

# Generalization and Generalizability Measures

Benjamin W. Wah, *Fellow, IEEE*

**Abstract**—In this paper, we define the generalization problem, summarize various approaches in generalization, identify the credit assignment problem, and present the problem and some solutions in measuring generalizability. We discuss anomalies in the ordering of hypotheses in a subdomain when performance is normalized and averaged, and show conditions under which anomalies can be eliminated. To generalize performance across subdomains, we present a measure called probability of win that measures the probability whether one hypothesis is better than another. Finally, we discuss some limitations in using probabilities of win and illustrate their application in finding new parameter values for TimberWolf, a package for VLSI cell placement and routing.

**Index Terms**—Anomalies in generalization, credit assignment problem generalization, machine learning, subdomains, probability of win, VLSI cell placement and routing.

## 1 INTRODUCTION

GENERALIZATION in psychology is the tendency to respond in the same way to different but similar stimuli [6]. Such transfer of tendency may be based on temporal stimuli, spatial cues, or other physical characteristics. Learning, on the other hand, may be considered as a balance between generalization and discrimination (the ability to respond to differences among stimuli). An imbalance between them may lead to negative results.

*Machine learning* in an area in artificial intelligence that extends knowledge, concepts, and understanding through one or more observations of instances of the concept [1]. The number of instances involved and the amount of information they carry will determine the learning method to be used.

Learning methods can be classified as data-intensive and knowledge-intensive (see Fig. 1). In *data-intensive methods*, symbolic concepts are learned using data-intensive similarity-based methods. The learner is shown a large number of related examples and is required to identify their similarities and generalize the concept embedded. Using this approach, Mitchell [25] defines *generalization* as a process that takes into account a large number of specific observations (inductive bias), and that extracts and retains the important features that characterize classes of these observations. He then casts generalization as a search problem, and alternative generalization methods as different search strategies.

An example of a data-intensive learning method is the learning of heuristics represented as a collection of production rules [42]. In this approach, learning modifies each of the rules based on decisions made by these rules and on positive and negative examples found. The process of 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 heuristic method is called

*credit assignment*. The former credit assignment is called *temporal*, and the latter, *structural*. Credit assignment is usually difficult when learning incrementally single concepts from examples, especially when learning multiple disjunctive concepts and when the learning data is noisy. In this case, a teacher may be needed to tell the learner the proper amount of credit to assign to a decision.

A second class of data-intensive learning methods are *decision-theoretic methods* that use statistical decision theory to discriminate probabilistic patterns exhibited in learning examples [28]. The major component in a decision-theoretic approach is the *loss function* that measures the loss when the learner categorizes a learning example incorrectly. It represents a statistical approach to credit assignment. By minimizing the total loss using statistical methods, it is sometimes possible to show asymptotic convergence of the concept to be learned. Examples of decision-theoretic methods include evolutionary programming [20], genetic algorithms [13], classifier systems [5], and artificial neural networks (ANNs) [22].

In contrast to using extensive training examples in data-intensive methods, *knowledge-intensive methods* rely on domain-specific knowledge to learn and to generalize. In *explanation-based learning*, the learner analyzes a single training example using domain knowledge and the concept under study to produce a generalization of the example and a deductive justification of the generalization [8], [26]. Knowledge-intensive methods work well when the concept to be generalized can be deduced from the domain knowledge.

To evaluate the quality of a learning and generalization method and to measure the degree to which learning and generalization has been achieved, *generalizability* measures have been developed. In the simplest case, they measure the number of positive and negative examples in learning. In more general cases, the degree to which an example satisfies a learned concept must be considered, and statistical techniques are employed to determine whether a learned concept can be generalized.

For example, in learning in feedforward ANNs, the effectiveness of an ANN that computes discrete  $\{0, 1\}$ -valued

• B.W. Wah is with the Department of Electrical and Computer Engineering and the Coordinated Science Laboratory, University of Illinois at Urbana-Champaign, Urbana, IL 61801. E-mail: b-wah@uiuc.edu. URL: <http://manip.crhc.uiuc.edu>. He is currently on leave at the Department of Computer Science and Engineering, Chinese University of Hong Kong, Shatin, Hong Kong.

Manuscript received 13 Sept. 1997; revised 20 Aug. 1998.  
For information on obtaining reprints of this article, please send e-mail to: [tkde@computer.org](mailto:tkde@computer.org), and reference IEEECS Log Number 108305.

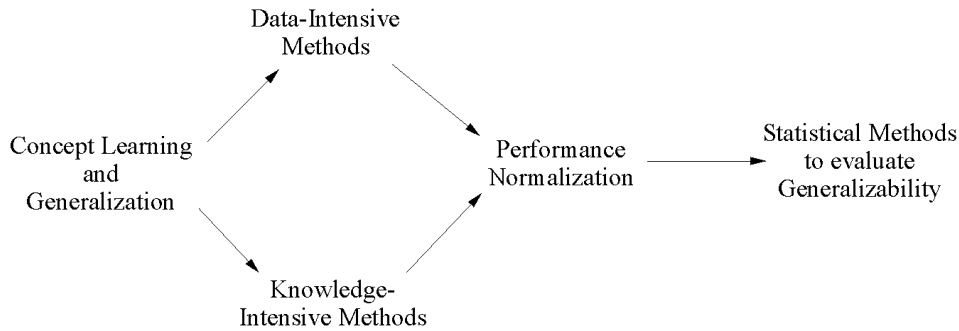


Fig. 1. The relationship between concept learning, generalization, and generalizability.

mappings can be evaluated by the network's ability to solve dichotomization problems using measures such as discrimination capacity, VC-dimension (named after Vapnik and Chervonenkis [36]), and efficiency of decision functions [23]. For an ANN that performs function approximation computing either discrete multiple-valued or continuous mappings, we can measure its quality using concepts such as combinatorial dimension, approximation error, and estimation error. Finally, the concept of *PAC* (probably approximately correct) learning [18] is useful for characterizing the time complexity of algorithms for learning both discrete and continuous mappings.

A related problem in generalizability is the normalization of learned results relative to a baseline. When the quality of a learned concept is measured numerically and depends on some attributes of the example, it may be necessary to normalize the measure with respect to that of a baseline before any statistical evaluations can be made. For instance, the quality measure of a learned concept may depend on the size of the learning example and needs to be normalized before results from multiple learning examples can be aggregated statistically. In this case, the generalizability of the learned concept may depend on the baseline and the statistical method used to aggregate performance measures. Anomalies in the ordering of hypotheses may happen when different normalization and aggregation methods are used. This is discussed in detail in Section 3.3.

In the next section, we summarize previous approaches in generalization and credit assignment. We then present, in Section 3, the general concept of generalizability, generalizability measures, and anomalies in generalization when performance measures are normalized and aggregated. We illustrate in Section 4 the application of the method in Section 3 to find new parameter values in TimberWolf.

## 2 CONCEPT GENERALIZATION USING INDUCTION

In this section, we summarize various strategies for generalization. Early work on inductive learning and generalization was done by Simon and Lea [31] who used training instances selected from some space of possible instances to guide the search for general rules. The process of inductive learning entails a mapping from the instance space to the rule space and involves experiment planning, instance selection, and result interpretation (or generalization). (See Fig. 2.)

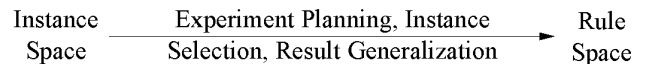


Fig. 2. The process of inductive learning and generalization.

### 2.1 The Generalization Problem

Generalization involves the extraction of information useful to guide the search of a rule space [1]. To simplify the search process, a good representation of the rule space must be chosen so that generalization can be carried out by inexpensive syntactic operations, such as turning constants to variables, dropping conditions, adding options, curve fitting and zeroing a coefficient.

The specific operators used may depend on the representation of the rule space. For instance, a production rule  $Z \rightarrow Z'$  can be used to represent either the backward form ( $Z'$  is the value of a state vector plus associated predicate) or the forward form ( $Z'$  is a computational rule). The evaluation of the execution of a rule constitutes credit assignment, whereas the creation of new rules involves generalization. The latter entails the identification of a subvector of variables relevant to the creation, the proper decision for the situation, and the reason for making the decision. Waterman [42] proposed a set of generalization operators that modify the defined symbolic values in a rule, eliminate one or more variables in a rule, and change action rules and error-causing rules.

Mitchell [25] defines generalization in the context of a language that describes instances and generalizations. Based on a set of positive and negative examples, predicates are matched from generalizations to instances. Hence, generalizations are defined within the provided language that are consistent with the presented training examples.

In general, generalization also requires a function to evaluate the positive and negative examples obtained in order to provide feedback (credit assignment). In the simplest case, the function counts the number of positive and negative examples. In decision-theoretic approaches, a loss function is used to measure the loss when the learner categorizes a learning example incorrectly. This is the approach taken in classifier-system and genetics-based learning that uses a fitness function. In reinforcement learning, the evaluation function may have to be learned independently (by some form of supervised learning) in order to provide proper temporal credit assignment [33], [27], [34].

The reinforcement function is particularly difficult to design when examples drawn from the problem space are not statistically related. This happens when the evaluation data depends on the size of the examples, or when the examples drawn belong to different problem subdomains. Some possible solutions to these issues are discussed in Section 3.3.

### 2.2 Generalization Strategies

As defined by Mitchell [25] and Mitchell et al. [26], generalization strategies can broadly be classified as data driven and knowledge-driven. (See Fig. 3.) Both paradigms use generate-and-test that generates alternative concepts, tests them on test cases, and constructs feedbacks (credit assignment) to aid the refinement of the concepts generated. The difference lies in the amount of tests performed: Data-driven methods do not rely on domain knowledge and often require extensive tests on the concepts under consideration before reliable feedbacks can be generated. In contrast, knowledge-driven methods rely on domain knowledge and one or a few tests to deduce new concepts.

Data-driven generalization strategies can be classified into depth-first search, breadth-first search, version-space, and decision-theoretic techniques [25].

A *depth-first strategy* starts from a single generalization as the current best hypothesis, tests it against each training example, and modifies the hypothesis in order to make it consistent with the training example. Its advantage is that it keeps a global picture in mind when modifying the hypothesis. However, it is usually expensive to backtrack when a negative training example is found. In this case, the new hypothesis generated must be tested against all previous training examples to make sure that they are consistent. Any inconsistencies will incur further backtracking.

A *breadth-first strategy*, on the other hand, generalizes from more specific hypotheses to more general ones. Initially, it starts from a set of the most specific hypotheses. Positive training examples allow the search to progress down the breadth-first tree, generating more general hypotheses, whereas negative training examples will prune the corresponding hypothesis from the search tree. The boundary of the search tree, therefore, represents the most general hypotheses generated so far that are consistent with the (positive) training examples. As a result, when a new (more general) hypothesis is generated, it only needs to be tested against all positive training examples to make sure that they are consistent with the current hypothesis. This is

the main advantage of a breadth-first search over a depth-first search.

A hybrid of depth-first and breadth-first strategies is a *version-space strategy*. The version space represents the set of hypothesis that are consistent with all the training examples. It defines two boundaries. The first boundary is obtained by depth-first search and bounds the acceptable level of specialization of hypotheses (those that are consistent with all the positive examples). The second boundary is obtained by breadth-first search and bounds the acceptable level of generality of hypotheses (those that are inconsistent with all the negative examples).

The fourth class of data-driven generalization strategies are the decision-theoretic techniques. These do not always use a single type of search method but may use a hybrid of search methods, such as depth-first and breadth-first searches, depending on the evaluation results. They rely on a loss function that measures the expected loss when the learner categorizes a learning example incorrectly. Although the loss function may be designed either based on formal statistical methods or heuristically, the generalization strategy can generally be shown to converge asymptotically to the desired concept. For instance, in genetic algorithms, Holland's Schema Theorem [13] shows that the number of structures in a knowledge base that share a given subset of components can be expected to increase or decrease over time at a rate proportional to the observed performance of the subset, eventually converging asymptotically to the optimal configuration.

In contrast to data-driven techniques, an explanation-based generalization strategy uses domain knowledge to generalize from an example, defining a concept that contains the example [26]. It analyzes a single example in terms of the domain knowledge and the goal concept and produces a proof (or explanation) that shows that the example is an instance of the goal concept. Here, the goal concept found satisfies the operability criteria, which is a predicate over concept definitions that specifies the form in which the concept must be learned. The proof tree in the process of generalization is constructed by replacing each instantiated rule by the associated general rule.

An explanation-based strategy can start from general concepts to derive specific ones, or vice versa. It consists of two phases: explanation and generalization. In the explanation phase, the relevant features of the training example are isolated in order to create an explanation structure that

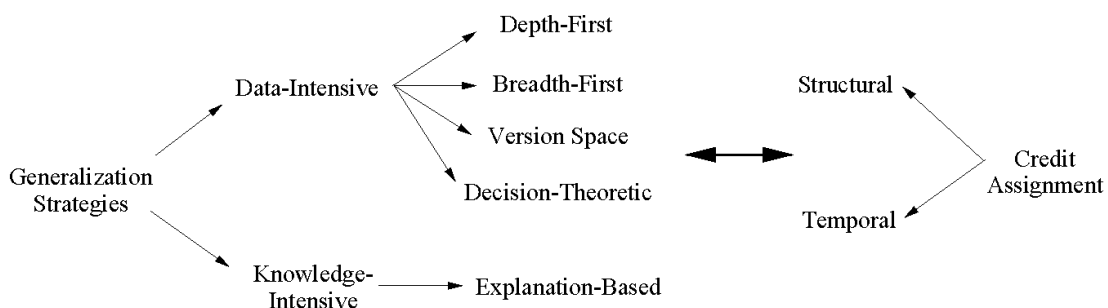


Fig. 3. A classification of generalization strategies.

terminates in an expression satisfying the operability criterion. In the generalization phase, a set of sufficient conditions are found to satisfy the explanation. This is done by regressing the goal concept through the explanation structure and by composing terms from different parts of the explanation to form a valid generalization.

One of the problems in explanation-based learning is that learning through multiple examples may result in multiple rules that cannot be combined into a single rule. This leads to gradual degradation in efficiency in the generalized rules. Another problem is that explanation-based generalization does not create new parameters; hence, parameters not explicitly expressed in the proof cannot be generalized. In this context, studies have been made to generalize the structure of an example proof such that a fixed number of rule applications in the proof is generalized into an unbounded number of applications [7].

### 2.3 Credit Assignment

Credit assignment entails the apportioning of feedback signals to individual decisions made in the past as well as rules/entities leading to a decision. The application of credit assignment requires a *world model* that captures the relationship among states, decisions, and feedback signals generated by the learning system or measured in the environment. This world model is explicitly defined in knowledge-rich applications, but may have to be inferred during learning and generalization when domain knowledge is not available. Credit assignment 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 times a decision was made and its feedback signal received.

There are two types of credit assignment: structural and temporal [33]. *Structural credit assignment* entails ways of using feedback signals to refine the individual components or rules of a hypothesis. This process is systematic in explanation-based learning as it involves rewriting one rule into another in the proof structure. In other learning approaches, credit assignment may not be possible when domain knowledge is missing. In this case, one may employ population-based learning [37] that maintains a population of competing hypotheses and delays choosing the best hypothesis to evaluate or new ones to create until more tests are performed on each of the alternatives.

*Temporal credit assignment*, on the other hand, entails the apportioning of temporal global feedback signals by the learning system to the 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. Temporal credit assignment is used extensively in reinforcement learning [33].

Credit assignment can be either implicit or explicit. An example of implicit credit assignment is done in Smith's LS-1 system [32] in which rules that are physically close together on the list representing the knowledge structure stand a good chance of being inherited as a group. On the other hand, in explicit credit assignment, explicit rules are defined for credit assignment. Examples of explicit temporal credit-assignment mechanisms are the profit sharing plan and the bucket brigade algorithm in classifier systems [14]. A hybrid of implicit and explicit credit assignment can also be defined [11].

## 3 GENERALIZABILITY MEASURES

To evaluate whether the goal of learning and generalization is achieved, generalizability measures are used to evaluate the quality of generalization. These measures are not limited to the field of machine learning but are used in performance evaluation of many other areas. For instance, in evaluating the speed of a computer, one generally defines a reference computer, such as the VAX 11/780, and computes the speedup of the computer with respect to the reference for a collection of benchmarks. Based on the evaluation results, one generalizes the speedup to benchmarks not tested in the evaluation.

Since different regions of the problem space of an application domain may have different characteristics, it may not be possible to evaluate generalization across all examples of a problem space. To this end, the problem space is decomposed into smaller partitions before generalization is evaluated. For instance, in evaluating a computer, one defines its speedups for different collections of benchmarks in order to reflect its performance under different applications.

In the partitioning of a problem space, we define a *problem subspace* as a user-defined partition of a problem space so that concepts/hypotheses for one subspace are evaluated independent of concepts/hypotheses in other subspaces. Such partitioning is generally guided by common-sense knowledge or by user experience in solving similar application problems. To identify a problem subspace, we need to know one or more attributes to classify test cases and a set of decision rules to identify the subspace to which a test case belongs. For instance, in evaluating the speedup of a computer, the partitioning of the class of all applications is guided by user experience into the class of scientific applications and the class of business applications.

Given a subspace of test cases, we define a *problem subdomain* as a partitioning of the subspace into smaller partitions so that the evaluation of a concept/hypothesis can be done quantitatively for all the test cases in a subdomain. Such partitioning is necessary because the statistical performance metrics computed (such as average or maximum) is not meaningful when the performance values are of different ranges and distributions. To continue from the previous example, the class of scientific benchmarks are further partitioned into subdomains according to their computational behavior, such as whether a program is CPU-bound or I/O-bound.

In the same way that test cases are partitioned into subspaces, we need to know the attributes to classify test cases

and a set of decision rules to identify the subdomain to which a test case belongs. This may be difficult in some applications because the available attributes may not be well defined or may be too large to be useful. For instance, the attribute to classify whether a benchmark program is CPU-bound or I/O-bound is imprecise and may depend on many underlying characteristics of the program.

After evaluating the performance of a hypothesis in each subdomain, we need to compare its performance across subdomains. For instance, one would be interested to know whether a computer has high speedups across both CPU-bound and I/O-bound applications. This comparison may be difficult because test cases in different subdomains of a subspace may have different performance distributions and cannot be compared statistically. We address this issue in Section 3.3. In the next subsection, we examine some formal results in generalizability for classification problems in one subdomain.

### 3.1 Formal Results on Learnability and Generalizability

Formal methods to deal with generalizability in learning with one performance measure have been studied extensively in computational learning theory. They center on the notion of PAC-learnability [35] of a concept class  $\mathbf{C}$  by a learning algorithm  $L$ , where a concept is defined as a subset of some instance space  $X$ . A *learner* tries to learn target concept  $C$ , finding out points of  $X$  (drawn randomly) whether they belong to the target concept. The goal of the learner is to produce with high probability ( $>1 - \delta$ ) a hypothesis that is close (within  $\epsilon$ ) to the target concept, assuming that the learner does not know the underlying distribution of the sample points. (The following definitions are from a survey paper by Kearns et al. [17].) A concept  $C$  produced by a learning algorithm  $L$  on input vector  $T$  is *approximately correct* if the error rate  $P(C \oplus T)$  is at most  $\epsilon$ . If, for any concept class  $\mathbf{C}$ , with probability distribution  $P$ , accuracy parameter  $\epsilon$ , and confidence parameter  $\delta$ , the probability that the output  $C$  is approximately correct is at least  $(1 - \delta)$ , then the learning algorithm is *probably approximately correct*; and,  $L$  is said to *PAC-learn*  $\mathbf{C}$ . A learning algorithm  $L$  is a polynomial PAC-learning algorithm for class  $\mathbf{C}$ , if  $L$  PAC-learns  $\mathbf{C}$  with both time complexity and sample complexity polynomial in  $\frac{1}{\epsilon}$  and  $\frac{1}{\delta}$ .

To understand bounds on estimation by a learning algorithm, we need to estimate the largest number of input-space points for which *almost every* possible dichotomy is achieved by some concept from a class  $\mathbf{C}$ . *VC-dimension* (named after Vapnik and Chervonenkis [36]) addresses this issue. VC dimension,  $V$ , of a concept class  $\mathbf{C}$  is the size of the largest set  $\mathbf{S}$  of input-space points such that for every subset  $\mathbf{U} \subseteq \mathbf{S}$ , there exists some concept  $C \in \mathbf{C}$  where  $\mathbf{U} = \mathbf{S} \cap C$ .  $C$  is some function realized by the concept; and  $\mathbf{C}$  is the set of all such functions realizable by that concept.

Sauer [29] notes that whenever the VC dimension of a function class is finite, the number of dichotomies grows subexponentially (actually, polynomially) in the number of points. The probability of a concept learned with a large estimation error producing correct outputs for a given set

of points goes rapidly to zero as the size of the set increases. A learning algorithm whose outputs are always consistent with the examples seen so far is called a *consistent PAC-learning algorithm*. If the VC dimension of a concept class is finite, then a consistent learning algorithm trained on a sufficiently large set of examples is likely to learn the correct concept.

Blumer et al. [4] have derived bounds on the number  $m(\epsilon, \delta)$  of examples needed by a consistent algorithm to PAC-learn a concept class  $\mathbf{C}$  having VC dimension  $d$ . This was improved by Ehrenfeucht et al. [10] to

$$\Omega\left(\frac{1}{\epsilon} \ln \frac{1}{\delta} + \frac{d}{\epsilon}\right).$$

Baum and Haussler [3] have used these results to relate the size of a neural network, the accuracy of the learned concept, and the number of examples needed in order to guarantee a particular degree of accuracy. Their analysis suggests that generalization can be improved by pruning unnecessary hidden units during learning. The reduced architecture has VC dimension not significantly larger than the VC dimension for an optimal number of hidden units. Baum and Haussler establish the following necessary and sufficient conditions for valid generalization for learning in neural networks of thresholded binary units.

- A network of  $N$  nodes and  $W$  weights, which after being trained on at least

$$O\left(\frac{W}{\epsilon} \log \frac{N}{\epsilon}\right)$$

examples, classifies at least  $(1 - \frac{\epsilon}{2})$  of them correctly, will almost certainly classify a fraction  $(1 - \epsilon)$  of future examples correctly.

- A fully connected feedforward network with one hidden layer, trained on fewer than  $\Omega(\frac{W}{\epsilon})$  examples will, for a dichotomy realizable by the network, fail to find the requisite set of weights for more than a fraction  $(1 - \epsilon)$  of future examples.

Haussler [12] shows that, for it to be likely that feedforward networks with sigmoidal units obtain a low estimation error, the number of examples must grow linearly with both the number of modifiable weights and the number of hidden layers. That is, either of the following desiderata demands a larger training sample:

- 1) lowering the estimation error;
- 2) increasing the confidence; and
- 3) learning with sigmoids having a higher slope.

Barron [2] shows that, for a feedforward network having  $n$  sigmoidal units and  $d$  input units and trained on  $N$  examples, the total mean squared error (approximation plus estimation) between the true function and the estimated function is bounded from above by

$$O\left(\frac{1}{n}\right) + O\left(\frac{nd}{N}\right) \log N.$$

In summary, the theory in learnability provides conditions and bounds on generalization that are useful when

certain restricted assumptions are met. Such assumptions may be difficult to ascertain in practice because it is difficult to characterize the set of test cases and hypotheses precisely. Under such conditions, heuristic methods to measure generalizability need to be developed. In the next two subsections, we present some results in this area.

### 3.2 Anomalies in Performance Normalization

In general learning problems, the raw performance results obtained in evaluating hypotheses on examples may depend on the size and characteristics of the examples and may not be directly comparable. For instance, Table 1 shows the CPU times of four computers in evaluating three benchmarks. Obviously, these performance values cannot be aggregated directly because they belong to different ranges and are of different distributions. To aggregate them statistically, we must normalize them first. In the following, we show five different normalization methods.

TABLE 1  
SUMMARY OF RAW CPU TIMES OF FOUR COMPUTERS  
IN EVALUATING THREE BENCHMARKS

Benchmark	Computer			
	$C_{75}$	$C_{76}$	$C_{86}$	$C_{99}$
$t_1$	30.19	30.31	26.21	40.61
$t_2$	43.12	43.34	34.09	24.65
$t_3$	71.93	72.49	104.51	98.41

- 1) *Average improvement ratio.* Using the performance values of one hypothesis as the baseline, we normalize each performance value of another hypothesis by computing its ratio with respect to that of the baseline when tested on the same example. The average of the improvement ratios is then used as the aggregate performance measure. The drawback of this approach is that different ordering of the hypotheses can be obtained, depending on the baseline hypothesis used. To illustrate this point, consider the performance data presented in Table 1. The second column of Table 2 shows the three anomalous orderings of the four computers based on their average normalized speedups using each computer as the baseline for normalization. This shows that generalization based on the average improvement ratios does not always lead to consistent conclusions.
- 2) *Average symmetric improvement ratio.* This is a normalization method we have developed before to avoid anomalies in inconsistent orderings of two

hypotheses due to the choice of the baseline hypothesis [39]. The idea is to avoid the problem in improvement ratios that put different weights in different ranges of normalized performance values. Note that degradations in the original improvement ratio are between zero and one, whereas improvements are between one and infinity. Consequently, when improvement ratios are averaged, degradations carry less weight than improvements.

The symmetric improvement ratio is defined as follows:

$$S_{sym+,i} = \begin{cases} S_{+,i} - 1 & \text{if } S_{+,i} \geq 1 \\ 1 - \frac{1}{S_{+,i}} & \text{if } 0 \leq S_{+,i} < 1 \end{cases} \quad (1)$$

$$\bar{S}_{sym+} = \frac{1}{m} \sum_{i=1}^m S_{sym+,i} \quad (2)$$

where  $S_{+,i}$  is the original improvement ratio on the  $i$ th test case. The symmetric improvement ratio has the property that improvements are in the range between 0 and infinity, and degradations are in the range between 0 and negative infinity. For two hypotheses, when we reverse the role of the baseline hypotheses, their symmetric improvement ratios only change in sign. Hence, symmetric improvement ratios avoid anomalies in performance orderings with two hypotheses.

However, anomalies in performance ordering are still present when more than two hypotheses are concerned. This is illustrated in Table 2 that shows three different orderings when different computers are used as the baseline. Hence, generalization based on the average symmetric improvement ratios may not lead to consistent conclusions.

- 3) *Harmonic mean performance.* This is defined as follows:

$$\bar{S}_h = \frac{m}{\sum_{i=1}^m 1/S_{+,i}} \quad (3)$$

Again, as illustrated in Table 2, anomalies in orderings are still present.

- 4) *Geometric mean performance.* This is defined as follows:

$$\bar{S}_g = \sqrt[m]{\prod_{i=1}^m S_{+,i}} \quad (4)$$

TABLE 2  
ANOMALOUS ORDERINGS OF COMPUTERS IN DECREASING AVERAGE NORMALIZED  
SPEEDUPS USING THREE DIFFERENT NORMALIZATION METHODS

Baseline	Average Improvement Ratio	Average Symmetric Improvement Ratio	Harmonic Mean
$C_{75}$	$C_{99}C_{86}C_{75}C_{76}$	$C_{99}C_{75}C_{76}C_{86}$	$C_{75}C_{76}C_{86}C_{99}$
$C_{76}$	$C_{99}C_{86}C_{75}C_{76}$	$C_{99}C_{75}C_{76}C_{86}$	$C_{75}C_{76}C_{86}C_{99}$
$C_{86}$	$C_{75}C_{76}C_{99}C_{86}$	$C_{75}C_{76}C_{85}C_{99}$	$C_{86}C_{75}C_{76}C_{99}$
$C_{99}$	$C_{75}C_{76}C_{86}C_{99}$	$C_{86}C_{99}C_{75}C_{76}$	$C_{99}C_{86}C_{75}C_{76}$

Taking the logarithm of both sides, we have:

$$\log \bar{S}_g = \frac{1}{m} \sum_{i=1}^m \log S_{+,i} = \frac{1}{m} \sum_{k=1}^m \log t_{b,k} - \frac{1}{m} \sum_{k=1}^m \log t_{h,k} \quad (5)$$

where  $t_{b,k}$  and  $t_{h,k}$  are, respectively, the  $k$ th performance values of the baseline and the hypothesis being normalized. Based on (5), an alternative way to view a geometric mean is that it is an arithmetic mean of the logarithms of raw performance values. The effect of the baseline hypothesis on the average normalized performance is reflected in the first constant term in (5). Hence, when the baseline is changed, only a constant term will be changed, and performance ordering is not affected. This is illustrated in the example in Table 1 in which the ordering  $C_{86}C_{75}C_{76}C_{99}$  is unchanged when the baseline is changed.

- 5) *Average normalized performance with respect to the median performance.* This belongs to a general class of methods that normalizes the performance values of hypotheses on each test case with respect to a test case-specific constant that is invariant as more hypotheses are evaluated. The specific method here uses the *median* performance value of all the hypotheses on each test case as the baseline for normalization. Unlike using a baseline hypothesis that may induce a different ordering when the baseline is changed, the median performance is invariant with respect to the hypotheses and test cases in a subdomain. Using this normalization method, the performance distributions of all the test cases will center around zero. This method is illustrated in the example in Table 1 in which the ordering is  $C_{76}C_{75}C_{99}C_{86}$ . In computing this ordering, we made a simplifying assumption that the median performance of each computer on the three benchmarks is the same as the median performance of the computer across all possible benchmarks.

A potential problem with this approach is the unavailability of the true median performance value of hypotheses for each test case. Hence, the sample median may have to be used instead. Unfortunately, estimated sample medians are inaccurate during learning because hypotheses may not be tested adequately, and sample medians are sensitive to the hypotheses tested. Solutions to this issue are still open at this time.

In summary, anomalies in performance normalization do not exist when either the baseline is fixed (as in the case of the median performance) or the effect of changing the baseline only results in changing a constant term in the (transformed) normalized performance (as in the case of the geometric mean performance). In other cases, it is possible for the order of the hypotheses to change when the baseline is changed. The necessary and sufficient conditions for anomalies to happen are still open at this time.

### 3.3 Generalizability Measures Across Subdomains

When hypotheses are tested across different subdomains of an application, their performance values, even after normalization, may have different ranges and different distributions. As a result, these performance values cannot be aggregated statistically, and the hypotheses cannot be compared directly and generalized across subdomains. In this section, we present a heuristic method to evaluate performance across subdomains in a range-independent way. We assume that the performance values of testing a hypothesis in a subdomain are independent and identically distributed. This assumption allows the values in a subdomain to be aggregated by statistical methods, such as averaging.

In the following, we present a method that uses the sample mean as a statistical estimate of the population mean. To address uncertainties in using sample means, we have studied a concept called *probability of win* [39],  $P_{win}$ , that compares two sample means and computes the probability that one sample mean is larger than another. This is similar to hypothesis testing in which we take random samples to test whether a property of a population is likely to be true or false [15]. Obviously, it may be difficult to test a hypothesis fully by testing the entire population of test cases or by testing only a single random sample.

There are four steps in general hypothesis testing:

- 1) Specify a significance level  $\alpha$ .
- 2) Specify the testing hypotheses that include both null hypothesis  $H_0$  and alternative hypothesis  $H_1$ .
- 3) Find the corresponding acceptance region using lookup tables.
- 4) Make decision on the sample value. If the sample falls in the acceptance region, then accept  $H_0$  and reject  $H_1$ ; otherwise, reject  $H_0$  and accept  $H_1$ .

The probability of win measures statistically how much better (or worse) the sample mean of one hypothesis is as compared to that of another. It resembles the significance level in general hypothesis testing, but there are two major differences. First, only one hypothesis  $\{H: \mu_1 > \mu_2\}$  is specified, without the alternative hypothesis. Further, in contrast to hypothesis testing, acceptance confidence is not given in advance but is evaluated based on sample values.

One advantage of  $P_{win}$  is that it is between zero and one and is independent of the actual performance difference across subdomains. Hence, it can be used to compare hypotheses in a uniform way across subdomains.

Consider the performance of  $H_i$  in subdomain  $j$ . (For convenience of formulation, subscript  $j$  is ignored in the following discussion.) Let  $\mu_i$  and  $\sigma_i$  be the true mean and true standard deviation of the mean normalized performance with respect to the baseline hypothesis  $H_0$ .<sup>1</sup> When  $n_i$  samples are taken, we can calculate the sample mean  $\bar{\mu}_i$  and sample standard deviation  $\bar{\sigma}_i$ . By Central Limit Theorem,

$$Pr(\bar{\mu}_i | \mu_i, \sigma_i, n_i) \approx \mathcal{N}\left(\mu_i, \frac{\sigma_i^2}{n_i}\right)$$

1. Any one of the normalization methods presented in Section 3.2 can be used.

where  $\mathcal{N}$  is the normal distribution function with mean  $\mu_i$  and standard deviation

$$\sqrt{\frac{\sigma_i^2}{n_i}}.$$

Let  $t$  be

$$t = \frac{\bar{\mu}_i - \mu_i}{\sigma_i / n_i}$$

where  $t$  has Student's  $t$ -distribution with  $n_i - 1$  degrees of freedom when the number of samples is less than 30 and the variance is unknown. The probability that this hypothesis is better than  $H_0$  with mean value zero<sup>2</sup> is

$$\Pr(H \text{ is true}) = \Pr\left(t \in \left(-\infty, \frac{\bar{\mu}_i}{\bar{\sigma}_i / \sqrt{n}}\right)\right) \quad (6)$$

$$= \int_{-\infty}^{\bar{\mu}_i} p(t \text{ is } t\text{-distributed}) dt \quad (7)$$

where the acceptance region of this hypothesis is

$$\left(-\infty, \frac{\bar{\mu}_i}{\bar{\sigma}_i / \sqrt{n}}\right).$$

Note that the right bound of the acceptance region is a random variable that depends on both the sample mean and the sample variance.

EXAMPLE 1. Table 3 illustrates the  $P_{win}$  for three hypotheses.

We see that  $P_{win}$  of  $H_1$  increases toward one when the number of samples increases. ( $H_1$  is better than  $H_0$ .) In contrast,  $P_{win}$  of  $H_2$  reduces to zero when the number of samples is increased. ( $H_2$  is worse than  $H_0$ .) Last,  $P_{win}$  of  $H_3$  reaches the maximum value 1.0, which means  $H_3$  is definitely better than  $H_0$ .

TABLE 3  
EXAMPLES ILLUSTRATING HOW  $P_{win}$  CHANGES  
WITH INCREASING NUMBER OF SAMPLES

$H_i$	$\mu_i$	$\sigma_i$	No. of Samples to Compute $P_{win}$		
			5	10	30
$H_1$	0.336	1.231	0.652	0.725	0.849
$H_2$	-0.129	0.222	0.202	0.097	0.012
$H_3$	0.514	0.456	0.940	0.991	1.000

Performance is normalized with respect to  $H_0$ .

Note that  $P_{win}$  considers both the mean and variance. Hence, when  $P_{win}$  of a hypothesis is close to 0.5, it is not clear whether the hypothesis is better than or worse than the baseline.

Given baseline hypothesis  $H_0$ , we now show  $P_{win}$  of  $H_i$  in subdomain  $j$  with respect to the average performance of  $H_0$ . Assuming sample mean  $\hat{\mu}_{i,j}$ , sample variance  $\hat{\sigma}_{i,j}^2$ , and  $n_{i,j}$  test cases,  $P_{win}$  is defined as follows:

$$P_{win}(i, j) = F_t\left(n_{i,j} - 1, \frac{\hat{\mu}_{i,j}}{\sqrt{\hat{\sigma}_{i,j}^2 / n_{i,j}}}\right) \quad (8)$$

2. If the average normalized performance of  $H_0$  is not 0, then appropriate shifting in the mean value to 0 can be performed.

where  $F_t(\nu, x)$  is the cumulative distribution function of Student's  $t$ -distribution with  $\nu$  degrees of freedom, and  $P_{win}(i, j)$  is the probability that the true performance (population mean) of  $H_i$  in subdomain  $j$  is better than that of  $H_0$ . When  $n_{i,j} \rightarrow \infty$ , we have

$$P_{win}(i, j) \approx \Phi\left(\frac{\hat{\mu}_{i,j}}{\sqrt{\hat{\sigma}_{0,j}^2 / n_{i,j}}}\right) \quad (9)$$

where  $\Phi$  is the standard cumulative normal distribution function [9].

It is important to point out that probabilities of win are used to evaluate whether a hypothesis is better than the baseline and is not meant to rank order all the hypotheses. Hence, when hypothesis are ordered using their probabilities of win and performance is normalized by any method in which the baseline can be changed, anomalies in performance ordering may still happen. As illustrated in [41], this phenomenon happens because not only the mean but the variance of the baseline are important in determining the ordering of the hypotheses. The variance of the performance values places another degree of freedom in the performance ordering, which can change the ordering when the variance changes. Consequently, the ordering may change when a baseline with a small variance is changed to one with a large variance (or vice versa). This is true even for normalization methods that do not have anomalies when hypotheses are ordered by their mean values, such as the geometric mean. In short, anomalies in performance ordering will not happen when hypotheses are ranked by their probabilities of win and when the baseline hypothesis is fixed, such as the case when the median performance of the hypotheses is used as the baseline.

We are now ready to define a generalizability measure across multiple subdomains. Because different subdomains have different statistical behavior, performance from different subdomains must be treated independently and cannot be combined.

There are two assumptions on the strategies presented here.

- We assume that the set of subdomains used in the design process are *representatives* of all the subdomains in the application. These subdomains behave in a statistically similar fashion to subdomains used in learning and in generalization.
- We assume that the relative importance of one subdomain as compared to another is unknown, and that the performance of hypotheses in subdomains may be dependent. Under these assumptions, we cannot aggregate performance values of hypotheses across subdomains. Our strategy is to select hypotheses so that their worst-case performance across all subdomains is better than a minimum level.

The objective of generalization here is to select a hypothesis that is better than the incumbent hypothesis over a problem domain. When there are multiple such hypotheses, our procedure should attempt to maximize the likelihood of selecting the best hypothesis among the given set. Define:

$$P_{WIN}(i) = \min_j P_{win}(i, j) \quad (10)$$



When there is a baseline hypothesis  $H_0$ , we apply one of the strategies in Section 3.2 to normalize the performance of a hypothesis in a subdomain with respect to the baseline hypothesis. We consider  $H_i$  to be better than  $H_0$  in subdomain  $j$  when  $P_{WIN}(i) > 0.5 + \Delta$ . Note that  $P_{WIN}(i)$  is independent of subdomain  $j$  and can be used in generalization if it were true across all subdomains, even those subdomains that were not tested in learning.

The following are three possible outcomes when comparing  $P_{WIN}(i)$  of  $H_i$  to  $H_0$ :

- 1)  $H_i$  is the only hypothesis that is better than  $H_0$  in all subdomains.  $H_i$  can then be chosen as the hypothesis for generalization.
- 2) Multiple hypotheses are better than  $H_0$  in all subdomains. Here, we should select one hypothesis that maximizes the likelihood of being better than  $H_0$  over the entire domain. This likelihood (or degree of confidence) can be adjusted by increasing  $\Delta$ , which is equivalent to placing a tighter constraint in each subdomain, hence eliminating some potential hypotheses that are found to be better than  $H_0$  under a looser constraint.
- 3) No hypothesis is better than  $H_0$  in all subdomains. Since no hypothesis is superior to  $H_0$ ,  $H_0$  is the most generalizable.

Alternatively, it is possible to find hypotheses such that  $P_{WIN} > 0.5 + \Delta$  where  $\Delta < 0$ . Such hypotheses have less certainty in performing better than  $H_0$  across all the subdomains. However, since  $P_{WIN}$  is based on the worst-case  $P_{win}$  across all the subdomains, hypotheses selected this way may still perform better the baseline in some subdomains. Such hypotheses should be considered as alternatives to  $H_0$ .

We have considered so far generalization based on one performance measure. In general, there may be multiple performance measures in an application, and generalization determines whether a hypothesis behaves consistently across all subdomains with respect to all the performance measures. The problem belongs to a general class of multiobjective optimization problems that can be solved in some special forms. In our approach, we propose to constrain all but one measures and to optimize the unconstrained measure subject to the constraints [39]. The constraints in this approach are defined with respect to the performance of an existing baseline hypothesis. This is similar to first normalizing the performance with respect to that of the baseline and formulating a constraint such that the normalized performance is larger than one. Again, care must be taken in normalization because anomalies in performance ordering may happen when certain normalization methods are used and the baseline is changed.

Probabilities of mean have been used to evaluate generalizability in various genetics-based learning and generalization experiments [39], [15], [34], [24], [41], [40], [16], [38]. These include the learning of load balancing strategies in distributed systems and multicomputers, the tuning of parameters in VLSI cell placement and routing, the tuning of fitness functions in genetics-based VLSI circuit testing, the automated design of feedforward neural networks, the design of heuristics in branch-and-bound search, range

estimation in stereo vision, and the learning of parameters for blind equalization in signal processing.

#### 4 EXAMPLE: VLSI PLACEMENT AND ROUTING

In this section we illustrate the use of generalizability measures in the design of heuristics for TimberWolf [30], a software package based on simulated annealing (SA) [19] to place and route various circuit components on a piece of silicon. The goal of the package is to minimize the chip area needed while satisfying constraints such as the number of layers of poly-silicon for routing and the maximum signal delay through any path. Its operations can be divided into three steps: placement, global routing, and detailed routing.

The placement and routing problem is NP-hard; hence, heuristics are generally used. SA used in TimberWolf is an efficient method to randomly search the space of possible placements. Although in theory SA converges asymptotically to the global optimum with probability one, the results generated in finite time are usually suboptimal. Consequently, there is a trade-off between the quality of a result and the cost (or computational time) of obtaining it.

In TimberWolf version 6.0, the version we have experimented, there are two parameters to control the running time (which indirectly control the quality of the result): *fast-n* and *slow-n*. The larger the *fast-n* is, the shorter time SA will run. In contrast, the larger the *slow-n* is, the longer time SA will run. Of course, only one of these parameters can be used at any time.

TimberWolf has six major components:

- 1) *cost function*,
- 2) *generate function*,
- 3) *initial temperature*,
- 4) *temperature decrement*,
- 5) *equilibrium condition*, and
- 6) *stopping criterion*.

Many parameters in these components have been well tuned manually in the last 10 years. However, their settings are generally heuristic because we lack domain knowledge to set them optimally. Moreover, the search of a single parameter set that works well across multiple circuits have been done in an ad hoc fashion. Our goal here is to show that, with a good generalization procedure, it is possible to find a single parameter set that improves both the cost and quality across multiple circuits.

Table 4 lists the parameters we have focused in our experiments and their corresponding default values. In addition, the package also uses a random seed that results in different performance when different seeds are used.

We have used seven benchmark circuits that were mostly from ftp.mcnc.org in /pub/benchmark [21] (*s298*, *s420*, *fract*, *primary1*, *struct*, *primary2*, and *industrial1*). To show that generalization works, we divided the circuits into two sets:

- 1) the first set is a learning set consisting of three smaller circuits (*s298*, *s420*, and *primary1*) that we used to experiment and find a generalized parameter set; whereas

TABLE 4  
THE PARAMETER SET IN TIMBERWOLF (VERSION 6)  
USED IN LEARNING AND GENERALIZATION

Parameter	Meaning	Default	Generalized
$P_1$	vertical path weight for estimating cost function	1.0	0.958
$P_2$	vertical wire weight for estimating cost function	1.0	0.232
$P_3$	orientation ratio	6	10
$P_4$	range limiter window change ratio	1.0	1.30
$P_5$	high temperature finishing point	23.0	10.04
$P_6$	intermediate temperature finishing point	81.0	63.70
$P_7$	low temperature finishing point	125.0	125.55
$P_8$	final iteration temperature	155.0	147.99
$P_9$	critical ratio that determines acceptance prob.	0.44	0.333
$P_{10}$	temperature for controller turn off	0.06	0.112

2) the second (the remaining four circuits, *fract*, *struct*, *primary2*, and *industrial1*) is a testing set that we used to test the new parameter set found.

In our experiments, we have studied only the standard-cell placement problem, noting that other kinds of placement can be studied in a similar fashion. We have also used *fast-n* values of 1, 5, and 10, respectively.

The *domain* in our study is the set of all performance values of possible circuits, large and small, for all combinations of parameters of TimberWolf. In this domain, we have found that different parameter values of the parameter set defined in Table 4 lead to different costs and qualities across different circuits as well as across different temperature schedules (*fast-n*). Hence, we cannot group the performance values of mappings of multiple circuits as well as multiple temperature schedules into one subdomain. In addition, we cannot subdivide the performance values of a circuit for a given temperature schedule into multiple subdomains because the only variable left in TimberWolf causing performance changes is the initial random seed. In short, we define a *subdomain* in this application as the set of performance values (quality and cost) of all the mappings for one circuit and one temperature schedule.

Since the quality and cost of a mapping generated by TimberWolf depend on the random seed used, we need to normalize the quality and cost of the mapping found using a new parameter set and a given random seed with respect to those of the mapping found using the default parameter set and the same random seed. In our experiments, we used the symmetric improvement ratio (2) as our normalization method and average the performance values over mappings of a circuit due to different random seeds. As the baseline for normalization is the default parameter set and is fixed, there are no anomalies in the ordering of parameter sets due to changing baselines.

To find parameter sets that can improve over the default parameter set, we need to have a method to systematically generate new parameter sets. Any method that explores the space of parameter sets in a systematic fashion is adequate. In our experiments, we applied TEACHER [39], a learning package based on genetic algorithms we have developed, to explore the space. TEACHER found 30 sets of parameter

values, 10 for each of the following three subdomains: *s298* with *fast-n* of 1, *s420* with *fast-n* of 5, and *primary1* with *fast-n* of 10. We used a fixed sequence of 10 random seeds in each subdomain to find its statistical performance of the mappings. Each learning experiment involved 1,000 applications of TimberWolf divided into 10 generations. Based on the best 30 sets of parameter values, we applied our generalization procedure to obtain one generalized parameter set. This generalized parameter set as well as the default parameter set are shown in Table 4.

Fig. 4 plots the quality (higher quality in the *y*-axis means reduced chip area averaged over 10 runs using the defined random seeds) and cost (average execution time of TimberWolf) between the generalized parameter set and the default parameter set on all seven circuits with *fast-n* of 1, 5, and 10, respectively. Note that all performance values in Fig. 4 are normalized using (1) with respect to those of *fast-n* of 10, and that the positive (resp., negative) portion of the *x*-axes shows the fractional improvement (resp., degradation) in computational cost with respect to the baseline parameter set using *fast-n* of 10 for the same circuit. Each arrow in this figure points from the average performance of the default parameter set to the average performance of the generalized parameter set.

Among the 21 subdomains (seven circuits and three temperature schedules), the generalized parameter set has worse quality than that of the default in only two subdomains, and has worse cost in four out of 21 subdomains. We see in Fig. 4 that most of the arrows point in a left-upward direction, implying improved quality and reduced cost.

Note that these experiments are meant to illustrate the power of our generalization procedure. We expect to see more improvement as we learn other functions and parameters in TimberWolf. Further, improvements in TimberWolf are important as the system is actually used in industry.

## 5 FINAL REMARKS

In this paper, we have defined the generalization problem, summarized various approaches in generalization, identified the credit assignment problem, and presented some solutions in measuring generalizability.

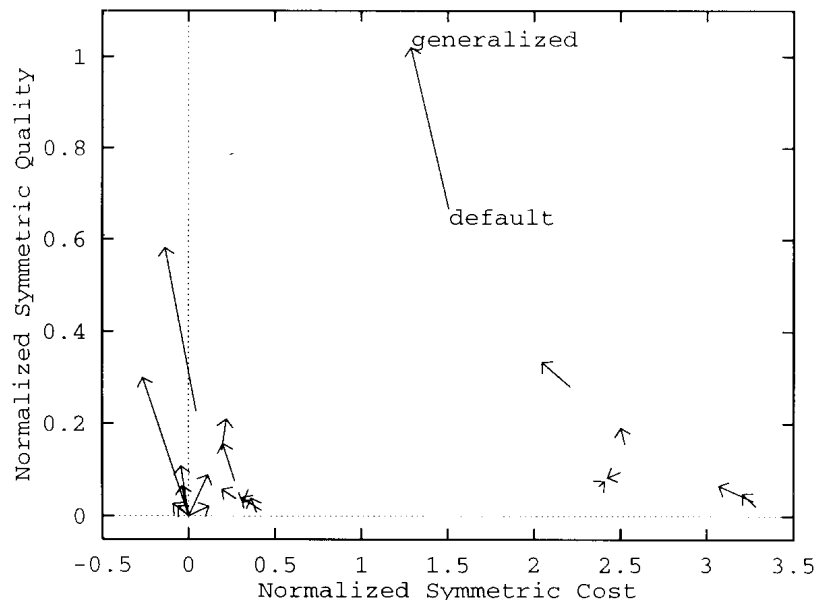


Fig. 4. Comparison of normalized average performance between the default and the generalized HMs. The plots are normalized with respect to the performance of applying the baseline HM on each circuit using  $fast-n = 10$ . (See (1).)

Existing formal results in measuring generalizability only address some restricted cases in which performance measures are from one common (possibly unknown) distribution. In general, the performance of applications may be measured by multiple metrics, and test cases may be grouped into subsets (or subdomains) such that each subset has a different performance distribution. Consequently, existing methods cannot be used to measure generalizability across subdomains of test cases.

We have presented some systematic methods to evaluate generalizability within a subdomain. To eliminate dependence on the size of a test case in a subdomain, we have shown various normalization methods to normalize performance with respect to a baseline hypothesis. Some of these methods can lead to anomalies in orderings when hypotheses are rank-ordered by the average normalized measure and the baseline is changed. Only when the baseline hypothesis is fixed (like using the median performance as the baseline) or when the effect of the baseline only exists as a constant in the average normalized measure (like using the geometric mean) that anomalies can be eliminated.

Finally, we have presented some methods to evaluate generalizability across subdomains. We have introduced a concept called *probability of win* that measures the probability that the sample mean of a hypothesis is better than the population mean of the baseline, given the number of samples tested and the variance of the samples. As probabilities of win are in the range between zero and one, they can be used to evaluate generalizability across subdomains. Unfortunately, probabilities of win cannot be used to rank-ordered hypotheses, even when performance is normalized and averaged using methods like the geometric mean. This happens because probabilities of win are used to evaluate whether a hypothesis is better than the baseline and is not meant to rank order hypotheses. The variance of the

performance values places another degree of freedom in the ordering, leading to a different ordering when a baseline with a different variance is used.

## ACKNOWLEDGMENTS

This research was supported by the National Science Foundation under Grant No. MIP 96-32316 and by the National Aeronautics and Space Administration under Grant No. NAG 1-613.

## REFERENCES

- [1] A. Barr and E.A. Feigenbaum, *The Handbook of Artificial Intelligence*, vols. 1-3, William Kaufmann, Los Altos, Calif., 1981.
- [2] A.R. Barron, "Approximation and Estimation Bounds for Artificial Neural Networks," *Proc. Fourth Ann. Workshop Computational Learning Theory*, pp. 243-249, Morgan Kaufmann, Palo Alto, Calif., 1991.
- [3] E.B. Baum and D. Haussler, "What Size Net Gives Valid Generalization?" D.Z. Anderson, ed., *Proc. Neural Information Processing Systems*, pp. 81-90, American Inst. of Physics, New York, 1988.
- [4] A. Blumer, A. Ehrenfeucht, D. Haussler, and M. Warmuth, "Classifying Learnable Geometric Concepts with the Vapnik-Chervonenkis Dimension," *Proc. 18th Symp. Theory Computing*, pp. 273-282, ACM, 1986.
- [5] L.B. Booker, D.E. Goldberg, and J.H. Holland, "Classifier Systems and Genetic Algorithms," *Artificial Intelligence*, vol. 40, no. 1, pp. 235-282, 1989.
- [6] *Encyclopaedia Britannica CD '95*, Encyclopaedia Britannica Inc., 1995.
- [7] W.W. Cohen, "Generalizing Number and Learning from Multiple Examples in Explanation Based Learning," *Machine Learning*, pp. 256-269, 1988.
- [8] G.F. DeJong and R.J. Mooney, "Explanation-Based Learning: An Alternative View," *Machine Learning*, vol. 1, no. 2, pp. 145-176, 1986.
- [9] J.L. Devore, *Probability and Statistics for Eng. and the Sciences*, Brooks/Cole, Monterey, Calif., 1982.
- [10] A. Ehrenfeucht, D. Haussler, M. Kearns, and L. Valiant, "A General Lower Bound on the Number of Examples Needed for Learning," D. Haussler and L. Pitt, eds., *Proc. Workshop Computational Learning Theory*, pp. 139-154, Morgan Kaufmann, Palo Alto, Calif., 1988.

- [11] J.J. Grefenstette, "Credit Assignment in Rule Discovery Systems Based on Genetic Algorithms," *Machine Learning*, vol. 3, nos. 2-3, pp. 225-246, Oct. 1988.
- [12] D. Haussler, "Generalizing the PAC Model: Sample Size Bounds from Metric Dimension-Based Uniform Convergence Results," *Proc. 30th Ann. Symp. Foundations Computer Science*, pp. 40-45. IEEE CS Press, 1989.
- [13] J.H. Holland, *Adaptation in Natural and Artificial Systems*. Univ. of Michigan Press, Ann Arbor, 1975.
- [14] J.H. Holland, "Properties of the Bucket Brigade Algorithm," J.J. Grefenstette, ed., *Proc. Int'l Conf. Genetic Algorithms and Their Applications*, pp. 1-7, Robotics Inst. of Carnegie Mellon Univ., Pittsburgh, Pa., 1985.
- [15] A. Ieumwananonthachai and B.W. Wah, "Statistical Generalization of Performance-Related Heuristics for Knowledge-Learn Applications," *Int'l J. Artificial Intelligence Tools*, vol. 5, nos. 1-2, pp. 61-79, June 1996.
- [16] A. Ieumwananonthachai and B.W. Wah, "Statistical Generalization of Performance-Related Heuristics for Knowledge-Learn Applications," D. Dasgupta and Z. Michalewicz, eds., *Evolutionary Algorithms in Eng. Applications*, pp. 293-313, Springer-Verlag, New York, 1997.
- [17] M. Kearns, M. Li, L. Pitt, and L. Valiant, "Recent Results on Boolean Concept Learning," *Machine Learning*, pp. 337-352, 1987.
- [18] M. Kearns, M. Li, L. Pitt, and L.G. Valiant, "Recent Results on Boolean Concept Learning," *Proc. Fourth Int'l Workshop Machine Learning*, Morgan Kaufmann, Palo Alto, Calif., 1987.
- [19] S. Kirkpatrick Jr., C.D. Gelatt, and M.P. Vecchi, "Optimization By Simulated Annealing," *Science*, vol. 220, no. 4,598, pp. 671-680, May 1983.
- [20] J.R. Koza, *Genetic Programming*, MIT Press, Cambridge, Mass., 1992.
- [21] "VLSI Layout and Synthesis Benchmarks," *Proc. Int'l Workshop Layout Synthesis*, 1992.
- [22] J.L. McClelland and D.E. Rumelhart, *Parallel Distributed Processing: Explorations in the Microstructure of Cognition*, vol. 1 on foundations, Bradford Books (MIT Press), Cambridge, Mass., 1985.
- [23] *Artificial Neural Networks: Concepts and Theory*, P. Mehra and B.W. Wah, eds., IEEE CSpres, Los Alamitos, Calif., 1992.
- [24] P. Mehra and B.W. Wah, "A Systematic Method for Automated Learning of Load-Balancing Strategies in Distributed Systems," *Int'l J. Systems Sciences*, 1997.
- [25] T.M. Mitchell, "Generalization As Search," *Artificial Intelligence*, vol. 18, pp. 203-226, 1982.
- [26] T.M. Mitchell, R.M. Keller, and S.T. Kedar-Cabelli, "Explanation-Based Generalization: A Unifying View," *Machine Learning*, vol. 1, no. 1, pp. 47-80, 1986.
- [27] D. Nguyen and B. Widrow, "The Truck Backer-Upper: An Example of Self-Learning in Neural Networks," *Proc. Int'l Joint Conf. Neural Networks*, vol. II, pp. 357-363, IEEE, 1989.
- [28] N.J. Nilsson, *The Math. Foundations of Learning Machines*, Morgan Kaufmann, Palo Alto, Calif., 1990.
- [29] N. Sauer, "On the Density of Families of Sets," *J. Combinatorial Theory*, vol. 13, pp. 145-147, 1972.
- [30] C. Sechen, *VLSI Placement and Global Routing Using Simulated Annealing*, Kluwer, Boston, 1988.
- [31] H.A. Simon and G. Lea, "Problem Solving and Rule Induction: A Unified View," *Knowledge and Cognition*, L.W. Gregg, ed., pp. 105-127, Erlbaum, Potomac, Md., 1974.
- [32] S.F. Smith, "Flexible Learning of Problem Solving Heuristics Through Adaptive Search," *Proc. Int'l Joint Conf. Artificial Intelligence*, pp. 422-425, Morgan Kaufmann, 1983.
- [33] R.S. Sutton, "Temporal Credit Assignment in Reinforcement Learning," PhD thesis, Univ. of Massachusetts, Amherst, Feb. 1984.
- [34] C.-C. Teng and B.W. Wah, "Automated Learning of the Minimal Configuration of A Feed Forward Neural Network," *IEEE Trans. Neural Networks*, vol. 7, no. 5, pp. 1,072-1,085, Sept. 1996.
- [35] L.G. Valiant, "A Theory of The Learnable," *Comm. ACM*, vol. 27, no. 11, pp. 1,134-1,142, Nov. 1984.
- [36] V.N. Vapnik and A.Y. Chervonenkis, "On The Uniform Convergence of Relative Frequencies of Events to Their Probabilities," *Theoretical Probability and its Applications*, vol. XVI, no. 2, pp. 264-280, 1971.
- [37] B.W. Wah, "Population-Based Learning: A New Method for Learning from Examples Under Resource Constraints," *IEEE Trans. Knowledge and Data Eng.*, vol. 4, no. 5, pp. 454-474, Oct. 1992.
- [38] B.W. Wah and A. Ieumwananonthachai, "Teacher: A Genetics-Based System for Learning and for Generalizing Heuristics," X. Yao, ed., *Evolutionary Computation*, World Scientific, 1998.
- [39] B.W. Wah, A. Ieumwananonthachai, L.C. Chu, and A. Aizawa, "Genetics-Based Learning of New Heuristics: Rational Scheduling of Experiments and Generalization," *IEEE Trans. Knowledge and Data Eng.*, vol. 7, no. 5, pp. 763-785, Oct. 1995.
- [40] B.W. Wah, A. Ieumwananonthachai, and Y.C. Li, "Generalization of Heuristics Learned in Genetics Based Learning," S.K. Pal and P. Wang, eds., *Genetic Algorithms and Pattern Recognition*, pp. 87-126, CRC Press, 1996.
- [41] B.W. Wah, A. Ieumwananonthachai, and T. Yu, "Genetics-Based Learning and Statistical Generalization," S. Tzafestas, ed., *Knowledge-Based Systems: Advanced Concepts, Tools, and Applications*, pp. 319-347, World Scientific, 1997.
- [42] D.A. Waterman, "Generalization Learning Techniques for Automating the Learning of Heuristics," *Artificial Intelligence*, vol. 1, nos. 1-2, pp. 121-170, 1970.



**Benjamin W. Wah** received his PhD degree in computer science from the University of California, Berkeley, in 1979. He is currently the Robert T. Chien Professor of Engineering and a professor in the Department of Electrical and Computer Engineering, the Coordinated Science Laboratory, and the Beckman Institute of the University of Illinois at Urbana-Champaign. During 1998 and 1999, he is also serving the Chinese University of Hong Kong as a professor of computer science and engineering. Previously, he was on the

faculty of Purdue University (1979-85); was program director at the National Science Foundation (1988-89); was Fujitsu Visiting Chair Professor of Intelligence Engineering, University of Tokyo (1992); and was McKay Visiting Professor of Electrical Engineering and Computer Science, University of California, Berkeley (1994). In 1989, he was awarded a University of Illinois University Scholarship and, in 1998, he received the IEEE Computer Society Technical Achievement Award. His current research interests are in the areas of nonlinear search and optimization, knowledge engineering, multimedia signal processing, and parallel and distributed processing. He was editor-in-chief of *IEEE Transactions on Knowledge and Data Engineering* from 1993 to 1996, and he now serves on the Editorial Board of *Information Sciences*, the *International Journal on Artificial Intelligence Tools*, and the *Journal of VLSI Signal Processing*. He has chaired a number of international conferences and is presently International Program Committee chair for the IFIP's year 2000 World Congress. A member of the IEEE Computer Society, he has served the Society in various capacities and is now its elected first vice president for publications. He is a fellow of the IEEE and the Society for Design and Process Science.