

Population-Based Learning: A Method for Learning from Examples Under Resource Constraints

Benjamin W. Wah, *Fellow, IEEE*

Invited Paper

Abstract—In this paper, we study a learning model for designing heuristics automatically under resource constraints. We focus on improving performance related heuristic methods (HM's) in knowledge lean application domains. We assume that learning is episodic, that the performance measures of an episode are dependent only on the final state reached in evaluating the corresponding test case and are independent of the intermediate decisions and internal states, and that the aggregate performance measures of the HM's involved are independent of the order of evaluation of test cases. The learning model is based on testing a population of competing HM's for an application problem, and switches from one to another dynamically, depending on the outcome of previous tests. Its goal is to find a good HM within the resource constraints, with proper trade-off between cost and quality. It complements existing point-based machine learning models that maintain only one incumbent HM, and that test the HM extensively before switching to alternative ones. It extends existing work on classifier systems by addressing issues related to delays in feedback, scheduling of tests of HM's under limited resources, anomalies in performance evaluation, and scalability of HM's. Finally, we describe our experience in applying the learning method on five application problems.

Index Terms—Heuristic method, knowledge-lean application, learning from examples, performance evaluation, point-based learning, population-based learning, resource scheduling.

I. INTRODUCTION

MANY application problems in decision, control, and optimization are *sequential problems* that involve making decisions one after another. The *problem solver* for such applications initiates multiple decisions at decision points in order to reach the final state. A *solution* in this context is taken as a sequence of decisions made for a given problem instance in order to reach a desired state. Often, such decisions are carried out by heuristics. *Heuristics*, in general terms, are “rules of thumb” or “common sense knowledge” [41] used in attempting the solution of a problem. Newell *et al.*, defined heuristics as “a process that may solve a given problem, but offers no guarantees of doing so” [38]. Pearl defined heuristics as “strategies using readily accessible though loosely applicable information to control problem-solving processes in human beings and machine” [41]. In this paper, we do

not distinguish between heuristics and strategies; both can be considered as a sequence of *decision rules* applied at decision points that select alternatives [11]. Note that decision rules are heuristics themselves. To avoid confusion between heuristics and elements of heuristics, we use the term *heuristic method* (HM) in this paper to denote a collection of *heuristic decision elements* (HDE's) or *heuristic decision rules* applied to solve a target problem.

HM's often work well but offer no guarantee on performance. In problems whose solutions are evaluated by the cost of getting the solution and the quality of the solution obtained, HM's perform a trade-off between cost and quality. One would expect the quality of the solution to be better if a HM takes more time. The level of trade-off, however, is not precisely defined, it is independent of how well specified the objective of the target problem that the HM solves. (The objective of the target problem may be well defined, or ill defined but measurable.) We define the *performance* of a HM as a collection, but an undefined numeric function, of *performance measures*, each of which is measurable at instants or over an interval. Note that performance measures do not necessarily correspond to cost and quality — elements of cost or quality may be measurable but not cost or quality itself. When performance measures are numeric values, the collection of values of one measure or a subset of measures may be condensed into a more concise form called an *aggregate performance measure* or aggregate performance function. Examples include the sample mean of values of the quality measure, and the sample mean of speedups, each of which is computed as the ratio of two values obtained from two completion time measures.

In applications whose solutions are evaluated by either quality or cost, but not both, there is no trade-off involved. For instance, a heuristic selection HM may be applied in a search algorithm for finding feasible or optimal solutions to a target problem. In this case, the quality of the method is 1 if an optimal or feasible solution is found, and 0 otherwise. The sole objective of the HM is to minimize the cost for finding the solution. However, this objective cannot be formulated mathematically; hence, the “optimal” method is unknown.

The problem of designing good HM's for solving an application problem can, therefore, be considered as a search for such HM's in the space of all possible HM's, without a well-defined objective function of performance measures. Our interest in this paper is on machine learning for the design of such HM's; this learning process can be viewed as an informed search of the space of possible HM's.

Manuscript received March 1, 1992; revised June 15, 1992. This work was supported in part by the National Aeronautics and Space Administration under Grant NCC 2-481 and by the National Science Foundation under Grant MIP 88-10584.

The author is with the Center for Reliable and High Performance Computing, Coordinated Science Laboratory, University of Illinois at Urbana-Champaign, Urbana, IL 61801.

IEEE Log Number 9203536.

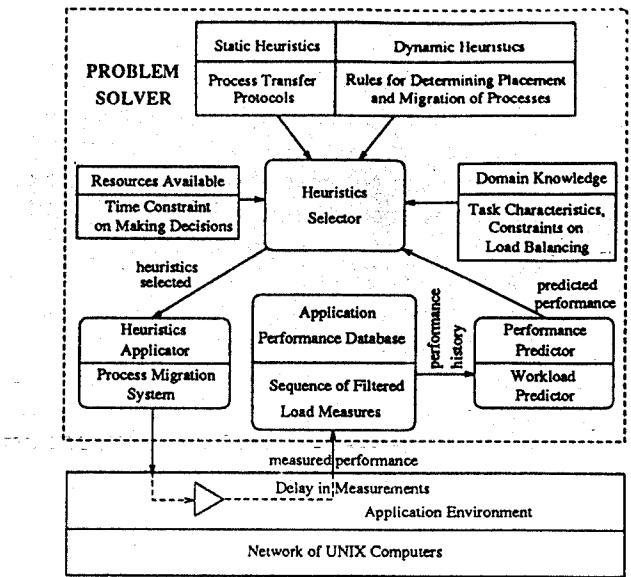


Fig. 1. HM's used in the problem solver as is illustrated by the load balancing problem defined in Table I. (The top caption in each box represents the generic function of the component; the bottom caption is the corresponding function in the load balancing problem.)

TABLE I
THE LOAD BALANCING PROBLEM (LB)

Definition:

Given a process with unknown execution time and a network of UNIX computers with stochastic but nonstationary background workload, determine where the process is to be executed so that the expected speedup with respect to a reference computer is the maximum.

Characteristics:

- The set of measurable parameters is not well defined. Further, there is no model relating the available measurable parameters, such as processor utilization, disk traffic, network traffic, and memory occupancy.
- The load balancing HM has an undefined relationship between cost (in terms of the overhead for collecting measurable parameters and for determining the location of execution) and quality (in terms of the speedup of the process mapped).

Fig. 1 shows the major components of a problem solver and the way they are related. Heuristics are used in the heuristics selector, heuristics applicator, performance predictor, and application performance database. The functions of each component are explained as follows, using the load balancing problem¹ (LB) defined in Table I as an example of the target application.

The *heuristics selector* selects the best heuristics to be applied in the problem solver, based on the domain knowledge available, resources remaining, and past or predicted aggregate performance measures used to rank the heuristics. The heuristics used are static if they depend on static values determined *a priori*, and dynamic otherwise. For instance, the static heuristics in LB are the process transfer protocols, whereas the dynamic heuristics are rules for determining when

¹The objective of load balancing is to find a site with the maximum *relative* speedup with respect to a reference computer for executing any given process. This does not require knowing the actual speedup, as the latter is unknown before the process completes.

migration should be initiated and the placement and migration of processes.

The *heuristics applicator* is a component that executes the heuristics selected by the heuristics selector. In LB, the heuristics applicator is a set of utilities either for preempting a process from execution and for moving it to another location, or for placing a newly arrived process at a given site for execution.

The *application environment* represents the current state of all variables used by the problem solver. This state is affected by the execution of the heuristics by the heuristics applicator, with possible delay in the generation of feedback signals (as is shown by the delay element in Fig. 1). In LB, the application environment is the UNIX operating system that provides either instantaneous performance measures, such as the utilization of the processor and the number of free memory pages at every real-time clock interrupt, or accumulated performance measures, such as the number of disk accesses and packets sent/received between two real-time clock interrupts.

The history of measured performance is captured in the *application performance database* and may be processed in order to reduce its size. If a performance measure is not stationary with time, then a window of this measure may need to be processed in order to facilitate performance prediction in the performance predictor. For instance, in LB, filtering can be applied on a window of CPU occupancies (which measure whether the CPU is occupied or is idle at every real-time clock interrupt) in order to allow the trend of processor utilization to be predicted by the performance predictor.

The last component in the problem solver is the *performance predictor*. This predicts the performance of the heuristics selected by the heuristics selector in the absence of their complete performance behavior. It models the performance of the problem solver, predicting the performance of the various heuristics if they were applied, and supplies these predictions to the heuristics selector. For instance, in LB, a performance predictor available in the UNIX operating system is the load average, which is a moving average of the sampled number of active processes over a window of time (say 2 s).

We present in this paper a machine learning model for learning HM's, based on a population of competing HM's. The major emphasis of population-based learning is on the systematic scheduling of tests of a possibly infinite set of alternative HM's under *resource constraints*. The ideal would be to learn all the heuristics used in the problem solver in Fig. 1 in one learning session; however, the complexity of such may be too high under the resource constraints given, and the learning system may have to focus on one or a small subset at a time. For instance, the learning system may focus to learn and improve the heuristics selector. In this case, the heuristics used there are to be learned, and the rules selected by it, fixed during learning.

This paper is organized as follows. We present in the next section terminologies and assumptions used, and define the scope where population-based learning is suitable. After surveying previous work on point-based learning in Section III, we present an integrated population-based and point-based learning model in Section IV and identify improvements over

the classifier system approach in genetic learning. We discuss in Section V issues related to the integrated learning model, which include the evaluation and normalization of performance measures of HM's, the development of metaheuristics to be used in the learning system, and the statistical scheduling of tests of HM's using limited resources. Section VI describes the application of population-based learning and experiments carried out for enhancing heuristics used in four applications: 1) load balancing of independent processes, 2) automated design of artificial neural networks, 3) depth perception in stereo vision, and 4) combinatorial search using a branch-and-bound algorithm. Throughout the paper, we illustrate our learning method on an application for mapping a set of communicating processes on a network of distributed-memory computers.

II. CHARACTERISTICS OF HEURISTIC METHODS LEARNED

In this section, we characterize HM's to be learned by population-based learning with respect to the domain knowledge available, overhead for evaluating a HM on a test case, and the nature of the HM's. Related issues considered include credit assignment and online versus offline learning.

In general, an application problem is solved by HM's of different types: some are derived from extensive background knowledge, while others are developed with little or no background knowledge. The learning and refinement of the former has been addressed extensively in traditional machine learning studies [32], which focus on representing domain-dependent and domain-independent knowledge, and on reasoning and deductions by using the knowledge stored. An example is illustrated by a learning method proposed by Pearl [40] which discovers heuristics by first consulting simplified models of the problem domain, and by systematic refinement and deletion of constraints from the original application specifications until a semi-decomposable model is found. The solution to the relaxed problem then constitutes a HM for the application problem.

A. Credit Assignment

In this paper, we are interested in the automated design and refinement of HM's for applications with little or no domain knowledge on the HM's concerned. Such *knowledge-lean applications* lack a detailed *world model* that captures the relationship among states, decisions, and *feedback signals* generated by the learning system or measured in the environment. As a result of this lack of domain knowledge, the inference of the world model during learning is either too complex or not fully understood.

The design of a learning system is further complicated when there may be delays in getting the feedback signals due to a decision. In this case, multiple subsequent decisions may have been made between the time one decision was made and its feedback signal received. For knowledge lean applications, the delay in feedback and the lack of a world model mean that there is little or no domain knowledge for *credit assignment*, which involves apportioning a feedback signal to individual decisions carried out in the past, as well as to decision elements

applied in each decision, in order to refine the HM. The former credit assignment is called temporal, and the latter, structural.

Structural credit assignment [53] entails ways of using feedback signals to refine the individual components or rules of a HM. A lack of domain knowledge in this area means that it is difficult to model the behavior of the HM with respect to the feedback signals obtained and to modify the HM in order to arrive at better ones. An example of such a case is when the HM is implemented as a rule-based system, where there is little domain knowledge for apportioning the feedback signal to individual rules and parameters of rules. To overcome this difficulty, the learning system can maintain not only one incumbent HM but a population of competing ones; the choice of the best one to evaluate or new ones to create is delayed until more tests are performed on the HM's generated. In this way, the testing of different HM's may be interleaved in time. Issues and approaches for such a population-based learning method is explored in detail in this paper.

Temporal credit assignment [53] entails the apportioning of temporal global feedback signals by the learning system to past decisions that affect these signals. When a decision is applied, its *temporal scope* is the interval of time during which its direct effect can be observed in the application environment. If the temporal scope is infinite and state changes are Markovian, then the effects due to a feedback signal will be attributed only to the most recent decision made in the past, and the effects of other decisions will be felt indirectly through intervening decisions and states. When the temporal scope is finite and state changes are dependent and non-Markovian, then an approximate temporal model is needed for temporal credit assignment.

In order for population-based learning to be feasible, the temporal scope of decisions must be made in such a way that it is possible to test HM's incrementally and to evaluate their performance when needed. We identify the following necessary characteristics of the temporal scope of decisions that must be possessed by HM's suitable to be learned by population-based learning. 1) Experiments in the target problem can be partitioned into *episodes* or *test cases* so that feedback signals due to a HM applied on a test case occur only when the evaluation of the test case is completed. Decisions in the HM's learned will, therefore, have a finite temporal scope. Note that state changes in the temporal scope may be dependent and non-Markovian. 2) The performance measures of an episode, such as quality and cost, are dependent only on the state and performance values measured at the end of the episode, and are independent of the sequence of intermediate decisions made and internal states generated. 3) The aggregate performance measures of a HM are independent of the order in which test cases are applied. This condition implies that decisions made in one test case do not affect decisions made in other test cases. As a result, computing the aggregate performance measures does not require knowing the effects of decisions for one test case and feedback signals for others. This condition is necessary in population-based learning because HM's may be tested on test cases in any order when multiple HM's are candidates for testing. An example of an aggregate performance measure that satisfies condition 3) is the average

of the performance values of applying a HM on a set of independent test cases.

B. Online Versus Off-Line Learning

Besides the conditions stated previously for credit assignment, learning is feasible only when it is possible to carry out experiments on different HM's with reasonable overhead. Depending on the overhead associated with learning experiments, learning methods can be classified as on-line and off-line.

On-line learning involves evaluating HM's on test cases during learning. It should be used when the HM to be learned is scalable between small and realistic test cases. In this case, the HM's can be tested on small test cases during learning with a relatively small overhead. After learning is completed, scalability is verified on the HM learned using realistic test cases.

Off-line learning involves evaluating the effects of all possible decision sequences of HM's on test cases ahead of time, and looking up the corresponding responses to a decision during learning. It should be used when 1) the HM to be learned is not scalable between small and realistic test cases, and the evaluation of realistic test cases is time consuming; and 2) the performance results of evaluating each realistic test case, using all possible decision sequences and starting from a given initial condition, are manageable and can be generated and stored ahead of time. The second condition means that it is possible to generate *off-line*, for various initial conditions of a test case, the performance results of applying all possible decision sequences on this test case. During learning, when a decision sequence is applied, its effects can be found by looking up the performance results stored, without actually carrying out the decision sequence. If this condition is not true, then the effects due to a decision sequence will have to be found by actually carrying out the decision sequence during learning. This results in a higher overhead in each learning experiment and, within a fixed amount of time, a less number of experiments performed on the HM's.

Offline learning has been applied in various forms in traditional machine learning. In all these cases, certain behavior of the system is established off-line before the heuristics are learned. For instance, in learning the weights of feed-forward artificial neural networks, a set of expected output patterns, each corresponding to a given input pattern, are generated off-line by instrumenting the target application. As another example, in problem solvers with a performance predictor (see Fig. 1), the learning of heuristics is simpler if the performance predictor can be learned off-line. This is the case in the load balancing problem defined in Table I, where the learning of process placement and migration HM's will be easier if the workload predictor is learned off-line.

C. Example

To illustrate the concepts we have defined so far, consider a target application of mapping a set of communicating processes on a network of distributed memory computers (see Table II). We assume that the HM is implemented as a set of rules (or HDE's) in the post-game analysis system [64], which

TABLE II
THE PROCESS MAPPING PROBLEM (PM) AND THE LEARNING OF ITS HM'S

Definition:

Assuming a set of communicating processes with known processing and communication requirements, and a network of idle distributed memory computers, the problem is to map the processes to the computers so that the completion time of executing the processes is the minimum. The target problem is NP-hard and has well-defined objective and constraints.

Characteristics:

- Measurable parameters are well defined and include the execution times of processes in between communications with other processes, the amount of data in a message, the network speed, and the average communication overhead.
- Relationship between cost (execution time) and quality (completion time of the processes mapped) of the heuristic mapping methods is undefined.

Credit assignment:

- Structural*: Difficult, as the HM is rule based [64] and the domain is knowledge lean.

- Temporal*: Simple, as the conditions prescribed for temporal credit assignment in population-based learning are satisfied.

Online learning:

HM's are tested on small and independent test cases during learning; scalability is verified after learning is completed.

is used to find efficient mappings of processes for networks of computers. The performance measures of the HM's involve cost and quality, and their trade-off is undefined; hence, the objective of designing good HM's is undefined as well. Table II also lists issues in credit assignment. Structural credit assignment is difficult for such heuristic mapping methods implemented as a rule-based system, as there is inadequate domain knowledge for relating performance feedback signals to improvements for HM's. Temporal credit assignment, on the other hand, satisfies the conditions prescribed for population-based learning. 1) Learning can be carried out using small and independent test cases with finite temporal scopes. 2) The performance measures of a HM applied on a test case (a process graph) depend only on the completion time of mapping the process graph on a given network of computers and the overhead for generating such a mapping. These measures are independent of the individual decisions and internal states in the decision sequence. 3) The aggregate performance measures (average completion time and average overhead for mapping a set of process graphs on the network of computers) are independent of the order of evaluating the test cases. Furthermore, HM's learned using small test cases are scalable (to some extent) to larger ones. Given these conditions, on-line population-based learning can be used.

D. Other Issues

In addition to the issues on feasibility of learning, HM's can be differentiated according to how difficult it is to learn them by population-based methods. We characterize this difficulty with respect to the nature of HM's to be learned and their possible margin for improvement.

First, population-based learning is suitable for learning performance related heuristics but not for learning correctness related ones. A HM is said to be *performance related* if the constraints of the target problem are trivially satisfied, and the goal of the HM is to find efficient solutions. In contrast, a HM is *correctness related* if the constraints of the problem are hard to satisfy, and the goal is to find a feasible as well as efficient solution that satisfies the constraints. In knowledge lean applications, HM's that are performance related are easier to learn than ones that are correctness related.

In learning correctness-related HM's, domain knowledge is needed for guiding the generation and verification of heuristics. Examples of such domain knowledge include the set of correct operators to apply in the problem solver, verification that the objective or optimality condition of the target problem is achieved, syntactic and semantic knowledge used to perturb existing HM's or to compose new ones, and in some cases, a way to prove the correctness of the HM learned. Given this domain knowledge, correctness related HM's can be learned by examples: the learning system learns these methods by trying to prove them through positive examples and disprove them through negative examples. This learning approach was taken in AQ [31], LEX2 [35], SAGE.2 [25], and TEACHER 1.0 [65]. The key in these learning systems is on techniques for transforming a HM into another when positive or negative examples are found. Note that the learning of correctness related HM's is not improved even when a population of competing HM's are maintained, as the performance measures available (in terms of the number of positive examples found for a HM) cannot be used to identify a promising HM to test next.

Second, population-based learning methods (and other learning methods as well) can be used for improving performance related HM's if there is a suitable margin for improving the aggregate performance measures targeted for enhancement. This margin is related to the quality and cost for finding solutions by existing HM's. In general, it will be difficult to improve an existing HM if the solution it finds has a quality very close to the maximum and a cost very close to the minimum. In this case, it will be easier to learn a HM with one or more of the performance requirements relaxed. For instance, if the existing HM has an average quality very close to the maximum, then it will be difficult to learn a new HM with a better average quality, but it may be easier to learn one with a lower average quality and lower average cost.

E. Summary

Table III summarizes the scope of HM's that can be learned by a population-based learning model. Problems that are not covered are either too difficult to be learned automatically or can be learned by existing machine learning methods. We are interested in finding better HM's for solving a set of instances of the target application, but not in finding a good HM for solving one particular problem instance, as this is less common in practice. Note that population-based learning addresses ways for coping with a lack of domain knowledge for structural credit assignment. It does not address ways

TABLE III
CHARACTERISTICS OF HM'S SUITABLE FOR POPULATION-BASED LEARNING

Application Environment:

- Knowledge lean, without a good world model.

Type of HM's:

- Performance related.

Structural credit assignment:

- Domain knowledge for structural credit assignment is missing.

Temporal credit assignment:

- Learning is episodic with finite temporal scopes. If state changes within the temporal scopes are non-Markovian, then an approximate temporal model will be needed for temporal credit assignment.
- Performance measures for evaluating a test case (episode) are dependent only on the final state reached in the episode and are independent of the intermediate decisions and internal states.
- Aggregate performance measures of a HM are independent of the order of evaluating the test cases.

On-line learning:

- Aggregate performance measures of HM's learned are scalable between small and realistic test cases. The overhead is low for evaluating a HM on a small test case.

Off-line learning:

- Aggregate performance measures of HM's learned are *not* scalable between small and realistic test cases; the overhead is large for learning with realistic test cases.
- The performance results of evaluating a realistic test case, using all possible decision sequences and starting from a given initial condition, are manageable and can be generated ahead of time.

for temporal credit assignment, which must be solved either implicitly due to the nature of the problem, or explicitly by assuming an approximate temporal model.

Examples of HM's that can be learned include thresholds of functions, numerical coefficients that are parts of a rule or function, numerical functions obtained by composition of other numerical functions, symbolic rules, and macro-operators formed by composition of basic operators. The simplest case, therefore, includes parameter tuning, while more complicated cases involve generating symbolic rules and relations as well as tuning parameters of rules in a rule-based system.

It is important to point out that if a HM has both performance- and correctness-related components, it is essential that the performance-related components, which are amenable to population-based learning, be explicitly indicated by the designers. Furthermore, the new HM's learned must be semantically sensible. This may either be enforced by explicit specifications used in the generation process or be verified in a post-learning phase. Examples of such specifications include the range of possible values of a variable and the way that HDE's can be modified and combined.

III. POINT-BASED LEARNING MODELS

In this section, we survey three models developed in traditional machine learning studies for learning HM's, and discuss why they cannot be used to learn HM's with the characteristics summarized in Table III. These existing models

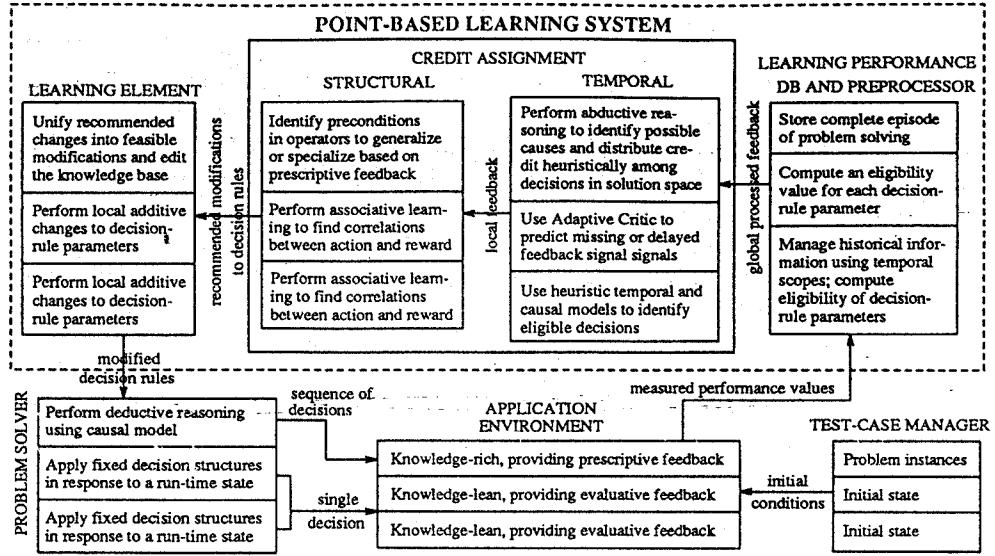


Fig. 2. Point-based learning models. (In each block, the box on the top represents actions carried out in Dietterich and Buchanan's model; the box in the middle, Minsky's model; and the box in the bottom, the hybrid model.)

examine one incumbent HM at a time; hence, they are termed *point-based learning models* in this paper. Fundamental work in this area was addressed by Minsky [33], Mitchell [34], [48], and Dietterich and Buchanan [7]. The basic principle in various learning models is based on a generate-and-test paradigm that generates plausible HM's, performs limited tests, and modifies the HM's according to the feedback signals obtained. Note that in order for a generate-and-test method to work, it is essential that experiments with different test cases be repeatable, allowing different HM's to be tested under identical experimental conditions. The three models are described briefly as follows and are summarized in Fig. 2.

A. Dietterich and Buchanan's Model

Dietterich and Buchanan's model [1], [7], and similar models proposed by Smith *et al.* [50] and Langley [25], belong to a class of models for learning HM's of target problems with well-defined objectives. They learn by supervised learning in a knowledge rich environment with *prescriptive feedback*, which carries explicit information about the desired operators and/or states; the learning system can use the feedback signals for computing an error in order to guide the modification of the HM tested. The HM's learned are generally in the form of static deterministic decision rules represented as preconditions of either operators or operator sequences.

The learning model, shown in Fig. 2, includes the learning performance database and preprocessor, credit assignment unit, and learning element. The problem solver, also called the performance element, and its initial conditions are shown as components outside the learning system. Actions carried out in Dietterich and Buchanan's model are indicated by boxes on the top in each block in Fig. 2. The *problem solver* finds solutions to problem instances by performing an explicit search for the goal state using deterministic decision rules. By reasoning using the *causal model*, which relates preconditions of decision rules to post-conditions, the problem solver generates a

sequence of decisions. These decisions affect the environment, whose effect is captured by the *learning performance database and preprocessor* as a complete solution sequence (or state changes). The preprocessed data are then used for credit assignment, which relates the external feedback to objectives and to determine the preconditions of operators responsible for the current feedback. Temporal credit assignment, based on *abductive reasoning*, identifies possible causes to certain feedback signals from the causal model, and distributes credit heuristically among candidate decisions in the solution sequence. Structural credit assignment, on the other hand, is performed by identifying preconditions in decision rules to generalize or specialize. The result of credit assignment is a set of recommended modifications to the decision rules, which are passed to the learning element. The *learning element* performs syntactic operations for parsing and for modifying HM's in symbolic representations. By using an internal model of the problem solver and knowledge for reasoning, the learning element generalizes stored solutions by deriving preconditions and macro-operator sequences, based on successfully solving a few episodes of the target problem. Knowledge in the learning element can be specified procedurally either in the form of a learning algorithm, such as that used in LEX [36], or declaratively in the form of editing rules, such as those in SAGE.2 [25], ODYSSEUS [61], and Cupr [23].

B. Minsky's Model

Minsky's Model [33] is an older and perhaps less restricted model of learning systems. It applies well for learning HM's of target problems with undefined objectives in knowledge lean environments, which tend to produce evaluative and possibly delayed feedback signals.

Evaluative feedback carries only implicit information about the desired behavior but explicit evaluation of the observed behavior. Such behavior is intrinsically *a posteriori*, being measured or generated after the behavior has occurred. It

requires a critic [60] who has some prior knowledge of the objective function and can assess the goodness of external states or sequences thereof. Scalar evaluative feedback signals are called *reinforcements* [33], and learning from such signals, *reinforcement learning*.

Evaluative feedback is often generated by dynamic and stochastic HM's. These HM's operate continuously using time-varying data from the application environment, make decisions based on probabilistic distributions, and employ statistical evaluation methods. A notable feature of the episodes is that they are hard to define. As a result, it is hard to retain and store a complete solution or a trace of execution of a HM, thereby requiring the eligibility values of decision making parameters to be computed on-line. Moreover, since a smaller amount of information is extracted in each cycle as compared to Dietterich and Buchanan's model, learning tends to take longer.

The model is again shown in Fig. 2; actions carried out are shown as boxes in the middle in each block. The *problem solver* applies static stochastic decision rules to generate a decision, whose effect on the environment may be delayed. The *learning performance database and preprocessor* computes incrementally the eligibility of decision rule parameters using the current feedback signal. Credit assignment is then performed on the global feedback signal. Temporal credit assignment is the major problem to be addressed due to the evaluative feedback and the time-varying nature of the data measured in the environment. It is carried out by the *adaptive critic* (also called the feedback predictor or the secondary reinforcement device), which learns on-line to produce rewards (or reinforcements) in the presence of delayed and evaluative feedback. It does not attempt to project into future states, as in Dietterich and Buchanan's model, due to the lack of a causal model. In existing implementations, the feedback signals and states are assumed to follow a Markovian model, which simplifies the maintenance of past states and temporal credit assignment. Structural credit assignment is performed next using associative learning that correlates the decision and the corresponding feedback signal. The result is a set of eligibility values of decision making parameters in the HM to be learned, which are then applied in the learning element. Examples of systems in this class include Klopf's drive reinforcement model [22] and Sutton and Barto's reinforcement model [54].

C. Hybrid Learning Model

There are problems that cannot be solved by either Dietterich and Buchanan's or Minsky's model alone. A problem in this class is characterized by a knowledge lean learning environment with an ill-defined objective and evaluative feedback, and a non-Markovian temporal scope that can be limited and defined heuristically by a temporal model. Minsky's model is not applicable in this case since it lacks the reasoning mechanism for distributing credit among explicitly stored past decisions. Dietterich and Buchanan's model is not applicable because it requires a well-defined objective with prescriptive feedback. These limitations lead to the development of a

The actions carried out in the hybrid model are shown as boxes in the bottom in each block in Fig. 2. Since the causal model is usually heuristic and less detailed than that in Dietterich and Buchanan's model, it does not provide post-conditions for finding the result of a decision sequence; rather, it is used mainly for identifying candidate decisions during temporal credit assignment. Decisions, in this case, are made one at a time, and temporal credit assignment is performed heuristically on a stored solution sequence by computing an eligibility value for each individual decision in the history. The hybrid model is characterized by Samuel's Checker Player [44], [45], EURISKO [26], Williams' REINFORCEMENT model [62], LS-1 [51], classifier system [16], and the truck-backer-upper problem of Widrow *et al.* [39].

D. Summary

To summarize, we present three point-based learning models in this section. For knowledge-rich environments with prescriptive feedback and a causal model, Dietterich and Buchanan's model can be used. This learning model is not applicable for knowledge lean environments, which are usually associated with evaluative (and sometimes delayed) feedback and problems with ill-defined objectives. For knowledge lean environments with Markovian state changes, Minsky's model can be applied. The Markovian property simplifies temporal credit assignment, allowing eligibility values that correspond directly to the decision making parameters to be computed online. For knowledge lean environments with non-Markovian state changes but with an approximate temporal model, the hybrid learning model can be used. This learning model works well when the heuristic temporal model is reasonably accurate, and when the time-varying parameters have stationary distributions.

For learning HM's with properties shown in Table III, only the hybrid point-based model can be applied if there was an approximate temporal model relating decisions to feedback signals. However, additional learning techniques must be developed for coping with the lack of domain knowledge for structural credit assignment. In the next section, we present a population-based learning model that generates new HM's without relying on domain knowledge for structural credit assignment, and that decides the HM to be evaluated depending on test results obtained during learning.

IV. INTEGRATED POPULATION-BASED AND POINT-BASED LEARNING MODEL

In this section, we present an integrated learning model that learns new HM's, starting from a pool of competing HM's and generating new ones when needed. This model is developed for learning HM's with properties described in Table III, whose problem solvers lack an internal model of the decision process for structural credit assignment. Such is the case when the HM's are represented as decision procedures.

Population-based learning generates new HM's using mostly domain independent knowledge and evaluative feedback signals: a new HM generated this way may have very different features from the original one. The learning system schedules

tests of HM's dynamically, without committing to a particular incumbent HM for an extended period of time. Structural credit assignment is, therefore, avoided entirely, and intermediate performance results may be used to alter the probability of selecting a HM for evaluation, but not necessarily to improve the incumbent HM. On the other hand, point-based learning can be applied when a particular HM is selected for evaluation. Using an approximate temporal model, a hybrid point-based learning method can explore some neighborhood of the initial HM, converging to the nearest locally best HM.

One of the major issues of population-based learning is the scheduling of a fixed amount of resources for testing HM's in order to find a good HM when the resources are expended. Resource scheduling is essential because the number of possible HM's for solving a target problem is very large and possibly infinite, and there are always resource constraints in a practical implementation. Under resource constraints, learning should not be focused on one incumbent HM at a time, as there is inadequate knowledge for structural credit assignment, and the penalty for following only one HM in the absence of such knowledge is large.

A. Components of the Learning Model

Fig. 3 shows the major components of the learning model. We describe briefly in this section the function of each component, leaving the discussion on their design issues to Section V. In this discussion, we assume that a HM is represented as a set of fixed decision rules and, in general, a set of fixed structures of HDE's.

In the integrated learning model, the time resource is divided into fixed-size pieces called *quanta*. At the beginning of a quantum, the *resource scheduler* selects either to test an existing HM in the population or to generate a new one. The decision on whether to generate a new HM depends on the resources remaining and the partial performance results of the HM's evaluated. If the decision is to generate a new HM, the *heuristics manager* will be invoked. Since the knowledge used for structural credit assignment is weak, the operators for composing new HM's are usually model free and syntactic in nature, and may perturb existing parameters and rules to obtain new ones. Examples of such operators have been studied in classifier systems and include mutation and crossover [6]. New HM's generated this way may be quite different from existing ones.

Once a HM is selected for testing, the *test case manager* is called to supply a set of test cases. These test cases may either be generated using domain knowledge of the application or be selected from an existing pool.

By using the set of test cases supplied by the test-case manager and the HM selected by the resource scheduler, the system learns by using the point-based hybrid learning model described in Section III. In the quantum of time allowed, the *problem solver* evaluates the test cases one at a time, each involving the generation of a decision sequence that searches for the goal state. The performance results measured in the *application environment* after each test case is evaluated are used for assigning credit to decisions carried out in the

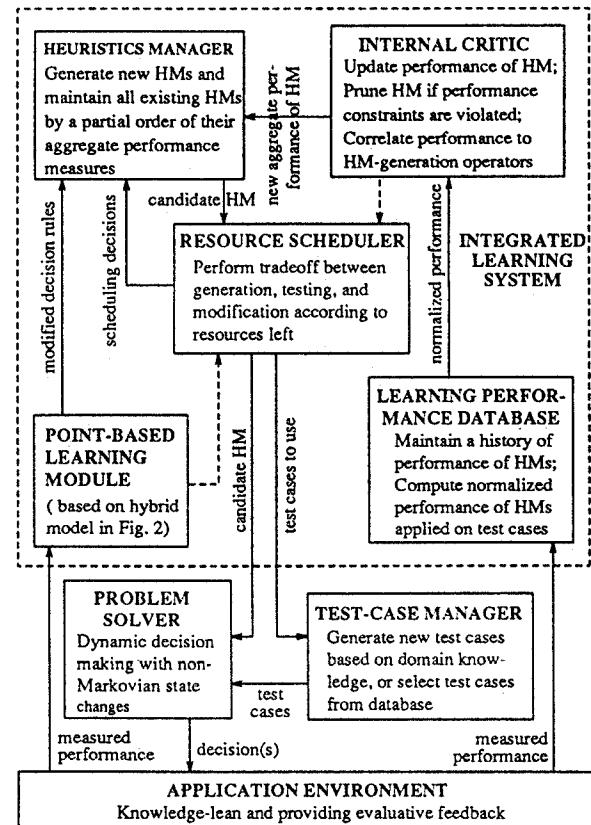


Fig. 3. Integrated population-based and point-based learning model. (Dashed arrows indicate completion signals sent to the resource scheduler.)

decision sequence (temporal credit assignment). The algorithm for temporal credit assignment is similar to the ones used in genetic algorithms [13], [17], [63], and generally involves modification of the probability that a HDE will be invoked in future tests. Structural credit assignment for modifying the elements or the parameters of a HDE is difficult, as there is insufficient domain knowledge.

Testing of the HM selected will be terminated by the resource scheduler when either the intermediate performance results are not satisfactory or the time quantum is expended. In either case, the performance results obtained on the test cases evaluated in this quantum are normalized and processed by the internal critic. By using associative learning methods, the internal critic correlates the performance of this HM to the operators used in generating this HM. It then passes the updated aggregate measures, together with possible changes to the HM concerned, to the heuristics manager. The cycle of learning is then repeated until the time allowed is exceeded.

In summary, there are two possible sequences of learning. Referring to Fig. 3, the first sequence is based on a point-based approach that calls the resource scheduler, problem solver, and the hybrid point-based learning method shown in Fig. 2. In this case, the problem solver evaluates a particular HM on a selected test case. Using evaluative feedback signals and a weak temporal credit assignment algorithm, the point-based learning system makes minor adjustments, and perhaps minor improvements, to the HDE's. The key issue here is on the design of the temporal credit assignment algorithm for assigning credit to each decision carried out in the decision

sequence. On the other hand, the second sequence is based on a population-based approach that traverses the resource scheduler, problem solver, learning performance database, internal critic, and heuristics manager. Without relying on domain knowledge for structural credit assignment, it makes major adjustments to the HM's tested, and interleaves testing among different HM's. The key issue here is the scheduling of resources for testing various HM's so that a good HM can be found when resources are expended.

The learning model shown in Fig. 3 has different emphasis, depending on whether learning is on-line or off-line. Recall from Section II that online learning is used when the HM to be learned is scalable between small and realistic test cases, and that off-line learning is used when it is possible to generate ahead of time the performance results of evaluating each realistic test case, using all possible decision sequences and starting from a given initial condition. If online learning is used, then test cases will actually be evaluated by invoking the problem solver and test case manager during learning, and may represent a significant overhead during learning. As a result, it is important to implement the problem solver and test case manager efficiently. Critical issues to be addressed include the study of good scheduling algorithms for allocating time and processor resources, and the design of an efficient heuristics manager and internal critic. In contrast, if off-line learning is used, then the problem solver and test-case manager will not be invoked during learning; consequently, resource scheduling is a secondary issue because the overhead is relatively small for looking up the performance results stored for a decision sequence. The critical issues to be studied are, therefore, the design of an efficient heuristics manager and internal critic.

B. Relationship to Classifier Systems

Our integrated learning model shares three basic features with the classifier system approach in genetic learning [6]. First, a HM in our model corresponds to a policy (which can be represented as a rule-based system) in a classifier system; similarly, a HDE in our model corresponds to a rule in a policy. Second, we assume that our application domain is knowledge lean, as is the case in classifier systems. As a result, we employ weak and domain independent methods for temporal credit assignment in the point-based learning component of our model. These methods are similar to the ones used in classifier systems, such as the bucket brigade algorithm [17], Wilson's reinforcement algorithm [63], and the profit sharing plan [13]. Third, the heuristics manager in our model employs domain independent operators for generating new HDE's and HM's; the counterparts are the mutation and crossover operators used in classifier systems for generating new rules and policies.

We have, however, studied three issues that were not addressed adequately in previous studies of classifier systems. First, we identify the temporal properties of HM's that must be satisfied in order for population-based learning to be possible. These conditions are summarized in Table III and apply to cases with or without delay between the initiation of a decision and the generation of its corresponding feedback. In contrast, policies in classifier systems are assumed to have no delay

between a decision and its feedback. We survey previous learning methods in Section III and describe why they cannot be applied for learning HM's for application problems studied in this paper. We also differentiate between applications that should be learned on-line and ones that should be learned off-line. Learning experiments demonstrating the above concepts are shown in Section VI.

Second, we study metrics for evaluating the performance of HM's. Specifically, we show in Section V-A that there may be anomalies in evaluation when multiple performance measures are combined into a single aggregate performance measure. We also propose new methods for resolving these anomalies during (Section V-B) and after (Section V-A) learning.

Third, we present in Section V-C, algorithms for scheduling learning experiments, based on dynamic performance data obtained in previous experiments. Resource scheduling is important when learning has to be carried out within a limited time. In contrast, resource scheduling algorithms developed for classifier systems are static in nature, employing schedules that predefine the duration of a generation, the number of policies tested in a generation, and the number for tests performed on a policy. The static scheduling approach is demonstrated by the work of Grefenstette *et al.*, for learning sequential decision rules [14].

V. ISSUES ON THE INTEGRATED LEARNING MODEL

In this section, we discuss issues on designing a system which implements the integrated learning model presented in the last section. The issues considered include performance evaluation of HM's, acquisition and development of meta-heuristics used in the learning system, and resource scheduling. Meta-heuristics used are further classified into the ones for the internal critic, heuristics manager, and test case manager.

A. Performance Evaluation of HM's

In this section, we study issues related to anomalies in performance evaluation of HM's and methods for coping with these anomalies.

1) Anomalies in Performance Evaluation: HM's are compared on the basis of their performance, which is related to the objective of the target problem to be solved. As is discussed in Section I, the objective of a HM, as well as its trade-off, are usually unknown with respect to its performance measures. A possible solution to cope with this issue is to derive the objective function as an aggregate performance function with tunable parameters, and to find a particular combination of values of these parameters that lead to the HM with the best trade-off. We applied this approach in our earlier work, using objectives based on parametrized functions of ratio of average quality to average cost [18] and ratio of normalized quality to normalized cost [19]. This approach fails for three reasons.

First, it is very difficult to find a particular combination of values of parameters in the aggregate performance function so that when applied in the integrated learning system, the HM with the best trade-off between cost and quality can be found.

Second, it is difficult to compare the performance of two HM's when they are evaluated by test cases of different sizes

TABLE IV
EXAMPLE TO ILLUSTRATE ANOMALIES IN PERFORMANCE EVALUATION OF HM's.
HM 1 IS BETTER THAN HM 2 WITH RESPECT TO AGGREGATE PERFORMANCE MEASURES \bar{Q}/\bar{C} AND $\bar{\pi}$; HM 2 IS BETTER THAN HM 1 WITH RESPECT TO \bar{P} AND $\bar{\sigma}$; THE AGGREGATE MEASURE \bar{S} SHOWS THAT BOTH HM'S HAVE SPEEDUP GREATER THAN 1 WITH RESPECT TO THE OTHER

HM <i>i</i>	Performance Measure	Test Case <i>j</i>					Average Measure Over All <i>j</i> 's
		1	2	3	4	5	
1	$Q_{1,j}$	10	20	15	10	30	$\bar{Q}_1 = 17$
	$C_{1,j}$	10	20	20	10	50	$\bar{C}_1 = 22$
	$P_{1,j} = \frac{Q_{1,j}}{C_{1,j}}$	1	1	0.75	1	0.6	$\bar{P}_1 = 0.87$
	$S_{1,j} = \frac{C_{2,j}}{C_{1,j}}$	2	1.4	0.5	1	0.2	$\bar{S}_1 = 1.02$
	$\pi_{1,j} = \frac{P_{1,j}}{P_{1,j} + P_{2,j}}$	0.67	0.70	0.43	0.5	0.29	$\bar{\pi}_1 = 0.52$
	$\sigma_{1,j} = \frac{S_{1,j}}{S_{1,j} + S_{2,j}}$	0.8	0.66	0.2	0.5	0.04	$\bar{\sigma}_1 = 0.44$
2	$Q_{2,j}$	10	12	10	10	15	$\bar{Q}_2 = 11.4$
	$C_{2,j}$	20	28	10	10	10	$\bar{C}_2 = 15.6$
	$P_{2,j} = \frac{Q_{2,j}}{C_{2,j}}$	0.5	0.43	1	1	1.5	$\bar{P}_2 = 0.89$
	$S_{2,j} = \frac{C_{1,j}}{C_{2,j}}$	0.5	0.71	2	1	5	$\bar{S}_2 = 1.84$
	$\pi_2 = \frac{P_{2,j}}{P_{1,j} + P_{2,j}}$	0.33	0.30	0.57	0.5	0.71	$\bar{\pi}_2 = 0.48$
	$\sigma_2 = \frac{S_{2,j}}{S_{1,j} + S_{2,j}}$	0.2	0.34	0.8	0.5	0.96	$\bar{\sigma}_2 = 0.56$

or parameters. In this case, the evaluation results may lead to a wrong conclusion that a HM evaluated on test cases of smaller size is better than other HM's evaluated on larger test cases.

Third, there are anomalies in the evaluation of HM's when two HM's are compared by using either different aggregate performance functions, or the same aggregate performance function with different parameters. Depending on the aggregate measures and the reference HM used, one HM can be shown to be better or worse than another.

To illustrate the above point, consider the hypothetical performance data of two HM's shown in Table IV. Let $Q_{i,j}$ and $C_{i,j}$ be, respectively, the quality and the cost of the *j*th case evaluated by HM *i*. The following observations are found from the data in Table IV.

1) $\bar{Q}_1/\bar{C}_1 > \bar{Q}_2/\bar{C}_2$ but $\bar{P}_2 > \bar{P}_1$: this shows that the evaluation based on ratios of averages can be inconsistent with the evaluation based on averages of ratios.

2) $\bar{S}_1 > 1$ and $\bar{S}_2 > 1$: this means that both HM's have speedups greater than 1 with respect to the other. Note that in computing the speedup of HM 1, HM 2 is used as the reference, whereas in computing the speedup of HM 2, HM 1 is used as the reference. What should have been done is

that only one HM is chosen as the reference and be used consistently throughout learning and evaluation.

3) $\bar{\pi}_1 > \bar{\pi}_2$ but $\bar{P}_2 > \bar{P}_1$: to overcome the problem that values of different test cases are weighted differently in computing the average, one possible approach is to normalize the performance values before taking the average. (The normalized performance value of a test case for HM 1 or 2 is found by dividing its performance value by the sum of the performance values of HM's 1 and 2.) In Table IV, $\bar{\pi}_i$ represents the average of normalized performance values for HM *i*. What we found is that the evaluation based on averages of normalized performance values can be inconsistent with the evaluation based on averages of raw performance values.

4) $\bar{\pi}_1 > \bar{\pi}_2$ but $\bar{\sigma}_2 > \bar{\sigma}_1$: this shows that the evaluation based on averages of normalized performance values can be inconsistent with the evaluation based on average normalized speedups.

2) *Methods for Resolving Anomalies*: Although the previous example is shown with hypothetical data, it points out that there may be anomalies in using aggregate functions for performance evaluation. To resolve these anomalies, we propose an approach based on consistent normalization and statistical analysis of the performance data. Our approach has the following steps.

1) Identify a *reference* or baseline HM upon which all other HM's are compared. This reference HM can be chosen to be the best existing HM for solving the given application problem.

2) Normalize each raw performance measure of a new HM *consistently* with respect to the same measure of the reference HM (evaluated on either the same set of test cases or test cases with the same distribution of performance) so that it is meaningful to compare two HM's based on their normalized measures. There are two approaches to normalization. First, all performance values of a measure of the new HM can be normalized with respect to a representative value of the same measure of the reference HM. Note that the representative value used must be invariant with time. To illustrate this approach, let the average quality and cost of the reference HM be \bar{Q}_0 and \bar{C}_0 , respectively. The normalized values of quality and cost measures are shown in the top half of Table V. The normalization constants used, namely, \bar{Q}_0 and \bar{C}_0 , are invariant when more tests are performed on the reference HM. Note that if \bar{Q}_0 and \bar{C}_0 cannot be obtained directly, then unbiased estimates of their population means will be used instead. Second, a raw performance value of a measure of the new HM can be normalized with respect to the corresponding raw performance value of the same measure of the reference HM. This is illustrated in the lower half of Table V. In this case, the ratio-of-quality (ROQ) and ratio-of-cost (ROC) are taken as the new performance measures. Note that the average ratio-of-quality (respectively, ratio-of-cost) measure for the reference HM is 1.

3) Compare two HM's based on individual normalized performance measures. If an aggregate performance measure is used, then it is essential that the aggregate measure be computed from test cases of a single normalized measure but not of multiple normalized measures. For instance, if cost and quality values are available for two HM's, then the HM's can be

TABLE V
TWO CONSISTENT WAYS FOR NORMALIZING QUALITY AND COST

Normalization for Quality and Cost Measures							Average
Reference HM	Quality	$Q_{0,1}/\bar{Q}_0$	$Q_{0,2}/\bar{Q}_0$...	$Q_{0,n}/\bar{Q}_0$	1	
	Cost	$C_{0,1}/\bar{C}_0$	$C_{0,2}/\bar{C}_0$...	$C_{0,n}/\bar{C}_0$	1	
New HM	Quality	$Q_{1,1}/\bar{Q}_0$	$Q_{1,2}/\bar{Q}_0$...	$Q_{1,n}/\bar{Q}_0$	\bar{Q}_1/\bar{Q}_0	
	Cost	$C_{1,1}/\bar{C}_0$	$C_{1,2}/\bar{C}_0$...	$C_{1,n}/\bar{C}_0$	\bar{C}_1/\bar{C}_0	
Normalization for ROQ and ROC Measures							Average
Reference HM	ROQ	$Q_{0,1}/Q_{0,1}$	$Q_{0,2}/Q_{0,2}$...	$Q_{0,n}/Q_{0,n}$	1	
	ROC	$C_{0,1}/C_{0,1}$	$C_{0,2}/C_{0,2}$...	$C_{0,n}/C_{0,n}$	1	
New HM	ROQ	$Q_{1,1}/Q_{0,1}$	$Q_{1,2}/Q_{0,2}$...	$Q_{1,n}/Q_{0,n}$	$1/n \sum_j Q_{1,j}/Q_{0,j}$	
	ROC	$C_{1,1}/C_{0,1}$	$C_{1,2}/C_{0,2}$...	$C_{1,n}/C_{0,n}$	$1/n \sum_j C_{1,j}/C_{0,j}$	

compared with respect to cost and quality measures separately, but not with respect to aggregate measures such as cost-quality ratios. This step is necessary in order to avoid anomalies when performance is evaluated using parametrized aggregate functions of multiple performance measures. It implies that a multidimensional representation of performance measures is needed when the number of measures is greater than one. We propose in the following a graphical representation that shows each measure as a separate dimension of a multidimensional plot.

Fig. 4 (a) shows a two-dimensional scatter plot for displaying the normalized cost-quality data of an existing HM found by post-game analysis [64] (see Table II). It also shows the normalized data for three new HM's we found earlier by using population-based learning [19]. In our previous study, point-based learning was not applied to modify a HM in successive tests of the HM. This plot is obtained by applying the following steps.

First, for each set of test cases evaluated by a HM, we determine whether there are any statistical outliers in their quality and cost values. We compute standardized scores $((x - \mu)/\sigma)$ for quality and cost data individually, and consider a data point as a univariate outlier when one of its standardized scores exceeds a threshold, say 8. Note that this computation does not assume any underlying distribution of the data. An additional method, which we did not use, is to find multivariate outliers by computing the Mahalanobis distance [55] for each data point. A data point is considered an outlier when its Mahalanobis distance has less than, say, 0.001 probability of happening. The computation of the Mahalanobis distance requires the underlying data to be multivariate normal; hence, if normality is not true, then the next step must first be applied before computing the distance.

Next, we determine the distribution of the bivariate data points. We first check for normality of each of the univariate distribution by computing measures such as the skewness and the kurtosis [55], or by applying goodness-of-fit tests such as the Kolmogorov-Smirnov and the Geary tests, or by applying the Shapiro-Wilk test [12]. The joint distribution can also be checked for bivariate normality by computing the chi-square plot or the gamma plot [21] and a number of other methods [12]. If the performance values do not have a normal distribution, then a distribution must first be assumed in order to transform the values into another set so that they have a normal distribution. Note that this and the previous steps may need to be repeated a number of times in order to assure that the data points have a normal distribution and to remove outliers from the data points.

Finally, the contour of a constant probability region is plotted for the bivariate distribution of normalized values in order to indicate the spread of the data. To uniquely characterize the region, we can compute its centroid, whose coordinates are found by calculating the averages of individual performance measures (assuming all points have equal weights). Note that the reference HM has centroid at (1,1). If the original normalized performance values are not normally distributed and a transformation was made in the last step, then the contour found in this step will have to be converted, using the inverse transformation, into one for the original data points.

Fig. 4 (a) shows the result after applying the above steps for the three new HM's learned. These HM's and the best existing HM (used as the reference) are evaluated with respect to test cases used during learning, test cases of the same size but not used during learning, and test cases of large size. The spread of the performance data for each HM evaluated on the same set of test cases is indicated by an elliptical region

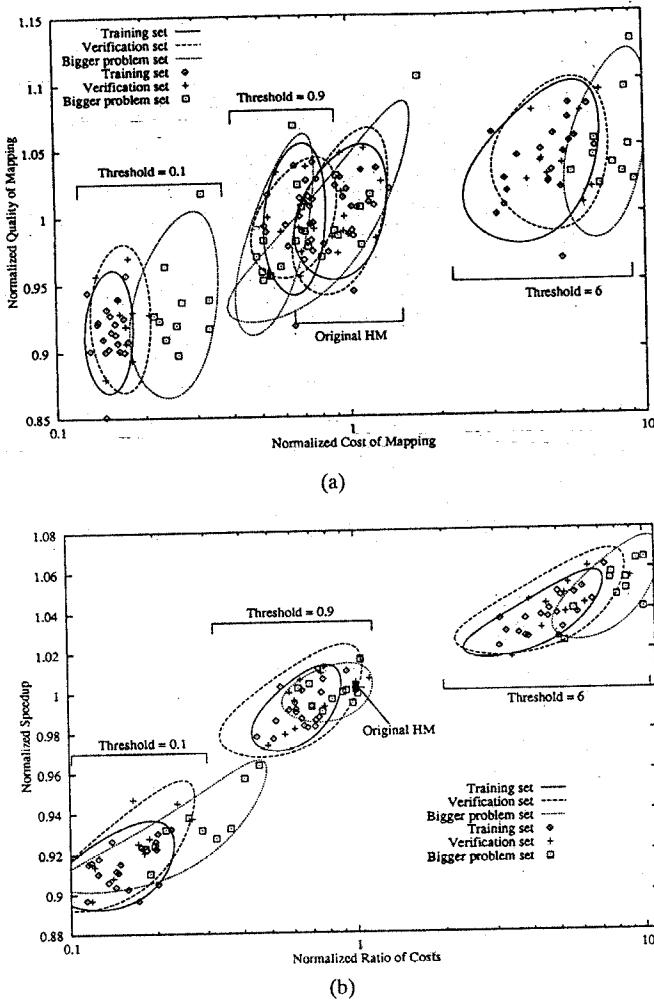


Fig. 4. Scatter plots comparing three HM's learned by population-based learning with respect to a reference HM found by post-game analysis [64]. (a) Using normalized quality and cost measures. (b) Using normalized speedup and ratio-of-cost measures.

containing statistically 90 % of the data points. Note that our normalization method allows performance to be compared with respect to test cases of different sizes and types.

The various thresholds in Fig. 4 (a) represent constraints on costs of mappings, normalized with respect to the average cost of the original post-game heuristics. A HM generated during learning is pruned if its average normalized cost exceeds the given threshold by 10 %. This issue is explored in more detail in the discussion on the internal critic in Section V-B.

Fig. 4 (a) clearly depicts the relationship of one HM with respect to another and their scalability. Here, we do not try to prove scalability but to verify it statistically. Note that it is possible to use test cases of different types in the learning system to ensure that the HM learned can be generalized. This is discussed in Section V-B.

The drawback of Fig. 4 (a) is that it only shows the distribution of performance values for different test cases but not the relationship for a particular test case between the new and the reference HM's. For instance, when we find a new HM that has a centroid better than the centroid of the reference HM, we cannot tell whether the quality (respectively, cost) of all or most of the test cases evaluated by the new HM

are better than that of the reference HM. To cope with this problem, it is necessary to use a measure that relates the quality (respectively, cost) of a new HM to the quality (respectively, cost) of the reference HM for a given test case. Examples of such a measure are the ROQ (or speedup) and the ROC measures shown in the lower half of Table V. Recall that the quality of the process mapping problem is taken as the completion time of the processes mapped.

Fig. 4 (b) shows a scatter plot using the speedup and the ROC as performance measures. Each data point represents the quality (respectively, cost) of evaluating a test case by a new HM normalized with respect to the quality (respectively, cost) of evaluating the same test case using the reference HM [64]. As is shown in Table V, every data point of the reference HM after normalization is at (1,1). The thresholds in Fig. 4 (b) have similar meaning as the ones in Fig. 4 (a), except that they are defined with respect to the ROC measure. Note that the objective of what constitutes a good HM in Fig. 4 (b) is different from that in Fig. 4 (a), and that the three HM's depicted in Fig. 4 (a) are different from the ones depicted in Fig. 4 (b).

B. Metaheuristics Used in the Learning System

The metaheuristics used in the learning system include the ones for the internal critic, heuristics manager, and test-case manager. We describe the possible approaches in each area, leaving the details to the references.

1) Internal Critic: The internal critic serves two purposes. First, it relates the global evaluative feedback to the performance objectives of the HM tested. Due to a lack of domain knowledge for structural credit assignment, it is not able to correlate the feedback to individual elements of the HM tested. Rather, it updates the performance measures of this HM and determines whether the HM is to be discarded or be inserted into the HM pool in the heuristics manager for further testing in the future. If the HM tested is to be inserted into the HM pool, the internal critic should also rank this HM with respect to other HM's in the pool.

The second function of the internal critic is to associate the aggregate performance of the HM tested with respect to the effectiveness of the operators used in generating this HM. This information allows the heuristics manager to choose effective operators for composing new HM's in the future.

An important issue to be considered in designing the internal critic is the criteria used for determining the relative "strength" of the HM tested. One possible way we had attempted was to compute a performance index for each HM according to a single parametrized performance function of multiple performance measures. When learning began, a unique combination of parameters were assigned for this function. The performance index of each HM was updated after tests were performed, and HM's were selected for testing by the resource scheduler according to their indexes. By repeating the learning process a number of times, each with a different combination of parameters, we hoped to find a combination that resulted in a HM with the proper trade-off between cost and quality. We had applied this approach for learning new HM's for the process

mapping problem shown in Table II [18], [19] and in tuning parameters in a stereo-vision algorithm (to be discussed in Section VI) [46]. Although we found improvements over the original HM's developed by human experts, this approach is fallible because, as is shown in the last subsection, whether a HM is better or not (that is, whether it has a higher relative "strength" or not) may depend on the parameters used in the single aggregate performance function.

To avoid anomalous ranking of HM's during learning due to the combination of parameters chosen, the objective used by the internal critic should be reformulated in such a way that it does not depend on the weight placed on any performance measure. An example of such an objective, as is discussed in Section V-A with respect to post-learning evaluation, is to constrain all (respectively, all but one) of the normalized performance measures and find a HM that satisfies all the performance constraints (respectively, satisfies all but one of the constraints and optimizes the one not constrained). A HM is pruned from further testing when one of its performance constraints is violated. Ranking of HM's used during learning is, therefore, consistent with the method proposed in the last subsection for the post-learning evaluation of HM's.

As an example, one possible objective in learning may be to find a new HM with the highest average quality, assuming that it is tolerable to have an average cost τ times as much as that of the original HM devised by human experts (chosen as the reference HM), where τ is a positive constant.

In this case, the internal critic computes two aggregate performance measures, each based on an individual normalized performance measure, after testing for outliers and normalizing the performance data at the end of a quantum. The aggregate measures, for instance, can be the average normalized quality and average normalized cost of the region of scattered performance points. A HM is deemed unacceptable and pruned from further learning if its average normalized cost exceeds the constraint τ . HM's shown in Fig. 4 were learned this way with τ taken as 1.1 times the threshold values.

Another objective in learning may be to find new HM's with smaller average costs and higher average quality than the ones of the original HM; in this case, a HM is pruned when either or both of the constraints are violated.

2) Heuristics Manager: The heuristics manager maintains existing HM's according to a partial order defined by the performance criteria used in the internal critic, and either selects an existing HM or generates a new one for testing when requested by the resource scheduler.

When a new HM is to be generated, the heuristics manager applies a combination of generation operators to generate one or more new HM's. The operators used are generally model free, since domain knowledge for structural credit assignment is missing. Common operators used are syntactic in nature, taking the structure of one HM and transforming it into another. It is important that the operators used be semantically sensible, generating new HDE's or HM's that are appropriate for the application. Checking for semantic sensibility is usually application dependent and is hard to generalize.

The generation operators can be classified into ones designed for modifying HDE's and ones for creating new HM's.

Operators that modify HDE's are generally grammar based, transforming a HDE from one form to another. Given a HDE, three such operators are possible: 1) *specialize*, in which one HDE is decomposed into multiple ones, each with more specialized preconditions, 2) *generalize*, in which two HDE's with similar behavior are combined into a general one with more general preconditions, and 3) *mutation*, in which a new HDE is created by perturbing parameter values of an existing HDE. These operators are similar to the ones used in classifier systems [14]. A partial history on the application of these operators may have to be kept to avoid repeatedly decomposing, recombining, and regenerating the same HDE.

Operators for creating new HM's, in the absence of a good model for structural credit assignment, are similar to the crossover operator used in classifier systems [14] for "breeding" new HM's. In creating a new HM, two promising HM's in the HM pool are first identified; HDE's from each are then clustered into two new HM's. The HDE's taken are either ones that, when applied individually or in a sequence, result in local improvement in performance, or ones that bring the solution closer to the final goal. Two contradicting goals need to be observed by the creation operator: 1) it should avoid creating HM's with very similar performance, which usually happens when a large number of HDE's are duplicated in the new HM's generated; 2) it should preserve HDE's that lead to high performance and use them as building blocks in the new HM's generated.

The generation process applies a combination of the operators described previously to arrive at new HM's. In the beginning, it uses no historical information on performance and chooses the operators randomly. As more tests are performed, the operators which are used, as well as the HDE's or HM's to operate on, can be based on success rates of different operators or operator sequences applied in the past.

3) : The test case manager is responsible for maintaining a pool of valid test cases and for generating new test cases if necessary. The test cases used for evaluating HM's should satisfy two requirements. First, they should have a common statistical behavior, thereby making it meaningful to evaluate the statistical performance of HM's. Second, their performance must be known with respect to the reference HM. This is needed because all performance measures are normalized with respect to the corresponding measures of the reference HM, as is discussed in Section V-A. Since the reference HM is likely to be less effective than the HM to be learned, evaluating the reference HM during learning may consume valuable learning time. As a result, the evaluation of the reference HM may have to be carried out before learning starts.

In designing the test case manager, two issues need to be considered. First, it must determine the size and type of test cases to be used. If the test cases are small, then many of them (and consequently HM's) can be evaluated during learning; but the HM derived from small test cases may not be scalable to larger ones. On the other hand, if the test cases are large, then fewer test cases and HM's can be evaluated in a fixed amount of time. Similar trade-offs need to be made in deciding whether test cases from different applications should be used. If the HM's to be learned need to be scalable and be general, it

may be necessary to use test cases of different sizes and from different applications. This requires more complex resource scheduling algorithms, as there may be large variations in the evaluation times for these test cases.

Second, the test case manager must decide for all HM's whether a random or a predefined order of test cases are to be used. A predefined order is preferable because the number of data points obtained during learning for a particular HM may be relatively small, and it may not be meaningful to compare two HM's statistically based on a small but different set of test cases. Moreover, the evaluation of two HM's using different test cases may require the evaluation of the reference HM on all of the test cases used — an additional overhead needed before learning begins. It is also important that test cases used be statistically independent, and that the generation of test cases be unbiased. Statistical correlation of performance values may be used for detecting any bias.

C. Resource Scheduling

Resource scheduling entails the allocation of a limited amount of resources for generating and for testing HM's. Resource scheduling algorithms can be static or dynamic. A static scheduling algorithm determines, before learning begins, the number of generations in the learning process, the number of HM's to be tested in each generation, and the number of tests performed on each HM. Examples of static scheduling algorithms have been studied in classifier systems [6], in which trade-offs have been evaluated with respect to different numbers of generate-and-test phases in a fixed amount of time [14] and different overheads in generation [10]. On the other hand, a dynamic resource scheduling algorithm determines dynamically, based on performance statistics obtained during learning, whether a new HM is to be generated or an existing HM is to be tested. This decision may depend on: 1) whether existing HM's that perform well have been tested to an adequate degree of confidence, 2) whether the performance of all existing HM's falls below a required minimum, 3) whether there are an insufficient number of HM's in the population, and 4) whether there are plenty of resources remaining. Dynamic scheduling algorithms are potentially better because they avoid testing HM's that have been found to be poor with some initial tests.

The top half of Table VI shows the characteristics of the general problem for dynamically scheduling the generations and the tests of HM's. Assume, without loss of generality, that HM's are characterized by quality and cost, and that the objective of learning is to optimize one of the measures, say quality, while constraining the other. (The performance values of feasible HM's are, therefore, a collection, called a *population*, of quality values, with constraints on costs satisfied.) In the general model, performance analysis is difficult due to the fact that many of the parameters of the model are unknown *a priori*, and that point-based learning and heuristics generation introduce unknown effects on the search space of HM's. A possible approach is to develop efficient resource scheduling algorithms based on a simplified model and to use them as heuristic algorithms for the general model. The assumptions

TABLE VI
GENERAL AND APPROXIMATE STATISTICAL MODELS

General Statistical Model

- Statistical distributions of populations are unknown and are possibly nonidentical.
- The distribution of means of populations is unknown.
- Applying point-based learning during testing of a HM leads to minor modification of the HM and a different population of performance values examined.
- Heuristics-generation algorithms (such as cross over, mutation, specialize, and generalize) incur a biased search of the space of HM's.
- Time to test a HM on a test case may be stochastic with an unknown and possibly population dependent distribution.

Approximate but Tractable Statistical Model

- All populations have normal distributions with unknown means but constant variances.
- The means of all populations have a known distribution function with unknown parameters.
- Point-based learning makes negligible changes to the HM tested, and hence, the population of performance values examined.
- Heuristics-generation algorithms incur a random search of the space of HM's.

The time to test a HM on a test case is constant. Alternatively, the average time may be used as an approximation.

for a simplified but tractable statistical model are shown in the lower half of Table VI.

1) *Related Work:* An approach for scheduling resources in learning was studied by Russell and Wefald [43], [59]. The major differences between their work and ours are that they assume a knowledge intensive domain (with at least a causal model of the environment defined), and a well-defined objective function for the application problem. The Rational Agents with Limited Performance Hardware (RALPH) project they developed aims to combine logicians' approach, which emphasizes correctness, and rationalists' approach, which emphasizes utility maximization. The compromise — bounded or limited rationality — is useful for real-time problem solving,² and can be stated as cost-quality trade-offs for the problem solver. They considered scheduling not only alternative moves (called "actions") but also search strategies (called "computations" or "deliberations") that could help improve the incumbent's expected utility. Learning was based on a search method called DTA* [59] which used a common evaluation function for choosing among alternatives. Their approach was to learn from past experience the probability of change in utility (respectively, cost) of a move *i*, given that a computation *j*, had been performed, for all *i* and *j*. Using the distributions learned, they derived the value estimated for a move or computation as the difference between its estimated benefit and its estimated cost.

A problem similar to the heuristics selection problem we address here has been studied extensively as *selection problems* in statistical inference [4], [15]. These problems can be classified into two types.

²Real-time problems in Russell and Wefald's studies [43] are defined as those in which the utility of an action is a decreasing function of the time spent in finding that action. The trade-off is represented as a time cost function, which gives the loss in utility caused by delay in performing an action. We do not assume such a utility function in this paper.

The first type was first studied by Stein in 1948 [52] and is called the *stopping problem*. It deals with finding the minimum number of tests in order to know which population among a finite number of populations is the best to within a certain degree of confidence. This result and subsequent extensions cannot be applied in our case because they considered a finite number of populations and an unbounded time.

The second type was pioneered by Bechhofer in 1954 [2], [3] and is called the *allocation problem*. It focuses on allocating a fixed number of tests among a given finite number of populations so that the probability of selecting the population with the maximum mean is maximized. Recent work in this area is predominantly on multistage methods in which the number of tests allowed are divided into stages, and promising populations selected in one stage are passed to the next stage for further testing. The allocation of tests can be static [5], which is determined *a priori* based on statistical distributions of populations, or can be dynamic [57], which is based on means and variances estimated during testing.

Results in the second type do not apply in our case because they deal with a finite and given number of populations before selection begins; our objective in HM selection, on the other hand, is to find a population (HM) with a large mean value but not necessarily the one with the largest mean, as there are infinitely many populations to start. However, existing statistical algorithms can be applied if the time allowed were divided into two parts: the first part is used for determining a finite number of populations to be tested in the second part, and the second part assigns tests according to algorithms developed in statistical allocation.

2) Resource Scheduling Algorithms Studied: Based on the assumptions defined in the lower half of Table VI, we have developed a multistage selection algorithm for selecting a promising population within a time constraint. (Time is the only resource to be scheduled.) In this subsection, we summarize our previous results in this area [18], [19].

In the general multistage algorithm, $G(T)$, stage i is characterized by a triplet consisting of 1) g_i , the particular allocation algorithm to be used for this stage, 2) t_i , the duration of the stage, and 3) n_i , the number of populations to be considered for testing during the stage [18]. It starts with a presampling stage that estimates the statistical parameters of the target application. This is chosen heuristically to be 10 % of the total time allowed. The next stage corresponds to coarse initial testing to weed out unworthy populations; each subsequent stage then evaluates populations not eliminated in previous stages and passes a small number of promising populations to the next stage. This general multistage algorithm covers both static and dynamic allocation algorithms.

Analysis of the general multistage algorithm centers on finding for each stage the joint probability distribution between \bar{x} , the sample mean of a population, μ , the sample mean of all population means, σ , the standard deviation of all population means, s , the number of samples drawn from each population, and any other variables used in selecting the next HM for testing [18]. Based on the assumptions stated in the lower half of Table VI, we have developed methods for analyzing both static and dynamic algorithms. However, analysis of dynamic

algorithms is generally intractable because each pick of a population for testing must be treated as a separate stage.

We have studied and evaluated the following scheduling algorithms [19]:

- single-stage round-robin algorithm, where a predetermined number of populations are sampled a constant number of times in the time³ allowed;
- two-stage round-robin algorithm, where analysis aims at finding the division of time between the first and the second stage and the parameters of the round-robin algorithm in each stage;
- two-stage round-robin/greedy algorithm, which uses (as heuristics) the same division of time found analytically for the two-stage round-robin algorithm;
- two-stage round-robin/minimum-risk algorithm, which is similar to the last algorithm except that it selects populations for testing in the second stage based on a risk function (defined as the expected squared-error loss of the estimated mean of each population).

For the last two algorithms, performance can be formulated but is too difficult to be solved in closed forms.

We have applied the above scheduling algorithms for learning new process mapping HM's for post-game analysis [64]. We first verified the assumptions stated in the lower half of Table VI with respect to a randomly selected HM for process mapping. We developed a learning system on a DECStation 5000/200 workstation. Assuming a total CPU time of 16 h, the available time is divided into two parts, with 12 h and a round-robin scheduling algorithms in the first part, and 4 h and a minimum-risk algorithms in the second part. The performance results of the three new HM's learned are shown in Fig. 4.

3) Possible Extensions: In general population-based learning, the scheduling algorithms we have studied so far are inadequate because they do not consider improvements caused by the heuristics generator and point-based learning. In both cases, the received evaluative feedback signals are used to guide either the generation of a new HM or the refinement of an existing HM; the search for new HM's is, therefore, biased. It is this bias that helps learning systems improve over time; however, the statistical modeling of this effect is very hard. We expect that heuristic scheduling algorithms will be used in the general case. Some possible heuristic approaches for addressing this problem are as follows: assuming that heuristic scheduling algorithms are defined with respect to the number of HM's to be generated in each generate-and-test phase, the duration of the testing period, and the number of tests performed on each HM.

In determining the duration of a generate-and-test phase and the corresponding number of populations to be generated, we observe that the populations tested in previous phases may affect the populations to be generated or to be tested in subsequent phases. If populations tested in each phase are generated randomly, then we have found that the optimal algorithm generates all the populations ahead of time before learning begins [20]. If populations are generated on demand, based on

³The time constraint is expressed as the total number of samples drawn, assuming that each draw takes unit time.

a posteriori distribution of the population means found earlier, then Bayesian analysis can be applied to capture the effect of changes in the *a posteriori* distribution of population means.

In finding how long a population is to be tested before switching to a different one, we observe that the distribution of populations tested changes with time when point-based learning is applied. In this case, we expect techniques in trend and time-series analyses to be useful for determining when testing should be stopped.

In summary, the scheduling of resources for an integrated learning system is complex because the space in which HM's are generated changes with time and may be hard to model. However, experimental results have shown that intermediate statistical measurements are useful for determining when to generate or what to test next, especially in a learning environment that lacks a world model. We expect that new insights can be gained by developing scheduling algorithms based on simplified models and applying these strategies as heuristics in practice.

VI. EXAMPLE APPLICATIONS

In this section, we illustrate the techniques developed in this paper on four applications: 1) load balancing of independent processes, 2) automated design of artificial neural networks, 3) depth perception in stereo vision, and 4) combinatorial search using a branch-and-bound algorithm. These applications have many common features that allow new HM's to be learned by an integrated population-based and point-based learning method. However, the first two applications have unique idiosyncrasies that need special considerations: 1) they both require the design of a performance predictor that predicts delayed feedback signals; 2) they both require off-line learning; and 3) the design of load balancing HM's further requires simplifying assumptions to be made in order to address the problem of overlapping and dependent temporal scopes. For the last two applications, on-line population-based learning can be applied in a similar way as in the mapping of communicating processes on distributed memory computers (see Table II). Note that our discussion only serves to illustrate difficulties involved in applying the learning method; implementations and performance evaluation of the learning experiments are generally application dependent.

A. Load Balancing of Independent Processes on a Distributed Computing System

In this subsection, we discuss issues and solutions for learning the design of adaptive and machine dependent HM's for the load balancing problem defined in Table I. There are three issues that we consider: the first is related to the development of an environment suitable for learning; and the last two are related to learning workload indexes and load placement HM's.

1) *Application Environment for Learning:* The first issue concerns the reproduction of an application environment suitable for testing alternative placement and migration HM's. In our case, a technique is needed for reproducing a background workload and for adjusting the background

workload in the presence of a test process. This is especially difficult in a time sharing operating system as the execution of multiple processes may be interleaved in time. Desirable but prohibitively expensive techniques are ones that can reproduce the execution of every machine level instruction. For practical reasons, we have developed a method that periodically samples the primitive measurable attributes in a computer system, such as disk traffic and memory occupancy (specifically, at every real-time clock interrupt), and that replays them later on by generating synthetic workload with the same behavior [29]. This sample-and-replay method has been found to be very accurate in approximating the physical-level system behavior.

2) *Learning Workload Characterization:* The second issue is on the design of a workload measure for identifying a location for a given process to be executed. This measure serves as the performance predictor in Fig. 1. It is a complex function of a large number of machine dependent, out dated, and possibly dependent measurable attributes in the system. Existing techniques are based on either analysis, such as queuing models, or heuristic metrics, such as the number of pending processes, resource utilization levels, and *ad hoc* averages of queue lengths of selected resources [66].

We have developed a systematic method for learning a machine dependent workload measure as a function of primitive measurable attributes [30]. Our design is based on an extension of Tesauro's comparator neural network [56] and generates an output indicating the *relative* speedup of a process with respect to a reference workstation, using inputs derived from windows of primitive measurable attributes before the process is executed. To simplify the discussion, consider a network of homogeneous workstations; in this case, the design of the comparator network is identical for all sites. In order to collect the training patterns for the comparator network, we run a set of test processes, one at a time, under a variety of controlled workload conditions on one workstation. A training pattern is then made up of the completion times of a test process run under two different workload conditions, which are characterized by windows of measurable attributes before the execution of the process. To reduce the amount of input information to the comparator neural network, we have developed a scheme for filtering the windows of workload attributes so that only their trends are used for predicting future workload conditions.

3) *Learning Load Balancing HM's:* Given the workload measures learned, the final issue is on learning automat placement and migration HM's. Learning involves the adjustment of parameters in a rule-based representation of the HM, which controls when migrations should be initiated. Structural credit assignment is difficult, as there is little background knowledge available (see Table VII). Temporal credit assignment is also difficult because process placement decisions may be related in an unknown fashion and tend to have overlapping and dependent temporal scopes. Moreover, on-line learning is impractical because there are a large number of alternatives for placing and for migrating a single process during the time that the process is executed.

In order to learn process placement and migration HM's by population-based learning, we need to make two simplifying

TABLE VII
LEARNING HM'S FOR THE LOAD BALANCING PROBLEM DEFINED IN TABLE I

Credit assignment of the original problem:

- Structural:** Difficult, as domain is knowledge lean.
- Temporal:** Episodic, but load balancing decisions may be dependent and have overlapping temporal scopes.

Simplified problem:

- Restrict the load balancing decision to one process at a time, independent of all future migration decisions for the given process; all other processes in the system are considered as background workload and are not modeled specifically.
- Correct wrong decisions by dynamic process migrations.

Off-line learning:

- Rules learned are not scalable between small and large processes.
- Large processes have large temporal scopes.
- The possible decisions at a given time are limited by the number of computers in the system.
- Alternative decisions can be tested by regenerating an identical system-wide workload.

assumptions so that these HM's satisfy the requirements stated in Table III for temporal credit assignment. First, we need to assure that performance measures for evaluating a test case are independent of the intermediate decisions and internal states. This condition will be valid if we restrict the temporal scope of a placement/migration decision to be only between this decision and the next placement/migration decision for any process in the system. This assumption implies that at every decision point in the system, the placements of all newly arrived processes and partially executed processes are considered anew. It is true when the overhead for migrating a process is negligible. Note that the overhead for migration can be accounted for indirectly by the migration rules, which migrate a process only when there is a reasonable difference in workloads between two computers. Second, we need to assure that the aggregate performance measures are independent of the order of evaluating the test cases. This condition will be valid if we assume that the future background workload in the system is not affected by a placement or migration decision. This assumption allows the placement and migration rules to be learned one process at a time in the presence of a predefined background workload. It is debatable, since one placement decision may affect other placement decisions, which in turn affect the background workload. However, in the absence of information about future workloads and arrivals, and the fact that errors in placements can be corrected later by migrations, this assumption is a reasonable approximation.

The previous assumptions reduce the temporal credit assignment to one process at a time: the response time feedback signal observed for one process is used to modify the placement and migration rules used in arriving at the decision for this process, and all other processes executed in the system are modeled as background workload. Learning is, therefore, episodic, and the performance measure of a test case (in terms of the completion time of the test process) is independent of intermediate decisions and states. Furthermore, the aggregate performance measure of the process placement and migration system (in terms of the average completion time of a set of test processes) is independent of the order of presentation of

TABLE VIII
THE CONFIGURATION DESIGN OF AN ANN

Definition:

Given a set of input-output patterns of the target application, the problem is to find a feed-forward ANN trained by the back-error propagation procedure so that the TSSE between the actual and the desired outputs of the ANN is below a maximum threshold.

Characteristics:

- The set of measurable parameters are ill defined.
- Relationship is undefined between the cost (or complexity) of the network and its quality (or TSSE) when training is terminated.

Credit assignment in training a given network:

- Structural: Defined by the back-error propagation procedure.
- Temporal: None, as feedback is immediate.

Credit assignment in the design process:

- Structural:** *Ad hoc*, due to a lack of domain knowledge.
- Temporal:** Episodic, but the temporal scope of a decision to modify a network is as long as the time it takes to train the network.

On-line learning for the design process:

- Use existing schemes [8], [9], [24], [37], [42], [47], [49] for modifying network configurations in the heuristics manager.
- Use a performance predictor for predicting missing feedback signals.

test cases.

The simplified problem can now be learned by off-line population-based learning. Learning begins by collecting the completion times of executing a given set of test processes, one at a time, on all possible sites in the network. By assessing the workloads for all sites before the execution of a process in the set (using the comparator neural network designed) and by knowing the completion times of the process for each site, we can apply population-based learning to tune the parameters used in the load balancing system. Learning can be carried out off-line because 1) all possible placement decisions at any time are limited by the number of computers in the system; 2) an identical system wide workload can be regenerated so that performance due to alternative decisions, starting from a given initial workload condition, can be evaluated before hand; and 3) for a given set of rules for detecting when migration is to be initiated, the completion times of the test process for all possible locations for migrating the process can be evaluated off-line.

B. Configuration Design of Artificial Neural Networks

In this section, we apply our integrated learning method for designing an appropriate configuration of a feed-forward artificial neural network (ANN) for an application problem. The target application is specified by a set of input-output patterns, and the design problem is to find an ANN that responds correctly to any input pattern in the set (and perhaps generalizes its outputs for new input patterns). The objective of the design process is to identify a "good" ANN configuration (but not necessarily train it to completion) using a fixed amount of computational resources. Table VIII summarizes the issues for this problem.

1) Existing Methods: Traditionally, a feed-forward ANN can be designed by first hypothesizing a network configuration with some initial weights, and by applying the back-error propagation procedure [28], [42] or its variations to "learn"

the correct weights for the network. It is generally assumed that the network is trained when its total sum-of-squared error (TSSE) between the actual and the desired outputs falls below a maximum threshold, even though training may have only converged to a local optimum. If training does not converge within a reasonable amount of time, then heuristic schemes will be applied to either pick a new configuration or modify an existing configuration before further training. Examples of these heuristic schemes can be classified into five types: the *ad hoc* method [42], the iterative refinement method [8], the start-big-and-remove method [24], [37], [47], the start-small-and-add method [9], [49], and the iterative refinement with backtracking method [58]. All these schemes except the last one are based on a greedy search that only examines one incumbent network at any time. The last scheme examines multiple partially trained networks simultaneously; however, it lacks a general measure for identifying promising configurations and a procedure for scheduling computational resources.

The traditional design process can be considered as a hierarchy of two levels of point-based learning. Let an ANN be a HM for solving the application problem. At the lower level, we can view the test cases as the set of training patterns specified for the application, a decision in a HM as the selection of a training pattern, and an episode as the time for training the network with one training pattern using the back-error propagation procedure. Structural credit assignment for the ANN selected for training is provided by the back-error propagation procedure. Temporal credit assignment, on the other hand, is not necessary because the feedback signal for each input pattern trained is available without delay.

At the higher level, the design process modifies the ANN selected for training when little improvement on TSSE is found. Many such existing schemes [8], [9], [24], [37], [42], [47], [49] are point based, focusing on one incumbent network at a time. These schemes use immediate evaluative feedback, based on past training performance of the network, to apportion credit structurally to components of the incumbent network. A decision is, therefore, an application of an operator for modifying the network, and its temporal scope ends when either the network is fully trained or time runs out. Temporal credit assignment is not addressed, as the final feedback (whether training converges or not) is related in a complex fashion to the intermediate decisions and states.

The previous discussion reveals four related issues on the automated design of an ANN. First, there is an infinite number of many possible networks, and each may take an exceedingly long time to train; hence, it is impossible to enumerate and train all of them to completion within a fixed amount of time, space, and resource constraints. Second, expert knowledge for predicting good networks is *ad hoc* in nature and may be application dependent, rendering it difficult to use in the design process. It is, therefore, essential to develop a way to capture and generalize this expert knowledge based on training examples for various applications. Third, the objective of the network to be designed is ill defined, based on a subjective trade-off between training time and network cost. Hence, it is important to have a design process that can propose alternate

configurations under a variety of cost quality trade-offs. Last, the aggregate performance measure of the ANN designed, in terms of whether training converges or not, is not known until either the network is fully trained or time runs out. Further, it is difficult to relate this aggregate measure to the large number of intermediate performance measures available during training; examples of the latter include the current TSSE, the training time expended so far, and the rate of change of TSSE.

2) Integrated Learning of ANN Configurations: The issues discussed previously can be resolved to a large extent by applying population-based and point-based learning in the design process. By maintaining a population of competing ANN's, the design system can defer decisions on the network to be trained, modifications to be performed, and new networks to be generated, until some partial training has been carried out on each. Population-based learning allows network modifications without relying on accurate domain knowledge for structural credit assignment. Existing schemes for modifying network configurations [8], [9], [24], [37], [42], [47], [49] can be used as operators in the heuristics manager; the particular operator to apply can be deferred until each network in the pool has been trained partially. On the other hand, point-based learning is carried out when a network is selected for training. As stated before, the back-error propagation algorithm provides structural credit assignment in the training process.

The unique difficulty with using population-based learning is that neither the conditions defined in Table III for on-line learning nor those defined for off-line learning are satisfied. The design process cannot be learned off-line because generating the state space of realistic ANN's amounts to finding an ANN for solving the target application. The conditions defined for on-line population-based learning are not satisfied because the temporal scope of each episode is very large, and training cannot be performed on small networks and generalized to larger ones. In order for on-line population-based learning to be carried out, it is necessary to develop a heuristic performance predictor that discerns promising networks from poor ones with only partial training. This predictor estimates for the internal critic the relative aggregate performance of a given ANN as compared to the performance of other competing ANN's in the population.

The performance predictor needed here is similar to the one for workload characterization in load balancing. In this case, inputs to the predictor consist of time varying patterns of TSSE values versus time for two different networks in the population, and the problem is to differentiate between them as to which one is more promising. Here, a network is more promising than another means that they both satisfy the given cost and complexity constraints, while the training time of the former is predicted to be smaller than the latter. The design of the performance predictor is simpler than that for the load balancing problem because traces of TSSE values versus training times for various ANN's can be collected ahead of time. Our initial results show that such a predictor can discern which one of two ANN's will converge faster with over 75 % accuracy after each has been trained for a reasonable number of epochs.

C. Online Learning for Stereo Vision and Combinatorial Search

In this section, we discuss the application of online population-based learning for two new applications. The way that population-based learning is applied here is similar to that for the mapping of processes on a network of distributed memory computers, which was discussed earlier [19] (see Table II and Fig. 4).

The first application deals with stereo vision, an inexpensive method for calculating the depth information in a scene. This is done by matching and triangulating the corresponding features in the scene taken from two different view points. A general stereo-vision algorithm refines iteratively the depth information by matching increasingly detailed features (tokens) found in the scene. This is necessary because with highly detailed scene descriptions, it is hard to find unique correspondence between the features without any initial depth knowledge. The iterative method is also considered to be a plausible approach used in biological vision systems [27].

There are two issues in stereo vision that are related to learning. First, there are a large number of parameters to be tuned in a stereo-vision algorithm. The values of these parameters, as well as the granularity of the features and the number of interactive refinements, are influenced strongly by the imaging domain. It is, therefore, essential to develop a systematic method for tuning these parameters in limited time. Second, the objective of a stereo-vision algorithm is ill defined with respect to measurable attributes in the system, such as speed of the algorithm, density of depth information determined, and accuracy of depth measurements. Hence, it is important to have a design process that can propose alternate parameter values for a variety of cost-quality trade-offs.

The HM's for this problem are suitable to be learned by an on-line integrated learning approach for similar reasons as the process mapping problem (see Table II). Our initial results show that highly accurate and efficient HM's can be found by population-based learning [46].

The second application deals with combinatorial search problems that can be solved by a branch-and-bound search. A branch-and-bound algorithm is an organized search of the space of all feasible solutions for solving a combinatorial optimization problem. By repeatedly partitioning the space of all feasible solutions into smaller and smaller subsets until either all feasible solutions are found, or infeasibility/dominance is encountered, the process guarantees that the optimal solution is found.

There are many existing HM's for exploring the search tree; examples of HM's amenable to learning are the selection function for defining what node to expand next, the decomposition function for decomposing a search node into descendent search nodes, and the heuristic pruning function for identifying unpromising search nodes. Each of these functions can be formulated in terms of parameters that can be measured during the search. For reasons similar to the ones shown in Table II, these functions can be learned by on-line integrated learning methods. Note that HM's in this application have only one performance measure, namely, the number of nodes searched before the solution is found, and that they are all performance

related. The population-based learning method is not suitable for learning correctness related heuristics, such as dominance relations [65].

Our initial results on learning new decomposition functions show that significant reduction in the search overhead can be attained. One problem with the integrated learning method for this application is that the functions learned may not be semantically meaningful: it is possible that functions learned are tailored to the specific test cases used. Consequently, it is important that these functions be interpreted after they are learned; if necessary, the learning process will need to be repeated using different test cases.

VII. CONCLUSIONS

In this paper, we have presented a learning-by-example method that integrates a point-based approach with a population-based approach for learning new HM's for problem solving. Traditional learning approaches are generally point-based in nature, focusing on modifications of an incumbent HM. In a population-based learning approach, a population of competing HM's are maintained; the choice of the HM to test depends on partial results of testing the HM's on test cases.

We have focused on learning performance-related HM's that involves trade-offs between the quality of solution found and overhead (or cost) incurred in finding the solution. We are interested in HM's that do not have a well-defined relationship between quality and cost. We have defined the characteristics of application problems whose HM's can be learned by the proposed learning method (see Table III). We have addressed issues of anomalies in performance evaluation of HM's, criteria for evaluating HM's during and after learning, acquisition of metaknowledge for the learning system, and design of efficient scheduling algorithms for finding good HM's in limited resources.

We have illustrated our learning method on five applications. HM's for three of these applications (processing mapping, stereo vision, and combinatorial search) can be learned in a straightforward way. Two other applications that require special considerations are load balancing of processes for a distributed network of workstations and configuration design of artificial neural networks. In the load balancing application, temporal scopes of episodes interfere with each other. This problem has been solved by making reasonable assumptions on the type of load balancing and migration HM's to be learned and by redefining the meaning of episodes. The temporal scopes of episodes can then be partitioned so that they do not interfere with each other. In the neural network design problem, the temporal scope of an episode is as long as the time it takes to find an appropriate network. To cope with this problem, an off-line design of a performance predictor is needed. The predictor predicts delayed feedback signals to be available in the future, and estimates for the internal critic the training time of a given neural network relative to another in the population. Note that when the predictor is very accurate, population-based learning degenerates to point based.

In summary, we have addressed in this paper the following key issues:

- methods for coping with a lack of domain knowledge for structural credit assignment in learning performance related HM's;
- temporal properties of HM's that must be satisfied in order for population-based learning to be possible;
- methods for resolving anomalous behavior in performance evaluation;
- efficient scheduling of computational resources for population-based learning.

Our solutions to these issues extends previous results on statistical allocation and learning in classifier systems. We have shown that population-based learning is a viable alternative to structural credit assignment, and that population-based learning methods are widely applicable and their conditions of applicability are characterized by the nature of their episodes.

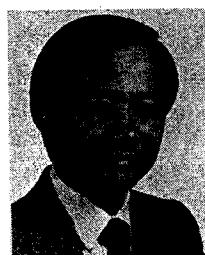
VIII. ACKNOWLEDGEMENTS

The author is indebted to many fruitful insights and help from research staff in his group in preparing this paper. In particular, the author would like to thank P. Mehra, A. Ieumwananonthachai, A. Aizawa, L. Chu, T. Karnik, and S. Schwartz for their discussions, and to A. Ieumwananonthachai for generating the results on learning new heuristic methods for process mapping. The author would like to thank Prof. H. Tanaka and Fujitsu Limited, who provided support when the author was visiting the University of Tokyo, during which part of this paper was written.

REFERENCES

- [1] A. Barr and E. A. Feigenbaum, *The Handbook of Artificial Intelligence*, Vols. 1-3. Los Altos, CA: William Kaufmann, 1982.
- [2] R. E. Bechhofer, C. W. Dunnett, and M. Sobel, "A two-sample multiple decision procedure for ranking means of normal populations with a common unknown variance," *Biometrika*, vol. 41, pp. 170-176, 1954.
- [3] R. E. Bechhofer, "A single-sample multiple decision procedure for ranking means of normal populations with known variances," *Ann. Math. Statist.*, vol. 25, no. 1, pp. 16-39, Mar. 1954.
- [4] R. E. Bechhofer, J. Kiefer, and M. Sobel, *Sequential Identification and Ranking Procedures*. Chicago, IL: University of Chicago, 1968.
- [5] R. E. Bechhofer, A. J. Hayter, and A. C. Tamhane, "Designing experiments for selecting the largest normal mean when the variances are known and unequal: Optimal sample size allocation," *J. Statist. Planning Inference*, vol. 28, pp. 271-289, 1991.
- [6] L. B. Booker, D. E. Goldberg, and J. H. Holland, "Classifier systems and genetic Algorithms," *Artif. Intell.*, North Holland, vol. 40, no. 1, pp. 235-282, 1989.
- [7] T. G. Dietterich and B. G. Buchanan, "The role of critic in learning systems," Tech. Rep. STAN-CS-81-891, Stanford Univ., CA, Dec. 1981.
- [8] S. E. Fahlman, "Faster-learning variations on back-propagation: An empirical study," in *Proc. Connectionist Models Summer School*, D. Touretzky, G. E. Hinton, and T. J. Sejnowski, eds., Palo Alto, CA, 1988, pp. 38-51.
- [9] S. E. Fahlman and C. Lebiere, "The cascade-correlation learning architecture," *Advances in Neural Information Processing Systems*, 2 ed., D. S. Touretzky, ed. Palo Alto, CA: Morgan Kaufmann, 1990, pp. 524-532.
- [10] J. M. Fitzpatrick and J. J. Grefenstette, "Genetic algorithms in noisy environments," in *Mach. Learning*, vol. 3, no. 2/3, pp. 101-120, Oct. 1988.
- [11] M. P. Georgeff, "Strategies in heuristic search," *Artif. Intell.*, vol. 20, pp. 393-425, 1983.
- [12] R. Gnanadesikan, *Methods for Statistical Data Analysis of Multivariate Observations*. New York: Wiley, 1977.
- [13] J. J. Grefenstette, "Credit assignment in rule discovery systems based on genetic algorithms," *Mach. Learning*, vol. 3, no. 2/3, pp. 225-246, Oct. 1988.
- [14] J. J. Grefenstette, C. L. Ramsey, and A. C. Schultz, "Learning sequential decision rules using simulation models and competition," *Mach. Learning*, vol. 5, pp. 355-381, 1990.
- [15] S. S. Gupta and S. Panchapakesan, "Sequential ranking and selection procedures," in *Handbook of Sequential Analysis*, B. K. Ghosh and P. K. Sen, eds. New York: Marcel Dekker, 1991, pp. 363-380.
- [16] J. H. Holland, *Adaptation in Natural and Artificial Systems*. Ann Arbor, MI: Univ. of Michigan, 1975.
- [17] ———, "Properties of the bucket brigade algorithm," in *Proc. Int. Conf. Genetic Algorithms and Their Applications*, 1985, pp. 1-7.
- [18] A. Ieumwananonthachai, A. N. Aizawa, S. R. Schwartz, B. W. Wah, and J. C. Yan, "Intelligent mapping of communicating processes in distributed computing systems," in *Proc. Supercomputing 91*, Albuquerque, NM, Nov. 1991, pp. 512-521.
- [19] ———, "Intelligent process mapping through systematic improvement of heuristics," *J. Parallel Distributed Comput.*, vol. 15, pp. 118-142, June 1992.
- [20] A. Ieumwananonthachai and B. W. Wah, "Parallel statistical selection in multiprocessors," in *Proc. Int. Conf. on Parallel Processing*, 1992, pp. 190-194.
- [21] R. A. Johnson and D. W. Wichern, *Applied Multivariate Statistical Analysis*. Englewood Cliffs, NJ: Prentice-Hall, 1982.
- [22] A. H. Klopf, "Drive-reinforcement learning: A real-time learning mechanism for unsupervised learning," in *Proc. Int. Conf. on Neural Networks*, vol. II, 1987, pp. 441-445.
- [23] M. M. Kokar and W. W. Zadrozny, "A logical model of machine learning," in *Proc. First Workshop on Change of Representation*, 1988.
- [24] J. K. Kruschke, "Creating local and distributed bottlenecks in hidden layers of back-propagation networks," in *Proc. Connectionist Models Summer School*, 1988, pp. 120-126.
- [25] P. Langley, "Learning to search: From weak methods to domain-specific heuristics," *Cognitive Sci.*, vol. 9, pp. 217-260, 1985.
- [26] D. B. Lenat, "Theory formation by heuristic search: The nature of heuristics II: Background and examples," *Artif. Intell.*, vol. 21, pp. 31-59, 1983.
- [27] D. Marr, *Vision*. New York: Freeman, 1982.
- [28] J. L. McClelland and D. E. Rumelhart, *Explorations in Parallel Distributed Processing: A Handbook of Models, Programs, and Exercises* vol. 3. Cambridge, MA: MIT, 1988.
- [29] P. Mehra and B. W. Wah, "Physical-level synthetic workload generation for load-balancing experiments," in *Proc. First Symp. on High Performance Distributed Computing*, Syracuse, NY, Oct. 1992.
- [30] P. Mehra and B. W. Wah, "Adaptive load-balancing strategies for distributed systems," in *Proc. 2nd Int. Conf. on Systems Integration*, June 1992, pp. 666-675.
- [31] R. S. Michalski, "A theory and methodology of inductive learning," in *Machine Learning*, R. S. Michalski, J. G. Carbonell, and T. M. M. Tioga, eds. Los Altos, CA: Morgan Kaufman, 1983.
- [32] R. S. Michalski, J. G. Carbonell, and T. M. Mitchell, eds., *Machine Learning: An Artificial Intelligence Approach*. Los Altos, CA: William Kaufmann, vols. 1 and 2, 1983 and 1985.
- [33] M. Minsky, "Steps toward artificial intelligence," *Computers and Thought*, E. A. Feigenbaum and J. Feldman, eds. New York: McGraw-Hill, 1963, pp. 406-450.
- [34] T. M. Mitchell, P. E. Utgoff, B. Nudel, and R. Benerji, "Learning problem-solving heuristics through practice," in *Proc. 7th Int. Joint Conf. on Artificial Intelligence*, 1981, pp. 127-134.
- [35] T. M. Mitchell, "Learning and problem solving," in *Proc. 8th Int. Joint Conf. on Artificial Intelligence*, Aug. 1983, pp. 1139-1151.
- [36] T. M. Mitchell, P. E. Utgoff, and R. B. Banerji, "Learning by experimentation: Acquiring and refining problem-solving heuristics," in R. S. Michalski, J. G. Carbonell, and T. M. Mitchell, eds. *Machine Learning*. Palo Alto, CA: Tioga, 1983.
- [37] M. C. Mozer and P. Smolensky, "Using relevance to reduce network size automatically," *Connection Sci.*, vol. 1, no. 1, pp. 3-16, 1989.
- [38] A. Newell, J. C. Shaw, and H. A. Simon, "Programming the logic theory machine," in *Proc. 1957 Western Joint Computer Conf.*, 1957, pp. 230-240.
- [39] D. Nguyen and B. Widrow, "The truck backer-upper: An example of self-learning in neural networks," in *Proc. Int. Joint Conf. on Neural Networks*, vol. II, 1989, pp. 357-363.
- [40] J. Pearl, "On the discovery and generation of certain heuristics," *AI Mag.*, pp. 23-33, Winter/Spring 1983.
- [41] ———, *Heuristics—Intelligent Search Strategies for Computer Problem Solving*. Reading, MA: Addison-Wesley, 1984.
- [42] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, "Learning internal representations by error propagation," in *Parallel Distributed Processing: Explorations in the Microstructure of Cognition*, D. E. Rumelhart

- and J. L McClelland, eds. Cambridge, MA: MIT, vol. 1, 1986, pp. 318–362.
- [43] S. Russell and E. Wefald, "Principles of metareasoning," *Artif. Intell.*, vol. 49, pp. 361–395, 1991.
- [44] A. L. Samuel, "Some studies in machine learning using the game of checkers," *IBM J. Res. Develop.*, vol. 3, pp. 210–229, 1959.
- [45] ———, "Some studies in machine learning using the game of checkers II—Recent progress," *IBM J. Res. Develop.*, vol. 11, no. 6, pp. 601–617, 1967.
- [46] S. R. Schwartz, "Resource constrained parameter tuning applied to stereo vision," M.Sc. thesis, Dep. Elect. Comput. Eng., Univ. of Illinois, Urbana, IL, Aug. 1991.
- [47] J. Sietsma and R. J. F. Dow, "Neural net pruning—Why and how?", in *Proc. Int. Joint Conf. on Neural Networks*, vol. 1, 1988, pp. 325–333.
- [48] D. Sleeman, P. Langley, and T. M. Mitchell, "Learning from solution paths: An approach to the credit assignment problem," *AI Mag.*, vol. 3, 1982, pp. 48–52.
- [49] G. Smith, "Back propagation with dynamic topology and simple activation functions," Tech. Rep. TR 90-12, School Inform. Sci. Technol., Flinders Univ. of South Australia, Adelaide, 1990.
- [50] R. G. Smith, T. M. Mitchell, R. A. Chestek, and B. G. Buchanan, "A model for learning systems," in *Proc. 5th Int. Joint Conf. on Artificial Intelligence*, Aug. 1977, pp. 338–343.
- [51] S. F. Smith, "Flexible learning of problem solving heuristics through adaptive search," in *Proc. Int. Joint Conf. on Artificial Intelligence*, 1983, pp. 422–425.
- [52] C. Stein, "The selection of the largest of a number of means," *Ann. Math. Statist.*, vol. 19, pp. 429, 1948.
- [53] R. S. Sutton, "Temporal credit assignment in reinforcement learning," Ph.D. dissertation, Univ. of Massachusetts, Amherst, MA, Feb. 1984.
- [54] R. S. Sutton and A. G. Barto, "Toward a modern theory of adaptive networks: Expectation and prediction," *Psychological Rev.*, vol. 88, no. 2, pp. 135–170, 1984.
- [55] B. G. Tabachnick and L. S. Fidell, *Using Multivariate Statistics*, second ed. New York: Harper and Row, 1989.
- [56] G. Tesauro, "Connectionist learning of expert backgammon evaluations," *Mach. Learning*, pp. 200–206, 1988.
- [57] Y. L. Tong and D. E. Wetzel, "Allocation of observations for selecting the best normal population," in *Design of Experiments: Ranking and Selection*, T. J. Santner and A. C. Tamhane, eds. New York: Marcel Dekker, 1984, pp. 213–224.
- [58] B. W. Wah and H. Kriplani, "Resource constrained design of artificial neural networks," in *Proc. Int. Joint Conf. on Neural Networks*, IEEE, San Diego, CA, June 1990, pp. 269–279.
- [59] E. H. Wefald and S.J. Russell, "Adaptive learning of decision-theoretic search control knowledge," *Mach. Learning*, pp. 408–411, 1989.
- [60] B. Widrow, N. K. Gupta, and S. Maitra, "Punish/reward: Learning with a critic in adaptive threshold systems," *IEEE Trans. Syst., Man, Cybern.*, vol. SMC-3, pp. 455–465, 1973.
- [61] D. C. Wilkins, "Knowledge base refinement using apprenticeship learning techniques," in *Proc. Nat. Conf. on Artificial Intelligence AAAI-88*, 1988.
- [62] R. J. Williams, "On the use of backpropagation in associative reinforcement learning," in *Proc. Int. Conf. on Neural Networks*, vol. I, July 1988, pp. 263–270.
- [63] S. W. Wilson, "Hierarchical credit allocation in classifier systems," in *Genetic Algorithms and Simulated Annealing, Research Notes in Artificial Intelligence*, L. Davis, ed. London, U.K.: Pitman, 1987.
- [64] J. C. Yan and S. F. Lundstrom, "The post-game analysis framework—Developing resource management strategies for concurrent systems," *IEEE Trans. Knowl. Data Eng.*, vol. 1, pp. 293–309, Sept. 1989.
- [65] C. F. Yu and B. W. Wah, "Learning dominance relations in combinatorial search problems," *IEEE Trans. Software Eng.*, vol. SE-14, pp. 1155–1175, Aug. 1988.
- [66] S. Zhou, "Performance studies of dynamic load balancing in distributed systems," Tech. Rep. UCB/CSD 87/376, Comput. Sci. Div., Univ. of California, Berkeley, CA, 1987.



Benjamin W. Wah (S'78-M'79-SM'85-F'91) received the Ph.D. degree in engineering from the University of California, Berkeley, CA, in 1979.

He is currently a Professor with the Department of Electrical and Computer Engineering and the Coordinated Science Laboratory, University of Illinois at Urbana-Champaign. Previously, he has served on the faculty of Purdue University, as a Program Director with the National Science Foundation, and as the Fujitsu Visiting Chair Professor of Intelligence Engineering, University of Tokyo. He

currently serves as the Associate Editor-in-Chief of the IEEE TRANSACTIONS ON KNOWLEDGE AND DATA ENGINEERING and will become the Editor-in-Chief in 1993. He also serves as an Area Editor of the *Journal of Parallel and Distributed Computing* and as an Editor for *Information Sciences* and the *International Journal on Artificial Intelligence Tools*. His major current research interests include areas which relate to parallel and distributed processing and artificial intelligence.

Dr. Wah has chaired a number of conferences, including the IEEE International Conference on Data Engineering, the International Conference on Parallel Processing, and the IEEE International Conference on Distributed Computing Systems. He currently serves on the Governing Board, Publications Board, Press Board, and Nominations Committee of the Computer Society.