

Learning Hierarchical Skills from Observation

Ryutaro Ichise,^{1,2} Daniel Shapiro,¹ and Pat Langley¹

¹ Computational Learning Laboratory
Center for the Study of Language and Information
Stanford University, Stanford CA 94305-4115 USA

² National Institute of Informatics, Tokyo 101-8430 Japan
{ichise,dgs,langley}@csl.stanford.edu

Abstract. This paper addresses the problem of learning control skills from observation. In particular, we show how to infer a hierarchical, reactive program that reproduces and explains the observed actions of other agents, specifically the elements that are shared across multiple individuals. We infer these programs using a three-stage process that learns flat unordered rules, combines these rules into a classification hierarchy, and finally translates this structure into a hierarchical reactive program. The resulting program is concise and easy to understand, making it possible to view program induction as a practical technique for knowledge acquisition.

1 Introduction

Physical agents like humans not only execute complex skills but also improve their ability over time. The past decade has seen considerable progress on computational methods for learning such skills and control policies from experience. Much of this research has focused on learning through trial and error exploration, but some has addressed learning by observing behavior of another agent on the task. In particular, research on *behavioral cloning* (e.g., Sammut, 1996) has shown the ability to learn reactive skills through observation on challenging control problems like flying a plane and driving an automobile.

Although such methods can produce policies that predict accurately the desirable control actions, they ignore the fact that complex human skills often have a hierarchical organization. This structure makes the skills more understandable and more transferable to other tasks. In this paper, we present a new approach to learning reactive skills from observation that addresses the issue of inferring their hierarchical structure. We start by specifying the learning task, including the training data and target representation, then present a method for learning hierarchical skills. After this, we report an experimental evaluation of our method that examines the accuracy of the learned program and its similarity to a source program that generated the training cases. In closing, we discuss related work and directions for future research on this topic.

2 The Task of Learning Hierarchical Skills

We define the task of learning skills in terms of its inputs and outputs:

- Given: a trace of agent behavior containing feature-action pairs;
- Find: a program that generates the same actions when presented with the same features.

Research on behavioral cloning (e.g., Anderson et al., 2000; Sammut, 1996) has already addressed this task, having developed methods that learn reactive skills from observation that are both accurate and comprehensible. However, complex skills can often be decomposed naturally into subproblems, and here we focus on capturing this hierarchical structure in an effort to produce even more concise and understandable policies.

We increase the generality of this learned structure by adopting the *separation hypothesis* (Shapiro & Langley, 2002), which asserts that differences in individual behavior are due to the action of distinct preferences over the same set of skills. That is, we all know how to perform common tasks like driving, but some prefer more safe, and others the more reckless options. This assumption separates the task of program acquisition into two parts, where the first is to acquire the necessary structure of skills, and the second is to resolve a (possibly numeric) representation of preference that explains individual choices. We address the first task here. The separation hypothesis simplifies the task of program acquisition because it implies that we should learn a non-deterministic mapping from the observed situation to a feasible set of actions, instead of aiming for a deterministic characterization of a single agent’s behavior. The resulting program will represent fewer distinctions, and thus will be easier to understand.

2.1 Nature of the Training Data

We assume that the learner observes traces of another agent’s behavior as it executes skills on some control task. As in earlier work on learning skills from observation, these traces consist of a sequence of environmental situations and the corresponding agent action. However, since our goal is to recover a non-deterministic mapping, we consider traces from multiple agents that collectively exhibit the full range of available options. Moreover, since we are learning reactive skills, we transform the observed sequences into an unordered set of training cases, one for each situation.

Traditional work in behavioral cloning turns an observational trace into training cases for supervised learning, treating each possible action as a class value. In contrast, we find sets of actions that occur in the same environmental situation and generate training cases that treat each observed action set as a class value. This lets us employ standard methods for supervised induction to partition situations into reactive but nondeterministic control policies.

2.2 Nature of the Learned Skills

We assume that learned skills are stated in ICARUS (Shapiro, 2001), a hierarchical reactive language for specifying the behavior of physical agents that encodes contingent mappings from situations to actions. Like other languages of this kind (Brooks, 1986; Firby, 1989; Georgeff et al., 1985), ICARUS interprets programs in a repetitive sense-think-act loop that lets an agent retrieve a relevant action even if the world changes from one cycle of the interpreter to the next. ICARUS shares the logical orientation of teleoreactive trees (Nilsson, 1994) and universal plans (Schoppers, 1987), but adds vocabulary for expressing hierarchical intent and non-deterministic choice, as well as tools for problem decomposition found in more general-purpose languages. For example, ICARUS supports function call, Prolog-like parameter passing, pattern matching on facts, and recursion. We discuss a simple ICARUS program in the following section.

2.3 An Icarus Plan for Driving

An ICARUS program is a mechanism for finding a goal-relevant reaction to the situation at hand. The primitive building block, or plan, contains up to three elements: an objective, a set of requirements (or preconditions), and a set of alternate means for accomplishing the objective. Each of these can be instantiated by further ICARUS plans, creating a logical hierarchy that terminates with calls to primitive actions or sensors. ICARUS evaluates these fields in a situation-dependent order, beginning with the objective field. If the objective is already true in the world, evaluation succeeds and nothing further needs to be done. If the objective is false, the interpreter examines the requirements field to determine if the preconditions for action have been met. If so, evaluation progresses to the means field, which contains alternate methods (subplans or primitive actions) for accomplishing the objective. The means field is the locus of all choice in ICARUS. Given a value function that encodes a user's preferences, the system learns to select the alternative that promises the largest expected reward.

Table 1 presents an ICARUS plan for freeway driving. The top-level routine, Drive, contains an ordered set of objectives implemented as further subplans. ICARUS repetitively evaluates this program, starting with its first clause every execution cycle. The first clause of Drive defines a reaction to an impending collision. If this context applies, ICARUS returns the Slam-on-brakes action for application in the world. However, if Emergency-brake is not required, evaluation proceeds to the second clause, which encodes a reaction to trouble ahead, defined as a car traveling slower than the agent in the agent's own lane. This subplan contains multiple options. It lets the agent move one lane to the left, move right, slow down, or cruise at its current speed. ICARUS makes a selection based on the long-term expected reward of each alternative. The remainder of the program follows a similar logic as the interpreter considers each clause of Drive in turn. If a clause returns True, the system advances to the next term. If it returns False, Drive would exit with False as its value. However, ICARUS supports a third option: a clause can return an action, which becomes the return value of the

Table 1. The ICARUS program for freeway driving.

<pre> Drive() :objective [*not*(Emergency-brake()) *not*(Avoid-trouble-ahead()) Get-to-target-speed() *not*(Avoid-trouble-behind()) Cruise()] Emergency-brake() :requires [Time-to-impact() <= 2] :means [Slam-on-brakes()] Avoid-trouble-ahead() :requires [?c = Car-ahead-center() Velocity() > Velocity(?c)] :means [Safe-cruise() Safe-slow-down() Safe-change-left() Safe-change-right()] Get-to-target-speed() :objective [Near(Velocity(), Target-speed())] :means [Adjust-speed-if-lane-clear() Adjust-speed-if-car-in-front()] Avoid-trouble-behind() :requires ;;faster car behind [?c = Car-behind-center() Velocity(?c) > Velocity()] :means [Safe-cruise() Safe-change-right()] Safe-cruise() :requires [Time-to-impact() > 2] :means [Cruise()] </pre>	<pre> Safe-slow-down() :requires [Time-to-impact(-2) > 2] :means [Slow-down()] Safe-speed-up() :requires [Time-to-impact(2) > 2] :means [Speed-up()] Safe-change-left() :requires [Clear-left()] :means [Change-left()] Safe-change-right() :requires [Clear-right()] :means [Change-right()] Adjust-speed-if-lane-clear() :requires [*not*(Car-ahead-center())] :means [Slow-down-if-too-fast() Speed-up-if-too-slow()] Adjust-speed-if-car-in-front() :requires [Car-ahead-center() *not*(Slow-down-if-too-fast())] :means [Speed-up-if-too-slow() Safe-cruise() Safe-slow-down()] Slow-down-if-too-fast() :requires [Velocity() > Target-speed()] :means [Safe-slow-down()] Speed-up-if-too-slow() :requires [Velocity() <= Target-speed()] :means [Safe-speed-up()] Slam-on-brakes() :action [match-speed-ahead()] </pre>
--	---

enclosing plan. For example, Avoid-trouble-behind might return Change-right, which would become the return value of Drive. Thus, the purpose of an ICARUS program is to find action.

At each successive iteration, ICARUS can return an action from an entirely different portion of Drive. For example, the agent might slam on the brakes on cycle 1, and speed up in service of Get-to-target-speed (a goal-driven plan) on cycle 2. However, if Emergency-brake and Avoid-trouble-ahead do not apply, and the agent is already at its target speed, ICARUS might return the Change-right action in service of Avoid-trouble-behind on cycle 3.

3 A Method for Learning Hierarchical Skills

Now that we have defined the task, we describe our method for learning hierarchical skills from behavioral traces. Our approach involves three distinct stages. The first induces unordered flat rules using a standard supervised learning technique that induces If-Then rules, each of which predicts an action set for a class of situations. To this end, we employ CN2 (Clark & Boswell, 1991) to generate

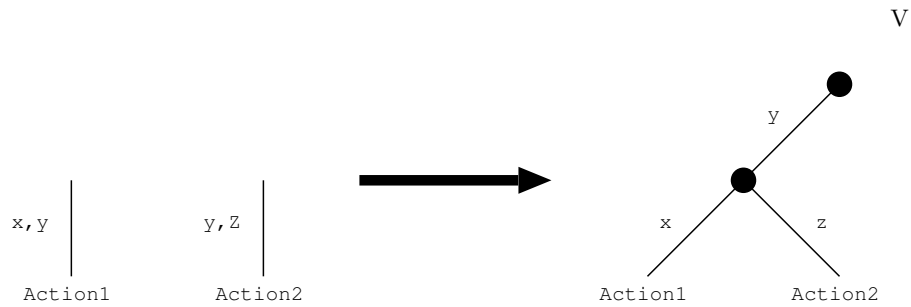


Fig. 1. Operator for promoting conditions.

a set of unordered production rules that determine the target class from attribute values. The second stage creates a classification hierarchy by combining tests that appear in multiple rules. When viewed as an action generator, this structure resembles a hierarchical program. The third stage transforms this representation into an ICARUS program, and simplifies it by taking advantage of ICARUS' semantics. This section discusses the second and third stages.

3.1 Constructing Hierarchies

The second stage of our approach to program induction generates a classification hierarchy. Our method operates by promoting conditions that appear in multiple rules. Consider the two rules:

- If x and y Then *Action1*
- If y and z Then *Action2*

Since the condition y appears in the both rules, we can promote it by creating a more abstract rule that tests the common precondition, using a technique borrowed from work on grammar induction (e.g., Langley & Stromsten, 2000). We illustrate this transformation in Figure 1. Here, the labels on arcs denote conditional tests, and the leaf nodes denote actions. The black circles indicate choice points, where one (or more) of the subsequent tests apply. These structures are interpreted from the top downwards. For example, the right side of Figure 1 classifies the current situation first by testing y , and then, if y holds, by testing x and z (in parallel) to determine which action or actions apply. This results in a more efficient classification process; y is only tested once and, if it does not hold, there is no reason to test x or z . This structure is similar to the decision trees output by C4.5, but more general in that it allows non-exclusive choice.

In addition to promoting conditions, we can promote actions within a classification hierarchy. Figure 2 provides a simple example, where *Action2* occurs at all leaf nodes within a given subtree. If the system is guaranteed to reach at least one of the leaf nodes³ we can associate *Action2* with the root node of the subtree. We represent such nodes with a hollow circle. This simplification applies even if the leaf nodes are at an arbitrary depth beneath the root of the subtree.

³ Here we mean that the tests in the subtree form a collectively exhaustive set.

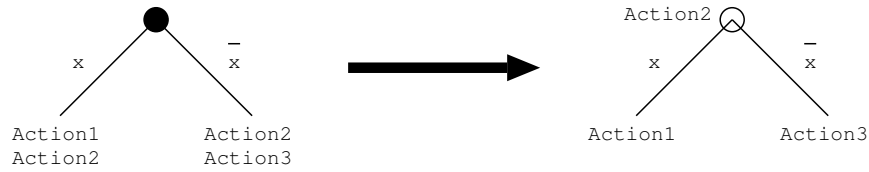


Fig. 2. Operator for promoting actions.

Condition promotion transforms the flat rules learned by CN2 into a classification hierarchy. However, since there are many possible ways to combine rules by promoting conditions, we have an opportunity to shape the final classification hierarchy by defining rule-selection heuristics. (Note that this degree of control would not be available if we had induced decision trees instead of rules.) The key idea is to merge rules with similar actions. In particular, we identify three heuristics that tend to combine rules with similar purposes and isolate rules that represent special cases. The resulting trees transform into understandable programs.

1. Select rules with the same action or same set of actions.
2. Select rules with subset relations among the actions.
3. Select rules with the same conditions.

Our algorithm considers these heuristics in priority order. If two rules select the same action class, they are the highest priority candidates for condition promotion. The operation will only be successful, of course, if the rules share conditions. If more than two rules select the same action class, the ones that share the largest number of conditions will be combined. The second heuristic applies if no two rules select the same action class. In this case, the algorithm looks for rules whose action sets bear a subset relation to one another, such as “Speed-up” and “Speed-up, Change-right”. If a single rule enters into many such pairings, the system takes the ones with the *smallest* number of actions on the theory that these rules express the most cohesive intent. Ties are broken by a similarity metric that maximizes the number of shared and thus promotable conditions. Finally, if no action sets bear subset relations, the system picks rules that share the largest number of conditions. The combination of any two rules yields a subtree with shared conditions on its top-level arc, and these conditions can enter into further promotion operations.⁴ The remaining conditions cannot be merged with any other rules. This process of rule selection and combination continues to exhaustion, merging top-level conditions to build multi-layered subtrees.

A simple example may help to clarify this algorithm. Consider the following three rules (whose abbreviations are defined in Table 2):

⁴ For the purpose of the rule-selection heuristics, the action set of a subtree is the union of the action sets in its leaf nodes, while the most similar subtrees share the largest number of top-level conditions.

IF	$TTIA < 52.18$	IF	$TTIA < 52.18$	IF	$TTIA < 52.18$
AND	$TTIA > 1.82$	AND	$TTIA > 1.82$	AND	$TTIA > 1.82$
AND	$CLR = True$	AND	$CLR = True$	AND	$CLR = False$
AND	$CLL = True$	AND	$CLL = False$	AND	$CLL = False$
THEN	Action =	THEN	Action =	THEN	Action =
	CHR, CHL, CRU, SLO		CHR, CRU, SLO		CRU, SLO

Although no two rules select identical actions sets, all three action sets bear subset relations. In this case, the algorithm will select the last two rules because their action sets are the smallest, and promote three conditions to obtain a new shared structure. Two of those conditions can be combined with the rule for CHR,CHL,CRU,SLO, yielding a three level subtree representing all three rules.

When the process of condition promotion terminates, we add a top-level node to represent the choice among subtrees. Then, we simplify the structure using the action promotion rule shown in Figure 2. This produces the rightmost subtree of the classification structure in Figure 3.

3.2 Constructing the ICARUS Program

We can simplify hierarchical classification structures by translating them into ICARUS and taking advantage of its representational power. The key idea is to recognize that the first phases of program induction always produce a mutually exclusive classification hierarchy, and thus that the branches can be ordered without loss of generality.

Consider the fourth and fifth subtrees of the top node in Figure 3. These represent a rule to avoid collisions, and the responses to a slower car in front (as discussed above). If ICARUS evaluates these in order, it can only reach the fifth branch if the fourth fails to return an action, meaning there is no imminent collision ($TTIA > 1.82$). We can use this knowledge to simplify the logical tests in the fifth subtree, producing the ICARUS subplans labeled R1, R2, R21 and R22 in Table 3. This completes the process of inducing a hierarchical control program from observational traces.

4 An Experiment in Hierarchical Behavior Cloning

Now that we have discussed our method for inducing hierarchical programs, we turn to an experiment that will let us evaluate the approach in a simple driving domain. In specific, we consider the problem of program recovery: we use the ICARUS program of Table 1 to generate trace data, and employ our induction method to recover a second ICARUS program that explains these data. We evaluate the results in terms of the accuracy and efficiency of the recovered program, as well as its conceptual similarity to the source program. We begin by describing the source data, and its transformation into the destination program.

4.1 Data on Driving Behavior

We used the ICARUS program in Table 1 to generate trace data. Since our goal was to recover the structure of a shared driving skill, we needed data from

Table 2. Notation used in example rules and hierarchies.

Actions		Conditions	
Abbreviation	Meaning	Abbreviation	Meaning
CRU	Cruise	CAC	Car Ahead Center
SLO	Slow Down	CBC	Car Behind Center
SPE	Speed Up	CLR	Clear Right
MAT	Match Speed Ahead	CLL	Clear Left
CHR	Change Right	TTIA	Time To Impact Ahead
CHL	Change Left	TTIB	Time To Impact Behind
		VEL	Velocity

multiple drivers whose preferences would collectively span all of the feasible behavior. Instead of creating these agents, we took the simpler approach of directly exercising every control path in the source program, while recording the feature set and the action set available at each time step. This produced a list of situation-action tuples that included every possible action response.

We enumerated five values of in-lane separation (both to the car ahead and behind), five values of velocity for each of the three in-lane cars, and the status of the adjacent lane (whether it was clear or not clear). We chose the particular distance and velocity numbers to produce True and False values for the relevant predicates in the driving program (e.g., time to impact ahead, velocity relative to target speed). This procedure also created multiple occurrences of many situation-action tuples (i.e., the mapping from distance and velocity onto time to impact was many-one). The resulting data had nine attributes. Four of these were Boolean, representing the presence or absence of a car in front/back, and whether the lanes to the right or left of the agent are clear. The rest were numerical attributes. Two of these represented time to impact with the car ahead or behind, two encoded relative velocity ahead or behind, and the last measured the agent’s own velocity.

Our formulation of the driving task assumes six primitive actions. We pre-processed the data to identify sets of these actions that occurred under the same situation. We obtained ten such sets, each containing one to four primitive actions. These sets define a mutually exclusive and collectively exhaustive set of classes for use in program induction.

4.2 Transformation into an ICARUS Program

We employed CN2 to transform the behavioral trace obtained from the ICARUS source program into a set of flat rules, and further transformed that output into a hierarchical classification structure using the condition and action promotion rules of Section 3.1. This produced the tree shown in Figure 3.

We simplified this tree by transforming it into an ICARUS program via a manual process (we expect to automate this in the future). We numbered the branches from left to right and considered them in the order 4,5,3,1,2. This

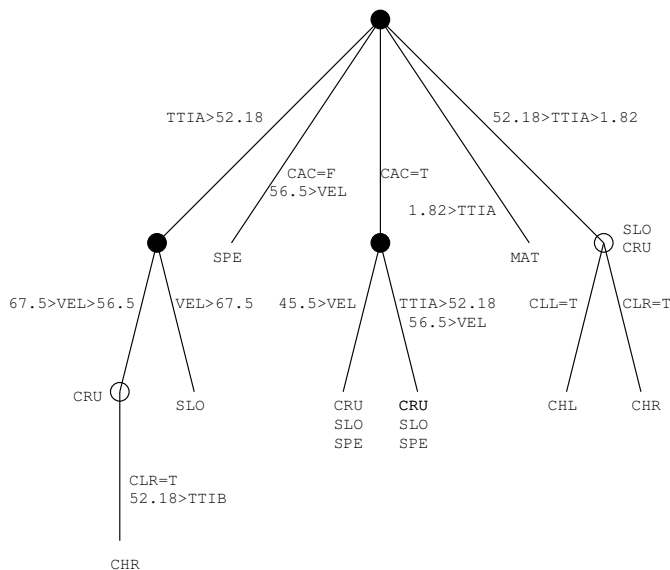


Fig. 3. The classification hierarchy obtained by our method.

ordering simplified the required conditions. Taken as a whole, these transformations recovered the ICARUS program shown in Table 3, completing the task of inducing a hierarchical program from observations.

4.3 Experimental Evaluation

We evaluated our learning method in several ways and at several stages in the transformation process. First, we examined the accuracy of the flat rules induced by CN2 to determine how much of the original behavior we were able to recover. (Since all of the subsequent processing steps preserve information, this measures the accuracy of the recovered program.) In addition, we compared the efficiency of the recovered ICARUS program relative to the flat rules produced by CN2. Finally, we evaluated the structure of the learned ICARUS program in a more subjective sense, by comparing it against the original ICARUS program that generated the data.

We measured the accuracy of the rules induced by CN2 by conducting a 10-fold cross validation. Since the trace data contained circa 4800 situation, action-set tuples, we employed ten training sets, each with about 4300 tuples. For each of these training sets, our method induced a program that had 100% accuracy on the corresponding test set. Moreover, even though the rules induced by the first stage were slightly different across the training runs, the resulting classification hierarchies were identical to the tree in Figure 3. Thus, our heuristics for rule combination regularized the representation.

We also compared the number of conditions that must be evaluated to select action in the induced flat rules and in the recovered ICARUS program. This

Table 3. The ICARUS program induced by our method.

Drive ()	:requires [NOT(R1) NOT(R2) NOT(R3) NOT(R4)]	R3 ()	:requires [VEL < 56.5] :means [SPE R31]
R1 ()	:requires [TTIA < 1.82] :means [MAT]	R31 ()	:requires [CAC = True] :means [SLO CRU]
R2 ()	:requires [TTIA < 52.18] :means [SLO CRU R21 R22]	R4 ()	:requires [NOT(R41)] :means [CRU R42]
R21 ()	:requires [CLL = True] :means [CHL]	R41 ()	:requires [VEL > 67.5] :means [SLO]
R22 ()	:requires [CLR = True] :means [CHR]	R42 ()	:requires [CLR = True TTIB < 52.18] :means [CHR]

provides a measure of the computational efficiency of the two representations. The flat rules required an average of 7361 evaluations to process the training data, while the learned ICARUS program employed 2216. Thus, the hierarchical representation requires only 30% of the effort.

When we compare the learned ICARUS program in Table 3 with the original program in Table 1 several interesting features emerge. First, the learned program is simpler. It employs 10 ICARUS functions, whereas the original program required 14. This was quite surprising, especially since the original code was written by an expert ICARUS programmer. Next, the learned program captures much of the natural structure of the driving task; the top-level routines call roughly the same number of functions, and half of those implement identical reactions. Specifically, R1 in Table 3 corresponds to Emergency-brake in Table 1, while R2 represents Avoid-trouble-ahead using a simpler gating condition. Similarly, R4 captures the behavior of Avoid-trouble-behind, although it adds the Slow-down operation found in Get-to-target-speed. R3 represents the remainder of Get-to-target-speed, absent the Slow-down action. The system repackaged these responses in a slightly more efficient way. The only feature missing from the learned program is the idea that maintaining target speed is an objective. We hope to address this issue in the future, as it raises the interesting problem of inferring the teleological structure of plans from observation.

5 Related Work on Control Learning

We have already mentioned in passing some related work on learning control policies, but the previous research on this topic deserves more detailed discussion. The largest body of work focuses on learning from delayed external rewards. Some methods (e.g., Moriarty et al., 1999) carry out direct search through the

space of policies, whereas others (e.g., Kaelbling et al., 1996) estimate value functions for state-action pairs. Research in both paradigms emphasizes exploration and learning from trial and error, whereas our approach addresses learning from observed behaviors of another agent. However, the nondeterministic policies acquired in this fashion can be used to constrain and speed learning from delayed reward, as we have shown elsewhere (Shapiro et al., 2001).

Another framework learns control policies from observed behaviors, but draws heavily on domain knowledge to interpret these traces. This paradigm includes some, but not all, approaches to explanation-based learning (e.g., Segre, 1987), learning apprentices (e.g., Mitchell et al., 1985), and programming by demonstration (e.g., Cypher, 1993). The method we have reported for learning from observation relies on less background knowledge than these techniques, and also acquires reactive policies, which are not typically addressed by these paradigms.

Our approach is most closely related to a third framework, known as *behavioral cloning*, that also observes another agent’s behavior, transforms traces into supervised training cases, and induces reactive policies. This approach typically casts learned knowledge as decision trees or logical rules (e.g., Sammut, 1996; Urbancic & Bratko, 1994), but other encodings are possible (Anderson et al., 2000; Pomerleau, 1991). In fact, our method’s first stage takes exactly this approach, but the second stage borrows ideas from work on grammar induction (e.g., Langley & Stromsten, 2000) to develop simpler and more structured representations of its learned skills.

6 Concluding Remarks

This paper has shown that it is possible to learn an accurate and well-structured program from a trace of an agent’s behavior. Our approach extends behavioral cloning techniques by inducing simpler control programs with hierarchical structure that makes them correspondingly easy for a person understand. Moreover, our emphasis on learning the shared components of skills holds promise for increased generality of the resulting programs.

Our technique for learning hierarchical structures employed several heuristics that provided a substantial source of inductive power. In particular, the attempt to combine rules for similar action sets tended to group rules by purpose, while the operation of promoting conditions tended to isolate special cases. Both techniques led to simpler control programs and, presumably, more understandable encodings of reactive policies.

We hope to develop these ideas further in future work. For example, we will address the problem of inferring ICARUS objective clauses, which is equivalent to learning teleological structure from observed behavior. We also plan to conduct experiments in other problem domains, starting with traces obtained from simulations and/or human behavior. Finally, we intend to automate the process of transforming classification hierarchies into ICARUS programs. This will let us search for criteria that generate the most aesthetic representation of skills.

References

- Anderson, C., Draper, B., & Peterson, D. (2000). Behavioral cloning of student pilots with modular neural networks. *Proceedings of the Seventeenth International Conference on Machine Learning* (pp. 25-32). Stanford: Morgan Kaufmann.
- Brooks, R. (1986). A robust layered control system for a mobile robot. *IEEE Journal of Robotics and Automation*, 2, 1.
- Clark, P., Boswell, R. (1991). Rule induction with CN2: Some recent improvements. *Proceedings of the European Working Session on Learning : Machine Learning*, LNAI, 482, 151–163.
- Cypher, A. (Ed.). (1993). *Watch what I do: Programming by demonstration*. Cambridge, MA: MIT Press.
- Firby, J. (1989). *Adaptive execution in complex dynamic worlds*. PhD Thesis, Department of Computer Science, Yale University, New Haven, CT.
- Georgeff, M., Lansky, A., & Bessiere, P. (1985). A procedural logic. *Proceedings of the Ninth International Joint Conference on Artificial Intelligence*. Morgan Kaufmann.
- Kaelbling, L. P., Littman, L. M., & Moore, A. W. (1996). Reinforcement learning: A survey. *Journal of Artificial Intelligence Research*, 4, 237–285.
- Langley, P., & Stromsten, S. (2000). Learning context-free grammars with a simplicity bias. *Proceedings of the Eleventh European Conference on Machine Learning* (pp. 220–228). Barcelona: Springer-Verlag.
- Mitchell, T. M., Mahadevan, S., & Steinberg, L. (1985). LEAP: A learning apprentice for VLSI design. *Proceedings of the Ninth International Joint Conference on Artificial Intelligence*, (pp. 573-580). Los Angeles, CA: Morgan Kaufmann.
- Moriarty, D. E., Schultz, A. C., & Grefenstette, J. J. (1999). Evolutionary algorithms for reinforcement learning. *Journal of Artificial Intelligence Research*, 11, 241–276.
- Nilsson, N. (1994). Teleoreactive programs for agent control. *Journal of Artificial Intelligence Research*, 1, 139–158.
- Pomerleau, D. (1991). Rapidly adapting artificial neural networks for autonomous navigation. *Advances in Neural Information Processing Systems 3* (pp. 429–435). San Francisco: Morgan Kaufmann.
- Quinlan, J. R. (1993). *C4.5: Programs for Machine Learning*. Morgan Kaufmann.
- Sammut, C. (1996). Automatic construction of reactive control systems using symbolic machine learning. *Knowledge Engineering Review*, 11, 27–42.
- Schoppers, M. (1987). Universal plans for reactive robots in unpredictable environments. *Proceedings of the Tenth International Joint Conference on Artificial Intelligence* (pp. 1039-1046). Morgan Kaufmann.
- Segre, A. (1987). A learning apprentice system for mechanical assembly. *Proceedings of the Third IEEE Conference on AI for Applications* (pp. 112–117).
- Shapiro, D., Langley, P., & Shachter, R. (2001). Using background knowledge to speed reinforcement learning in physical agents. *Proceedings of the Fifth International Conference on Autonomous Agents* (pp. 254–261). Montreal: ACM Press.
- Shapiro, D. (2001). *Value-driven agents*. PhD thesis, Department of Management Science and Engineering, Stanford University, Stanford, CA.
- Shapiro, D., & Langley, P. (2002). Separating skills from preference: using learning to program by reward. *Proceedings of the Nineteenth International Conference on Machine Learning* (pp. 570–577). Sydney: Morgan Kaufmann.
- Urbancic, T., & Bratko, I. (1994). Reconstructing human skill with machine learning. *Proceedings of the Eleventh European Conference on Artificial Intelligence* (pp. 498–502). Amsterdam: John Wiley.