity of decentralized control of markov decision processes. In *16th Conference on Uncertainty in Artificial Intelligence*.
- Bertsekas, D. P. 1987. *Dynamic Programming: Deterministic and Stochastic Models*. Prentice-Hall, Englewood Cliffs, NJ.
- Bertsekas, D. P., and Tsitsiklis, J. N. 1996. *Neuro-Dynamic Programming*. Belmont, MA: Athena Scientific.
- Boyer, X., and Koller, D. 1998. Tractable inference for complex stochastic processes. In G. F. Cooper and S. Moral, editors, *Proceedings of the Fourteenth Conference on Uncertainty in AI*. San Francisco, CA, Morgan Kaufmann, pp. 33–42.
- Bradtke, S. J., and Duff, M. O. 1995. Reinforcement learning methods for continuous-time Markov decision problems. In G. Tesauro, D. S. Touretzky, and T. Leen, editors, *Advances in Neural Information Processing Systems: Proceedings of the 1994 Conference*, Cambridge, MA: MIT Press, pp. 393–400.
- Branicky, M. S., Borkar, V. S., and Mitter, S. K. 1998. A unified framework for hybrid control: Model and optimal control theory. *IEEE Transactions on Automatic Control* 43: 31–45.
- Brooks, R. A. 1986. Achieving artificial intelligence through building robots. Technical Report A.I. Memo 899, Massachusetts Institute of Technology Artificial Intelligence Laboratory, Cambridge, MA.
- Crites, R. H. 1996. *Large-Scale Dynamic Optimization Using Teams of Reinforcement Learning Agents*. Ph.D. thesis, Amherst, MA: University of Massachusetts.
- Crites, R. H., and Barto, A. G. 1998. Elevator group control using multiple reinforcement learning agents. *Machine Learning* 33: 235–262.
- Das, T. K., Gosavi, A., Mahadevan, S., and Marchallick, N. 1999. Solving semi-Markov decision problems using average reward reinforcement learning. *Management Science* 45: 560–574.

- Dean, T. L., and Kanazawa, K. 1989. A model for reasoning about persistence and causation. *Computational Intelligence* 5: 142–150.
- Dietterich, T. G. 2000. Hierarchical reinforcement learning with the maxq value function decomposition. *Journal of Artificial Intelligence Research* 13: 227–303.
- Digney, B. 1996. Emergent hierarchical control structures: Learning reactive/hierarchical relationships in reinforcement environments. In P. Meas and M. Mataric, editors, *From Animals to Animats 4: The Fourth Conference on Simulation of Adaptive Behavior*. Cambridge, MA: MIT Press.
- Digney, B. 1998. Learning hierarchical control structure from multiple tasks and changing environments. In *From Animals to Animats 5: The Fifth Conference on Simulation of Adaptive Behavior*. Cambridge, MA: MIT Press.
- Driessens, K., and Dzeroski, S. 2002. Integrating experimentation and guidance in relational reinforcement learning. In *Machine Learning: Proceedings of the Nineteenth International Conference on Machine Learning*.
- Fikes, R. E., Hart, P. E., and Nilsson, N. J. 1972. Learning and executing generalized robot plans. *Artificial Intelligence* 3: 251–288.
- Finc, S., Singer, Y., and Tishby, N. 1998. The hierarchical hidden Markov model: analysis and applications. *Machine Learning* 32(1): July.
- Forestier, J.-P., and Varaiya, P. 1978. Multilayer control of large Markov chains. *IEEE Transactions on Automatic Control* AC-23: 298–304.
- Ghavamzadeh, M., and Mahadevan, S. 2001. Continuous-time hierarchical reinforcement learning. In *Proceedings of the Eighteenth International Conference on Machine Learning*.
- Grudic, G. Z., and Ungar, L. H. 2000. Localizing search in reinforcement learning. In *Proceedings of the 18th National Conference on Artificial Intelligence (AAAI-00)*, pp. 590–595.
- Harel, D. 1987. Statecharts: A visual formalism for complex systems. *Science of Computer Programming* 8: 231–274.
- Hengst, B. 2002. Discovering hierarchy in reinforcement learning with hexq. In *Machine Learning: Proceedings of the Nineteenth International Conference on Machine Learning*.
- Hernandez, N., and Mahadevan, S. 2001. Hierarchical memory-based reinforcement learning. *Proceedings of Neural Information Processing Systems*.
- Howard, R. A. 1971. *Dynamic Probabilistic Systems: Semi-Markov and Decision Processes*. New York: Wiley.
- Huber, M., and Grunert, R. A. 1997. A feedback control structure for on-line learning tasks. *Robotics and Autonomous Systems* 22: 303–315.
- Iba, G. A. 1989. A heuristic approach to the discovery of macro-operators. *Machine Learning* 3: 285–317.
- Jaakkola, T., Jordan, M. I., and Singh, S. P. 1994. On the convergence of stochastic iterative dynamic programming algorithms. *Neural Computation* 6: 1185–1201.
- Jonsson, A., and Barto, A. G. 2001. Automated state abstraction for options using the U-tree algorithm. In *Advances in Neural Information Processing Systems: Proceedings of the 2000 Conference*, Cambridge, MA: MIT Press, pp. 1054–1060.
- Kaelbling, L., Littman, M., and Cassandra A. 1998. Planning and acting in partially observable stochastic domains. *Artificial Intelligence* 101.
- Kaelbling, L. P., Littman, M. L., and Moore, A. W. 1996. Reinforcement learning: A survey. *Journal of Artificial Intelligence Research* 4: 237–285.
- Klopf, A. H. 1974. Brain function and adaptive systems—A heterostatic theory. Technical Report AFCRL-72-0164, Air Force Cambridge Research Laboratories, Bedford, MA, 1972. A summary appears in *Proceedings of the International Conference on Systems, Man, and Cybernetics*, IEEE Systems, Man, and Cybernetics Society, Dallas, TX.
- Klopf, A. H. 1982. *The Hedonistic Neuron: A Theory of Memory, Learning, and Intelligence*. Washington, D.C.: Hemisphere.
- Koenig, S., and Simmons, R. 1997. Xavier: A robot navigation architecture based on partially observable Markov decision process models. In D. Kortenkamp, P. Bonasso, and R. Murphy, editors, *AI-based Mobile Robots: Case-studies of Successful Robot Systems*. Cambridge, MA: MIT Press.
- Kokotovic, P. V., Khalil, H. K., and O'Reilly, J. 1986. *Singular Perturbation Methods in Control: Analysis and Design*. London: Academic Press.

- Korf, R. E. 1985. *Learning to Solve Problems by Searching for Macro-Operators*. Boston, MA: Pitman.
- Littman, M. 1994. Markov games as a framework for multi-agent reinforcement learning. In *Proceedings of the Eleventh International Conference on Machine Learning*, pp. 157–163.
- Mahadevan, S. 1996. Average reward reinforcement learning: Foundations, algorithms, and empirical results. *Machine Learning* 22: 159–196.
- Mahadevan, S., Marchalleck, N., Das, T., and Gosavi, A. 1997. Self-improving factory simulation using continuous-time average-reward reinforcement learning. In *Machine Learning: Proceedings of the Fourteenth International Conference*.
- Makar, R., Mahadevan, S., and Ghavamzadeh, M. 2001. Hierarchical multi-agent reinforcement learning. In J. P. Müller, E. Andre, S. Sen, and C. Frasson, editors, *Proceedings of the Fifth International Conference on Autonomous Agents*, pp. 246–253.
- McCallum, A. K. 1996. *Reinforcement Learning with Selective Perception and Hidden State*. Ph.D. thesis, University of Rochester.
- McGovern, A. 2002. *Autonomous Discovery of Temporal Abstractions from Interaction with An Environment*. Ph.D. thesis, University of Massachusetts.
- McGovern, A., and Barto, A. 2001. Automatic discovery of subgoals in reinforcement learning using diverse density. In C. Brodley and A. Danyluk, editors, *Proceedings of the Eighteenth International Conference on Machine Learning*, San Francisco, CA: Morgan Kaufmann, pp. 361–368.
- Minsky, M. L. 1954. *Theory of Neural-Analog Reinforcement Systems and its Application to the Brain-Model Problem*. Ph.D. thesis, Princeton University.
- Naidu, D. S. 1988. *Singular Perturbation Methodology in Control Systems*. London: Peter Peregrinus Ltd.
- Nourbakhsh, I., Powers, R., and Birchfield, S. 1995. Dervish: An office-navigation robot. *AI Magazine* 16(2): 53–60.
- Parr, R. 1998. *Hierarchical Control and Learning for Markov Decision Processes*. Ph.D. Thesis, Berkeley CA: University of California.
- Parr, R., and Russell, S. 1998. Reinforcement learning with hierarchies of machines. In *Advances in Neural Information Processing Systems: Proceedings of the 1997 Conference*. Cambridge, MA: MIT Press.
- Perkins, T. J., and Barto, A. G. Lyapunov design for safe reinforcement learning. *Journal of Machine Learning Research*. To appear.
- Perkins, T. J., and Barto, A. G. 2001. Lyapunov-constrained action sets for reinforcement learning. In C. Brodley and A. Danyluk, editors, *Proceedings of the Eighteenth International Conference on Machine Learning*. San Francisco, CA: Morgan Kaufmann, pp. 409–416.
- Precup, D. 2000. *Temporal Abstraction in Reinforcement Learning*. Ph.D. thesis, Amherst, MA: University of Massachusetts.
- Precup, D., and Sutton, R. S. 1998. Multi-time models for temporally abstract planning. In: *Advances in Neural Information Processing Systems: Proceedings of the 1997 Conference*. Cambridge MA: MIT Press, pp. 1050–1056.
- Precup, D., Sutton, R. S., and Singh, S. 1998. Theoretical results on reinforcement learning with temporally abstract options. In *Proceedings of the 10th European Conference on Machine Learning, ECML-98*. Springer Verlag, pp. 382–393.
- Puterman, M. L. 1994. *Markov Decision Problems*. New York: Wiley.
- Rohanimesh, K., and Mahadevan, S. Structured approximation of stochastic temporally extended actions. In preparation.
- Rohanimesh, K., and Mahadevan, S. 2001. Decision-theoretic planning with concurrent temporally extended actions. In *Proceedings of the Seventeenth Conference on Uncertainty in Artificial Intelligence*.
- Ross, S. 1983. *Introduction to Stochastic Dynamic Programming*. New York: Academic Press.
- Rummery, G. A., and Niranjan, M. 1994. On-line q-learning using connectionist systems. Technical Report CUED/F-INFENG/TR 166, Cambridge University Engineering Department.
- Samuel, A. L. 1959. Some studies in machine learning using the game of checkers. *IBM Journal on Research and Development* 3: 211–229. Reprinted in E. A. Feigenbaum and J. Feldman, editors, *Computers and Thought*, New York: McGraw-Hill, pp. 71–105.
- Samuel, A. L. 1967. Some studies in machine learning using the game of checkers. II—Recent progress. *IBM Journal on Research and Development* 11: 601–617.

- Schwartz, A. 1993. A reinforcement learning method for maximizing undiscounted rewards. In *Proceedings of the Tenth International Conference on Machine Learning*. Morgan Kaufmann, pp. 298–305.
- Shatkay, H., and Kaelbling, L. P. 1997. Learning topological maps with weak local odometric information. In *IJCAI* (2), pp. 920–929.
- Singh, S., and Bertsekas, D. 1997. Reinforcement learning for dynamic channel allocation in cellular telephone systems. In *Advances in Neural Information Processing Systems: Proceedings of the 1996 Conference*. Cambridge, MA: MIT Press.
- Singh, S., Jaakkola, T., Littman, M. L., and Szepesvári, C. 2000. Convergence results for single-step on-policy reinforcement-learning algorithms. *Machine Learning* 38: 287–308.
- Singh, S. P. 1992. Reinforcement learning with a hierarchy of abstract models. In *Proceedings of the Tenth National Conference on Artificial Intelligence*. Menlo Park, CA: AAAI Press/MIT Press, pp. 202–207.
- Singh, S. P. 1992. Scaling reinforcement learning algorithms by learning variable temporal resolution models. In *Proceedings of the Ninth International Machine Learning Conference*. San Mateo, CA: Morgan Kaufmann, pp. 406–415.
- Stone, P., and Sutton, R. S. 2001. Scaling reinforcement learning toward RoboCup soccer. In C. Brodley and A. Danyluk, editors, *Proceedings of the Eighteenth International Conference on Machine Learning*. San Francisco, CA: Morgan Kaufmann, pp. 537–544.
- Sugawara, T. and Lesser, V. 1998. Learning to improve coordinated actions in cooperative distributed problem-solving environments. *Machine Learning* 33: 129–154.
- Sutton, R. S. 1996. Generalization in reinforcement learning: Successful examples using sparse coarse coding. In D. S. Touretzky, M. C. Mozer, and M. E. Hasselmo, editors, *Advances in Neural Information Processing Systems: Proceedings of the 1995 Conference*. Cambridge, MA: MIT Press, pp. 1038–1044.
- Sutton, R. S., and Barto, A. G. 1981. Toward a modern theory of adaptive networks: Expectation and prediction. *Psychological Review* 88: 135–170.
- Sutton, R. S., and Barto, A. G. 1998. *Reinforcement Learning: An Introduction*. Cambridge, MA: MIT Press.
- Sutton, R. S., Precup, D., and Singh, S. 1999. Between mdps and semi-mdps: A framework for temporal abstraction in reinforcement learning. *Artificial Intelligence* 112: 181–211.
- Tan, M. 1993. Multi-agent reinforcement learning: Independent vs. cooperative agents. In *Proceedings of the Tenth International Conference on Machine Learning*, pp. 330–337.
- Tesauro, G. J. 1992. Practical issues in temporal difference learning. *Machine Learning* 8: 257–277.
- Tesauro, G. J. 1994. TD-gammon, a self-teaching backgammon program, achieves master-level play. *Neural Computation* 6(2): 215–219.
- Theocharous, G. 2002. *Hierarchical Learning and Planning in Partially Observable Markov Decision Processes*. Ph.D. Thesis, Michigan State University.
- Theocharous, G., and Mahadevan, S. 2002. Approximate planning with hierarchical partially observable Markov decision process for robot navigation. In *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*.
- Theocharous, G., Rohanimanesh, K., and Mahadevan, S. 2001. Learning hierarchical partially observable markov decision process for robot navigation. In *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*.
- Thrun, S. B., and Schwartz, A. 1995. Finding structure in reinforcement learning. In G. Tesauro, D. S. Touretzky, and T. Leen, editors, *Advances in Neural Information Processing Systems: Proceedings of the 1994 Conference*. Cambridge MA: MIT Press, pp. 385–392.
- Tsitsiklis, J. N., and Van Roy, B. 1997. An analysis of temporal-difference learning with function approximation. *IEEE Transactions on Automatic Control* 42: 674–690.
- Watkins, C. J. C. H. 1989. *Learning from Delayed Rewards*. Ph.D. thesis. Cambridge, U.K.: Cambridge University.
- Watkins, C. J. C. H., and Dayan, P. 1992. Q-learning. *Machine Learning* 8: 279–292.
- Weiss, G. 1999. *Multiagent Systems: A Modern Approach to Distributed Artificial Intelligence*. Cambridge, MA: MIT Press.
- Werbos, P. J. 1977. Advanced forecasting methods for global crisis warning and models of intelligence. *General Systems Yearbook* 22: 25–38.

- Werbos, P. J. 1987. Building and understanding adaptive systems: A statistical/numerical approach to factory automation and brain research. *IEEE Transactions on Systems, Man, and Cybernetics* 17: 7–20.
- Werbos, P. J. 1992. Approximate dynamic programming for real-time control and neural modeling. In D. A. White and D. A. Sofge, editors, *Handbook of Intelligent Control: Neural, Fuzzy, and Adaptive Approaches*. New York: Van Nostrand Reinhold, pp. 493–525.
- Woods, W. A. 1970. Transition network grammars for natural language analysis. *Communications of the ACM* 13: 591–606.