

---

# Reinforcement Learning with Hierarchies of Machines

---

**Ron Parr and Stuart Russell**

Computer Science Division, UC Berkeley, CA 94720  
{parr,russell}@cs.berkeley.edu

## Abstract

We present a new approach to reinforcement learning in which the policies considered by the learning process are constrained by hierarchies of partially specified machines. This allows for the use of prior knowledge to reduce the search space and provides a framework in which knowledge can be transferred across problems and in which component solutions can be recombined to solve larger and more complicated problems. Our approach can be seen as providing a link between reinforcement learning and “behavior-based” or “teleo-reactive” approaches to control. We present provably convergent algorithms for problem-solving and learning with hierarchical machines and demonstrate their effectiveness on a problem with several thousand states.

**Category:** reinforcement learning. **Preference:** plenary.

## 1 Introduction

Optimal decision making in virtually all spheres of human activity is rendered intractable by the complexity of the task environment. Generally speaking, the only way around intractability has been to provide a hierarchical organization for complex activities. Although it can yield suboptimal policies, top-down hierarchical control often reduces the complexity of decision making from exponential to linear in the size of the problem. For example, hierarchical task network (HTN) planners can generate solutions containing tens of thousands of steps [4], whereas “flat” planners can manage only tens of steps.

HTN planners are successful because they use a plan library that describes the decomposition of high-level activities into lower-level activities. This paper describes an approach to learning and decision making in *uncertain* environments (Markov decision processes) that uses a roughly analogous form of prior knowledge. We use *hierarchical abstract machines* (HAMs), which impose constraints on the policies considered by our learning algorithms. HAMs consist of nondeterministic finite state machines whose transitions may invoke lower-level machines. Nondeterminism is represented by *choice states* where the optimal action is yet to be decided or learned. The language allows a variety of prior constraints to be expressed, ranging from no constraint all the way to a fully specified solution. One useful intermediate point is the specification of just the general organization of behavior into a layered hierarchy, leaving it up to the learning algorithm to discover exactly which lower-level activities should be invoked by higher levels at each point.

The paper begins with a brief review of Markov decision processes (MDPs) and a description of hierarchical abstract

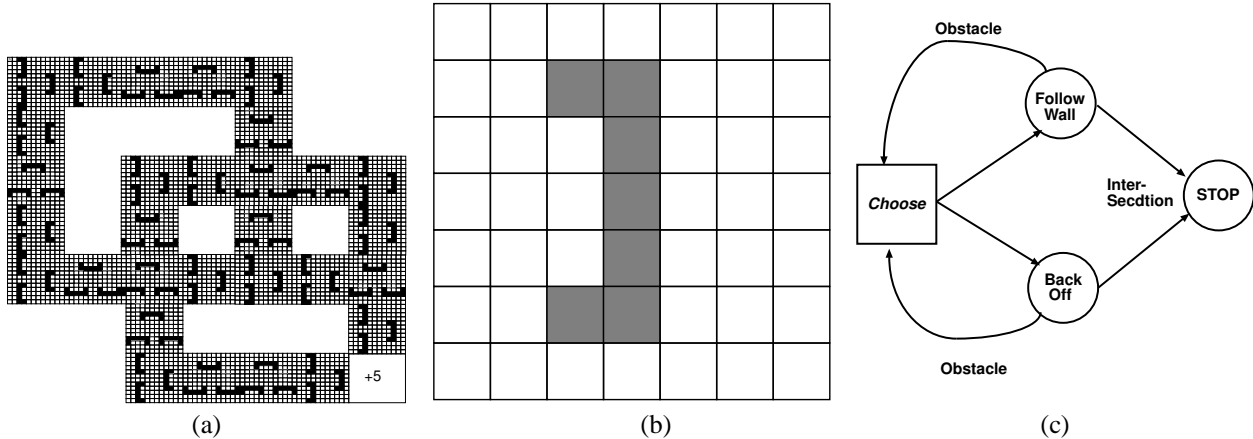


Figure 1: (a) An MDP with  $\approx 3600$  states. (b) Close-up showing a typical obstacle. (c) Nondeterministic finite-state controller for negotiating obstacles.

machines. We then present, in abbreviated form, the following results:

- Given any HAM and any MDP, there exists a new MDP such that the optimal policy with respect to the new MDP is optimal in the original MDP among those policies that satisfy the constraints specified by the HAM. This means that even with complex machine specifications we can still apply standard decision-making and learning methods.
- An algorithm exists that determines this optimal policy, given an MDP and a HAM.
- On an illustrative problem with 3600 states, this algorithm yields dramatic performance improvements over standard algorithms applied to the original MDP.
- A reinforcement learning algorithm exists that converges to the optimal policy, subject to the HAM constraints, with no need to construct the HAM-induced MDP.
- On the sample problem, this algorithm learns dramatically faster than standard RL algorithms.

We conclude with a discussion of related approaches and ongoing work.

## 2 Markov Decision Processes

We assume the reader is familiar with the basic concepts of MDPs (see, e.g., [10]). We will use the following notation: An MDP is a 4-tuple,  $(\mathcal{S}, \mathcal{A}, T, R)$  where  $\mathcal{S}$  is a set of *states*,  $\mathcal{A}$  is a set of *actions*,  $T$  is a *transition model* mapping  $\mathcal{S} \times \mathcal{A} \times \mathcal{S}$  into probabilities in  $[0, 1]$ , and  $R$  is a *reward function* mapping  $\mathcal{S} \times \mathcal{A} \times \mathcal{S}$  into real-valued rewards. Algorithms for solving MDPs can return a *policy*  $\pi$  that maps from  $\mathcal{S}$  to  $\mathcal{A}$ , a real-valued *value function*  $V$  on states, or a real-valued *Q-function* on state–action pairs. In this paper, we focus on infinite-horizon MDPs with a discount factor  $\beta$ . The aim is to find an optimal policy  $\pi^*$  (or, equivalently,  $V^*$  or  $Q^*$ ) that maximizes the expected discounted total reward of the agent.

Throughout the paper, we will use as an example the MDP shown in Figure 1(a). Here  $\mathcal{A}$  contains four primitive actions (up, down, left, right). The model  $T$  specifies that each action succeeds 80% of time, while 20% of the time the agent moves in an unintended perpendicular direction. The agent begins in a start state in the upper left corner. A reward of 5.0 is given for reaching the goal state and the discount factor  $\beta$  is 0.999.

## 3 Hierarchical abstract machines

An abstract machine can be viewed as a *constraint* on policies. For example, the machine described as “repeatedly choose right or down” eliminates from consideration all policies that go up or left. The machine “repeatedly choose

right,” applied to any particular MDP, allows only a single policy, whereas the machine “repeatedly choose any action” does not constrain the policy at all. The idea of the HAM language is to allow specification of constraints on behavior at whatever level of detail and specificity is appropriate for the set of problems to be faced.

The HAM language allows one to specify more interesting machines than those in the preceding paragraph. Generally speaking, a HAM consists of a set of named, nondeterministic, finite-state machines.<sup>1</sup> Each machine is defined by a set of states, a transition function, and a start function that determines the initial state of the machine. HAM states are of four types:

**Action** states execute an action in the environment.

**Call** states execute another HAM as a subroutine.

**Choice** states nondeterministically select a next machine state.

**Stop** states halt execution of the machine and return control to the previous call state.

The transition function determines the next HAM state after an action or call state as a (possibly stochastic) function of the current HAM state and some features of the resulting environment state. HAMs will typically use a partial description of the environment to determine the next state. While this makes HAMs suitable for partially observable domains, for the purposes of this paper we make the standard assumption that the agent has access to a complete description as well.

As an example of the kinds of HAM controllers we can express in this language, Figure 1(c) shows one element of the HAM we used for the MDP in Figure 1. This element is used for traversing a hallway while negotiating obstacles of the kind shown in Figure 1(b). It runs until the end of the hallway or an intersection is reached. When it encounters an obstacle, a choice point is created to choose between two possible next states. One calls the backoff machine to back away from the obstacle and then move forward until the next one. The other calls the follow-wall machine to try to get around the obstacle. The follow-wall machine is very simple and will be tricked by obstacles that are concave in the direction of intended movement; the backoff machine, on the other hand, can move around any obstacle in this world but wastes time backing away from some obstacles unnecessarily and should be used sparingly.

Our complete “navigation HAM” involves a three-level hierarchy, somewhat reminiscent of a Brooks-style architecture but with hard-wired decisions replaced by choice states. The top level of the hierarchy is basically just a choice state for choosing a hallway navigation direction from the four coordinate directions. This machine has control initially and regains control at intersections or corners. The second level of the hierarchy contains four machines for moving along hallways, one for each direction. Each machine at this level has a choice state with four basic strategies for handling obstacles. Two back away from obstacles and two attempt to follow walls to get around obstacles. The third level of the hierarchy implements these strategies, using the primitive actions to do so.

The transition function for this HAM assumes that an agent executing the HAM has access to a short-range, low-directed sonar that detects obstacles in any of the four axis-parallel adjacent squares and a long-range, high-directed sonar that detects larger objects such as the intersections and the ends of hallways. The HAM uses these partial state descriptions to identify feasible choices. For example, the machine to traverse a hallway northwards would not be called from the start state because the HAM high-directed sonar would detect a wall to the north.

In summary, a HAM represents an abstract description of how an agent should act in an environment. Our navigation HAM represents an abstract plan to move about by repeatedly selecting a direction and pursuing this direction until an intersection is reached. Each machine for navigating in the chosen direction also represents an abstract plan for moving in a particular direction while avoiding obstacles. The next section defines how a HAM interacts with a specific MDP and how to find an optimal policy that respects the HAM constraints.

## 4 Defining and solving the HAM-induced MDP

The key concept in our formal analysis is that of *applying* a HAM  $H$  to an MDP  $M$  to yield an *induced* MDP  $H \circ M$ . For example, if  $H$  is “repeatedly choose right or down” and  $M$  is the MDP in Figure 1(a), then  $H \circ M$  is an MDP with the same states but with just two possible actions in each state. If  $H$  is “repeatedly choose right,”  $H \circ M$  is a degenerate

---

<sup>1</sup>HAMs for reinforcement learning, which does not require state enumeration, can also use machines with a countably infinite number of states, e.g., Turing machines.

MDP with only one action—effectively a Markov chain. If  $H$  is “repeatedly choose any action,” then  $H \circ M = M$ . The *choice points* in the induced MDP are those states where the agent still has to choose an action.

When the HAM has internal state, things get more complex. A somewhat simplified description of the construction of  $H \circ M$  is as follows:

- The set of states in  $H \circ M$  is the cross-product of the states of  $H$  with the states of  $M$ .
- For each state in  $H \circ M$ , the action set is just the null action unless the corresponding HAM state is a choice state, in which case the action set is the set of choices, or an action state, in which case the action is the specified primitive action.
- The transition model is taken from  $M$  for primitive actions and taken from  $H$  otherwise.
- The reward is taken from  $M$  for primitive actions, otherwise it is zero.

With this construction, we have the following (proof omitted):

**Lemma 1** For any Markov decision process  $M$  and any<sup>2</sup> HAM  $H$ , the induced system  $H \circ M$  is a Markov decision process.

**Lemma 2** If  $\pi$  is an optimal policy for  $H \circ M$ , then the primitive actions specified by  $\pi$  constitute the optimal policy for  $M$  that is consistent with  $H$ .

Of course,  $H \circ M$  may be quite large if it really contains the full cross-product of the environment and HAM states. Fortunately, there are two things that will make the problem much easier in most cases. The first is that not all pairs of HAM states and environment states will be possible, i.e., reachable from an initial state. The second is that the actual complexity of the induced MDP is determined by the number of choice points, i.e., states of  $H \circ M$  in which the HAM component is a choice state. This leads to the following theorem:

**Theorem 1** For any MDP  $M$  and HAM  $H$ , let  $\mathcal{C}$  be the set of choice points in  $H \circ M$ . Then there exists an MDP, which we will call  $reduce(H \circ M)$ , with states  $\mathcal{C}$  such that the optimal policy for  $reduce(H \circ M)$  corresponds to the optimal policy for  $M$  that is consistent with  $H$ .

**Proof sketch** We begin by applying Lemma 1 and then observing that in states of  $H \circ M$  where the HAM component is not a choice state, only one action is permitted. If we allow different transitions in the MDP to have different discount factors (a straightforward extension of the standard MDP framework) then we can remove these states from the extended state space with no change in the optimal policy or value function for the extended problem. We can thus remove all non-choice states from  $H \circ M$  until only the choice points remain.  $\square$

This theorem formally establishes the mechanism by which the constraints embodied in a HAM can be used to simplify an MDP. As an example of the power of this theorem, and to demonstrate that this transformation can be done efficiently, we applied our navigation HAM to the problem described in the previous section. Figure 2(a) shows the results of applying policy iteration to the original model and to the transformed model. Even when we add in the cost of transformation (which, with our rather underoptimized code, takes 866 seconds), the HAM method produces a HAM-optimal policy in less than a quarter of the time required to find the optimal policy in the original model. The actual solution time is 185 seconds versus 4544 seconds.

We conclude this section by pointing out one other important property of the HAM approach. The HAM model transformation produces an MDP that is an accurate model of the application of the HAM to the original MDP. Unlike typical approximation methods for MDPs, the HAM method can give strict performance guarantees. The solution to the transformed model  $reduce(H \circ M)$  is the optimal solution from within a well-defined class of policies and the value assigned to this solution is the true expected value of applying the concrete HAM policy to the original MDP.

## 5 Reinforcement learning with HAMs

HAMs can be of even greater advantage in a reinforcement learning context, where the effort required to obtain a solution typically scales very badly with the size of the problem. An agent operating with a HAM can use the

---

<sup>2</sup>To preserve the Markov property, we require that if a machine has more than one possible caller in the hierarchy, that each appearance is treated as a distinct machine. This is equivalent to requiring that the call graph for the HAM is a tree.

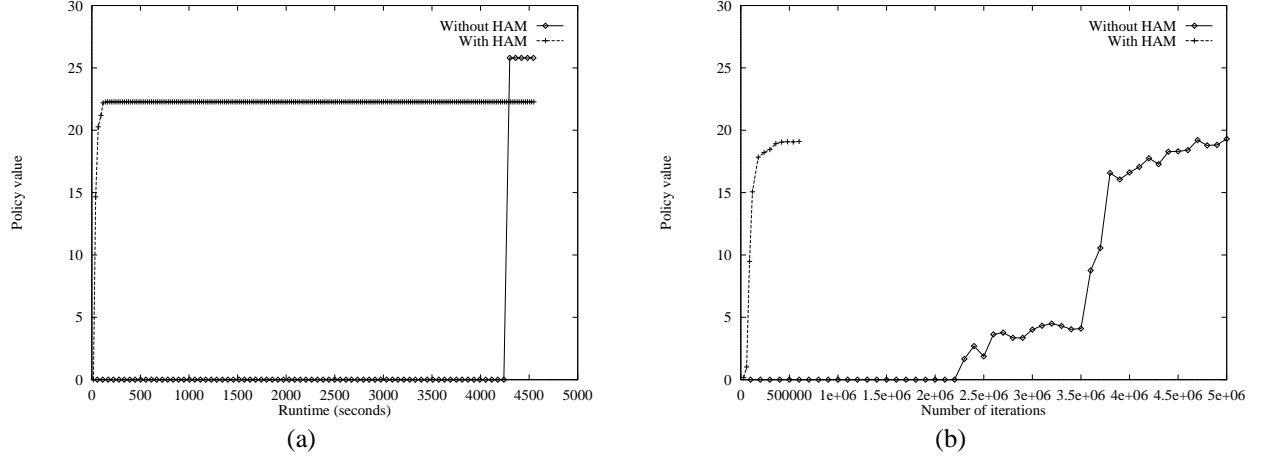


Figure 2: Experimental results showing policy value as a function of runtime on the domain shown in Figure 1. (a) Policy iteration with and without the HAM. (b) Q-learning with and without the HAM (averaged over 10 runs).

constraints implicit in the machines to focus its exploration of the state space, reducing the “blind search” phase that reinforcement learning agents must endure while learning about a new environment. The agent will also learn more quickly for the same reason policy iteration is faster in the HAM-induced model; they agent is effectively operating in a reduced state space.

We now show that it is possible for an agent to use a variation of Q-learning called HAMQ-learning that learns directly in the reduced state space *without performing the model transformation* described in the previous section. This is significant because the explicit transformation requires knowledge of the environment model, something that is not usually known *a priori* in reinforcement learning contexts.

A HAMQ-learning agent keeps track of the following quantities:  $t$ , the current environment state;  $n$ , the current machine state;  $s_c$  and  $m_c$ , the environment state and machine state at the previous choice point;  $a$ , the choice made at the previous choice point; and  $r_c$  and  $\beta_c$ , the total accumulated reward and discount since the previous choice point. It also maintains an extended Q-table,  $Q([s, m], a)$ , which is indexed by an environment-state/machine-state pair and by an action taken at a choice point.

For every environment transition from state  $s$  to state  $t$  with observed reward  $r$  and discount  $\beta$ , the HAMQ-learning agent does the following:

$$\begin{aligned} r_c &\leftarrow r_c + \beta_c r \\ \beta_c &\leftarrow \beta \beta_c \end{aligned}$$

For each transition to a choice point, the agent does

$$\begin{aligned} Q([s_c, m_c], c) &\leftarrow Q([s_c, m_c], c) + \alpha [r_c + \beta_c * V([t, n]) - Q([s_c, m_c], c)] \\ r_c &\leftarrow 0 \\ \beta_c &\leftarrow 1 \end{aligned}$$

This leads to the following theorem:

**Theorem 2** For any MDP  $M$  and any HAM  $H$ , HAMQ-learning will converge to the optimal action choice for every choice point in  $reduce(H \circ M)$  with probability 1.

**Proof sketch** We note that the expected reinforcement signal from the HAMQ-learning rules is the same as the expected reinforcement signal that would be received if the agent were acting directly in the transformed model of Theorem 1 above. Thus, Theorem 1 of [8] can be applied to prove the convergence of the HAMQ-learning agent, provided that we enforce suitable constraints on the exploration strategy and the update parameter decay rate.  $\square$

We ran some experiments to measure the performance of HAMQ-learning on our sample problem. Exploration was

achieved by selecting actions according to the Boltzman distribution with a temperature parameter for each state. We also used an inverse decay for the update parameter  $\alpha$ . Figure 2(b) compares the learning curves for Q-learning and HAMQ-learning. HAMQ-learning appears to learn much faster: Q-learning required 9,000,000 iterations to reach the level achieved by HAMQ-learning after 270,000 iterations. Even after 20,000,000 iterations, Q-learning did not find a policy as good as the HAM-optimal policy.<sup>3</sup>

## 6 Related work

This necessarily brief section discusses approaches to hierarchical decomposition in offline and online MDP methods. As yet, no systematic understanding exists of the space of possible approaches, and various approaches may apply in various contexts. We can, however, identify two common threads: *environment decomposition* and *behavior decomposition*.

*State aggregation* (see, e.g., [12]) is an example of environment decomposition in which “similar” states are clustered together and assigned the same value. This approach effectively reduces the size of the state space. Since it treats distinct states as a single state, it can violate the Markov property, leading to oscillation, divergence, and the loss of performance guarantees. Moreover, state aggregation in this form is hard to apply in many cases where natural bases for aggregation such as physical location can cluster together states with sharply varying values. Other work on environment decomposition includes that of Dean and Lin [6] and Bertsekas and Tsitsiklis [2], who showed that some MDPs are loosely coupled and hence amenable to divide-and-conquer algorithms.

In reinforcement learning, Dayan and Hinton [5] have proposed *feudal* RL which specifies an explicit subgoal structure, with fixed values for each subgoal achieved, in order to achieve a hierarchical decomposition of the state space. This seems natural in some domains, but it is not clear how the subgoal values should be related to the overall reward function. Singh [11] focuses on the simultaneous learning of policies to achieve subgoals and an ordering on these subgoals that achieves a higher level goal. Subgoal structure is encoded in a reward function that forces the agent to solve the subgoals in the proper sequence.

Behavioral decomposition seems to avoid the problem of creating non-Markov processes that plagues state-based approaches. Sutton [13] proposes *temporal abstractions*, which concatenate sequences of state transitions together to permit reasoning about temporally extended events, and which can thereby form a behavioral hierarchy. Lin’s somewhat informal scheme [9] also allows agents to treat entire policies as single actions. Both of these approaches can be seen as special cases of our framework.

The design of hierarchically organized, “layered” controllers was popularized by Brooks [3]. His designs use a somewhat different means of passing control, but our analysis and theorems apply equally well to his machine description language. The “teleo-reactive” agent designs of Benson and Nilsson [1] are even closer to our HAM language. Both of these approaches assume that the agent is completely specified, albeit self-modifiable. The idea of partial behavior descriptions can be traced at least to Hsu’s *partial programs* [7], which were used with a deterministic logical planner.

## 7 Conclusions and future work

We have presented Hierarchies of Abstract Machines (HAMs) as a principled means of constraining the set of policies that are considered for a Markov decision process and we have demonstrated the efficacy of this approach in a simple example for both policy iteration and reinforcement learning. Our results show very significant speedup for decision-making and learning—but of course, this reflects the provision of knowledge in the form of the HAM. The HAM language provides a very general method of transferring knowledge to an agent and we have only scratched the surface of what can be done with this approach. We believe that our formal techniques and algorithmic approach can also be applied to a variety of other languages for describing behavioral constraints.

The similarities between HAMs and finite-state controllers make them a natural and comfortable vehicle for conveying many types of knowledge. We believe that if desired, subgoal information can be incorporated into the HAM

---

<sup>3</sup>Speedup techniques such as eligibility traces could be applied to get better Q-learning results; such methods apply equally well to HAMQ-learning.

structure, unifying subgoal-based approaches with the HAM approach. Moreover, the HAM structure provides a natural decomposition of the HAM-induced model, making it amenable to the divide-and-conquer approaches of [6] and [2].

There are opportunities for generalization across all levels of the HAM paradigm. Value function approximation can be used for the HAM induced model and inductive learning methods can be used to produce HAMs or to generalize their effects upon different regions of the state space. Gradient-following methods can also be used to adjust the transition probabilities of a stochastic HAM.

HAMs also lend themselves naturally to partially observable domains. They can be applied directly when the choice points induced by the HAM are states where no confusion about the agent's true state is possible. The application of HAMs to more general partially observable domains is more complicated and is a topic of ongoing research. We also believe that the HAM approach can be extended to cover the average-reward optimality criterion.

We expect that successful pursuit of these lines of research will provide a formal basis for understanding and unifying several seemingly disparate approaches to control, including behavior-based methods. It should also enable the use of the MDP framework in real-world applications of much greater complexity than hitherto attacked, much as HTN planning has extended the reach of classical planning methods.

## References

- [1] S. Benson and N. Nilsson. Reacting, planning and learning in an autonomous agent. In K. Furukawa, D. Michie, and S. Muggleton, editors, *Machine Intelligence 14*. Oxford University Press, Oxford, 1995.
- [2] D. C. Bertsekas and J. N. Tsitsiklis. *Parallel and Distributed Computation: Numerical Methods*. Prentice-Hall, Englewood Cliffs, New Jersey, 1989.
- [3] Rodney A. Brooks. A robust layered control system for a mobile robot. *IEEE Journal of Robotics and Automation*, 2:14–23, 1986.
- [4] K. W. Currie and A. Tate. O-Plan: the Open Planning Architecture. *Artificial Intelligence*, 52(1):49–86, November 1991.
- [5] Peter Dayan and Geoffrey E. Hinton. Feudal reinforcement learning. In Stephen Jose Hanson, Jack D. Cowan, and C. Lee Giles, editors, *Neural Information Processing Systems 5*, pages 361–368, San Mateo, California, 1993. Morgan Kaufmann.
- [6] Thomas Dean and Shieu-Hong Lin. Decomposition techniques for planning in stochastic domains. In *Proceedings of the Fourteenth International Joint Conference on Artificial Intelligence (IJCAI-95)*, Montreal, Canada, August 1995. Morgan Kaufmann.
- [7] Yung-Jen Hsu. Synthesizing efficient agents from partial programs. In *Methodologies for Intelligent Systems: 6th International Symposium, ISMIS '91, Proceedings*, pages 142–151, Charlotte, North Carolina, October 1991. Springer-Verlag.
- [8] T. Jaakkola, M.I. Jordan, and S.P. Singh. On the convergence of stochastic iterative dynamic programming algorithms. *Neural Computation*, 6(6):1185–1201, 1994.
- [9] Long-Ji Lin. *Reinforcement Learning for Robots Using Neural Networks*. PhD thesis, Computer Science Department, Carnegie-Mellon University, Pittsburgh, Pennsylvania, 1993.
- [10] Martin L. Puterman. *Markov decision processes: Discrete stochastic dynamic programming*. Wiley, New York, 1994.
- [11] S. P. Singh. Transfer of learning by composing solutions of elemental sequential tasks. *Machine Learning*, 8(3):323–340, May 1992.
- [12] Satinder P. Singh, Tommi Jaakkola, and Michael I. Jordan. Reinforcement learning with soft state aggregation. In G. Tesauro, D. Touretzky, and T. Leen, editors, *Advances in Neural Information Processing Systems 7*, pages 361–368, Cambridge, Massachusetts, 1995. MIT Press.
- [13] R. Sutton. Temporal abstraction in reinforcement learning. In *Proceedings of the Twelfth International Conference on Machine Learning*, Tahoe City, CA, July 1995. Morgan Kaufmann.