

History of Success and Current Context in Problem Solving

Combined Influences on Operator Selection

MARSHA C. LOVETT AND JOHN R. ANDERSON

Carnegie Mellon University

Problem solvers often have multiple operators available to them but must select just one to apply. We present three experiments that demonstrate that solvers use at least two sources of information to make operator selections in the building sticks task (BST): information from their past history of using the operators and information from the current context of the problem. Specifically, problem solvers are more likely to use an operator the more successful it has been in the past and the closer it takes the current state to the goal state. These two effects, respectively, represent the learning and performance processes that influence solvers' operator selections. A computational model of BST problem solving, developed within the ACT-R theory (Anderson, 1993), provides the unifying framework in which both types of processes can be integrated to predict solvers' selection tendencies. © 1996 Academic Press, Inc.

Most problems can be approached in multiple ways but solved by only a few. Problem solving can be viewed, then, as finding one of the few paths that leads from a problem's initial state to its goal state through some space of possible intermediate states (Newell & Simon, 1972). In this framework, problem solvers move from one state to another by applying a problem-solving operator. However, at each step in the solution, there are typically several operators that can be applied, thus forcing a selection (either implicitly or explicitly). As such, problem solving consists of a sequence of choice points at which microcosmic decisions must be made. To understand how people solve problems, then, it is important to understand how these selections are made.

Two basic approaches have been taken with respect to this question. One approach emphasizes how solvers process current information in the problem, and the other emphasizes how solvers rely on their past experience. Not surprisingly, research that emphasizes current information has shown that

This research was supported by Grant 92-53161 from the National Science Foundation and Grant N00014-90-J-1489 from the Office of Naval Research to John Anderson. The first two experiments represent two experiments from the first author's dissertation. We thank Bruce Burns, Kevin Dunbar, Chris Genovese, Abraham S. Luchins, Edith Luchins, Stellan Ohlsson, and an anonymous reviewer for their thoughtful comments on an earlier draft of this manuscript. Correspondence and reprint requests concerning this article should be addressed to Marsha C. Lovett, Department of Psychology, Carnegie Mellon University, Pittsburgh, PA 15213, USA.

problem solvers' operator selections are strongly influenced by features of the problem they are currently solving (e.g., Anzai & Simon, 1979; Atwood, Masson, & Polson, 1980; Atwood & Polson, 1976; Cooper & Regan, 1980; Jeffries, Polson, Razran, & Atwood, 1977; Larkin, 1981). For example, solvers often use *hill climbing*, a heuristic that leads them to make choices that take the current state of a problem closer to the goal state. In contrast to this current information approach, research that emphasizes past experience aims to capture the learning that occurs during problem solving and thereby demonstrate the impact of past experience on operator selection (e.g., Luchins, 1942; Luchins & Luchins, 1959; Reder 1987, 1988; Thorndike, 1932). Here, a typical result is that solvers tend to select an operator (or operator sequence) that was successful on past problems.

Although current and past information have generally been studied separately, they must both be processed to some degree in all problem-solving tasks. For example, consider a person faced with the task of navigating to a particular destination. At each choice point along the way (e.g., going left, right, or straight at each intersection), the person may be influenced by the current surroundings (e.g., to pick the direction that appears to head closest to the destination). If this person also has some relevant navigating experience (perhaps even on the same route), knowledge derived from past experience may represent another influence on the solver's selections (e.g., to follow the route that worked on the last trip). When current and past information suggest different selections, however, the solver is faced with a difficult trade-off: how to combine current information and past experience to decide on a particular path? The existence of this kind of situation in real-world problem solving suggests that to understand problem solving and operator selection fully, we must understand how current information and past experience interact and combine.

The two emphases on current information and past experience map onto another standard distinction in problem solving: the study of problem-solving performance versus skill acquisition. These two areas have generally been treated separately: Research on performance tends to focus on the current information people use within each problem, whereas research on skill acquisition tends to focus on the experience and knowledge people gain and use across multiple trials. Such divergent research paradigms make it difficult to integrate models of problem-solving performance with those of skill acquisition, and yet both areas could benefit from such a marriage.

In this paper, we seek to integrate processing of current information (performance) with processing of experience-based information (learning) to develop a single model of operator selection. We will show that *both* types of information play an important role in the selections people make as they solve problems, and we will present a computational model that integrates current and experience-based information in a way that leads to operator selections quite similar to human solvers'.

The key to this integration is a conceptualization of current information and past experience in a way that allows both to be quantified and combined meaningfully. We accomplish this by specifying two variables that represent the two types of information. To study the effects of current information, we focus on the impact of a particular current problem feature—distance between the current state of the problem and the goal state. To study the effects of experience, we focus on the impact of solvers' histories of success and failure with the available operators. While previous work has shown that various measures of distance-to-goal influence operator selections (Atwood *et al.*, 1980; Atwood & Polson, 1976; Jeffries *et al.*, 1977; Newell & Simon, 1972), no detailed analysis has been completed that describes the quantitative relationship between solvers' histories of success and their operator selections. Our goal is to use these two variables together, both in experiments and in a computational model, to explore in detail how they individually *and jointly* affect operator selection.

Background on Distance-to-Goal and History-of-Success

The effect of distance-to-goal on operator selection is a robust problem-solving phenomenon. For example, college students solving a variant of the water jars task (used by Luchins, 1942) tend to select the move (out of three moves available, on average) that will take them closest to the goal state (Atwood & Polson, 1976). This problem-solving heuristic, called hill climbing, and its close cousin, means-ends analysis,¹ have been used to explain solvers' tendency to select operators that will take them closer to the goal (or the current subgoal) in a variety of domains (Atwood & Polson, 1976; Bhaskar & Simon, 1977; Jeffries *et al.*, 1977; Klahr, 1985; Larkin, 1981). Models of problem solving that include such metrics for evaluating how close various operators will take the current problem state to the goal state are thus quite successful at making the same operator selections that solvers do (e.g., Atwood & Polson, 1976; Atwood *et al.*, 1980; Larkin, 1981; Newell & Simon, 1972). These models of the processing of current information, however, are generally rather static and cannot learn additional information about the operators being used.

The effect of history of success, in contrast, takes into account the learning that occurs during past problem-solving experiences and its impact on operator selection. A classic study of the history-of-success variable was performed by Luchins (1942; Luchins & Luchins, 1959). In a series of experiments using the water jars task, subjects in an experimental group were trained on problems all solved by the same complex sequence of operator applications.

¹ The means-ends heuristic is more complex than the hill-climbing heuristic. It involves recursively setting subgoals to reduce large differences between the current state and the goal state and applying operators that will allow those subgoals to be achieved.

These subjects experienced many successes with that operator sequence. When given test problems that could be solved in a simpler, more efficient manner, they continued to use the same complex operator sequence. Luchins labeled this perseveration of the previously successful solution method the *Einstellung* effect. Subjects in a control group who did not receive the training problems chose simpler, more efficient solutions to solve the same test problems. These results demonstrate that solvers' histories of success, and not features of the current problem, can dominate the operator-selection process.

Reder (1987, 1988) has also presented evidence for the influence of previous success on subsequent selections by giving different groups of subjects different reading comprehension questions answerable either by (a) a plausibility judgment strategy or (b) a retrieval strategy. Subjects who were initially given a mix of 80% plausibility and 20% retrieval questions were more likely to use plausibility judgment on later questions, whereas subjects who were initially given a 20% plausibility, 80% retrieval mix were more likely to use retrieval. As in Luchins's studies, these results demonstrate that subjects were influenced by what operators had worked in the past.

It is important to note that two of the above empirical results are based on the water jars task, and yet each one focuses on a distinct approach of the two mentioned above. This tendency toward separation is also true with respect to models of problem solving, such that many models are developed to emphasize either problem-feature information or history-of-success information but not both. The vast majority of models, then, cannot answer questions such as: What do solvers do when the operator that takes the current state closest to the goal has a relatively unsuccessful past record or vice versa? To model this kind of tradeoff, we need a quantitative theory of operator selection that not only more directly relates history of success to operator selections but also explains how solvers combine these two sources of information to make operator selections.

An ACT-R Model of Operator Selection

The ACT-R theory (Anderson, 1993) provides a natural way for both distance-to-goal and history-of-success information to be combined in operator-selection decisions. To discuss this combination, however, we first need to describe how the relevant knowledge is represented. In ACT-R, much of skill learning is represented by the formation and tuning of *production rules*. Each production rule (or production) has a set of conditions describing when it can apply and an action that is executed when it does apply. A production represents the solver's (procedural) knowledge about taking those actions under those conditions. The notion that each production represents a separate unit of procedural knowledge is central to the ACT* and ACT-R theories (Anderson, 1983; Anderson, 1993). In addition, ACT-R's main tenet is that the processes acting on these productions (and on other declarative structures) are tuned to the statistical structure of the environment (Anderson, 1990).

These two ideas provide the framework in which an ACT-R model of operator selection uses both distance-to-goal and history-of-success information, as we describe next.

In an ACT-R model of operator selection, the process of selecting among several problem-solving operators corresponds to the process of selecting among several production instantiations. (A production instantiation is simply a production mapped onto the specifics of the current situation; it represents a single move in problem solving.) This selection among production instantiations is based on the model's evaluation of which one has the highest probability of success and lowest expected cost. Since the ACT-R theory prescribes that this process should be adapted to the statistical structure of the environment, the model is sensitive to the fact that actions which end up closer to the goal tend to lead to achievement of the goal with greater probability and lower cost. Thus, distance-to-goal information plays a key role in the operator-selection process due to its importance as a predictor of success and cost in the real world. Consequently, an ACT-R model of operator selection evaluates a production instantiation more highly the closer that move takes the current state to the goal state. The model we propose in this paper is endowed with a hill-climbing distance metric similar to other models' (Atwood *et al.*, 1980; Atwood & Polson, 1976).

ACT-R's evaluation of various production instantiations also allows for the influence of history of success on operator selection. Since the production is the unit of procedural knowledge in ACT-R, history-of-success information is recorded at the production level. That is, a production's "history of success" includes the number of times past instantiations of that production led to success as well as the number of times past instantiations of that production led to failure, but it does not include information about the problems or contexts in which those successes and failures occurred. This history-of-success information is important because it indicates how likely a production is to lead to success, averaged over all the past uses of that production. The idea underlying production-based learning of success is that the conditions associated with each production represent a constrained set of situations over which history-of-success information can be appropriately aggregated; past success with a production is predictive of future success with *that* production. Thus, according to the ACT-R theory, the evaluation of a potential move will be higher the better the corresponding production's history of success.

This specification of ACT-R leads to several predictions for operator selection. First, solvers will be more likely to select a particular step the closer it takes the current state to the goal state (where "closer" is defined by the model's distance metric). Second, solvers will be more likely to select a problem-solving step the more successful past instantiations of the corresponding production have been in the solver's history. These first two predictions stem from ACT-R's rational analysis of how operator selections would be made if they were optimized to the structure of the environment (Anderson,

1990). Third, experiences of success or failure with a production on one particular problem type will impact subsequent selections of that production for all new problems (as long as the same production is applicable). This predicted “generalization” effect is a result of ACT-R’s claim that success information is maintained at the production level without storage of problem-specific information. Fourth, latencies to make different problem-solving moves should decrease as a power function of the amount of practice with each production. (This point will be discussed further after the details of our model are given.)

In the experiments below, we attempt to test these predictions through a systematic study of how solvers use and integrate distance-to-the-goal and history-of-success information in operator selection. The first experiment below is a demonstration that solvers are indeed sensitive to both types of information in a task we have developed called the building sticks task (BST). This experiment also replicates the *Einstellung* effect (Luchins, 1942) and shows that an ACT-R model of operator selection can account for this effect. Experiments 2 and 3 provide more comprehensive experimental manipulations of the relevant variables and test all four predictions above. In addition to qualitative tests of these predictions, we describe the specific ACT-R model of BST problem solving that we have implemented and perform quantitative tests of the model’s fit to the data.

EXPERIMENT 1

Experiment 1 is an extension of Luchins’s (1942) “*Einstellung*” experiments. In those original experiments, Luchins showed that subjects would continue to use the same solution method that had worked well on previous problems, even when it was no longer appropriate. One task he used to study this effect is called the water jars task (WJT). A WJT problem is defined by its three jar capacities (e.g., $A = 21$, $B = 127$, and $C = 3$) and a desired quantity (e.g., 100). The solver’s goal is to obtain the desired quantity by filling and pouring from the three jars. For the example above, this goal can be achieved by filling jar B from a tap, pouring from jar B to fill jar C twice (emptying jar C between fillings), and pouring from what is left in jar B to fill jar A. This solution is denoted B-2C-A. In Luchins’s original experiments, solvers were given WJT problems in the numerical format above (without real jars or liquid).

In his first experiment, Luchins (1942) found that, of the subjects who had used the complex solution method successfully for five training problems in a row, 81% used the same method on test problems that could be solved more simply (e.g., $A + C$), and 54% failed to solve a test problem that could not be solved by the original method. After this “failure” problem, 63% of subjects used the previously successful method on subsequent problems that could be solved more simply. In contrast, subjects in a separate control group, who received only the test problems, never used the complex solution method.

These results suggest that subjects in the former group were responding to their histories of success: They were very likely to reuse the method that had succeeded many times, and then after it failed, they were somewhat less likely to use it.

However, the structure of the WJT as implemented in these studies makes it impossible to measure subjects' exact histories of success with different operators; rather, they must be inferred. For instance, WJT solvers were not required to record their intermediate steps on paper, and they could easily use mental arithmetic to make various unobserved solution attempts before writing their final answer. Moreover, because there was no explicit feedback in the WJT, the success or failure that these subjects "experienced" did not always conform to objective success or failure. Subjects could (and did) make arithmetic mistakes that led them to *believe* a particular operator sequence had succeeded when, in fact, it had not. These factors make a study of the more subtle effects of history-of-success on operator selection infeasible within that version of the WJT.

To avoid these difficulties, we developed the BST, a task that is isomorphic to the WJT. The BST was designed to allow for a more dense and complete protocol of subjects' operator selections as well as a more accurate record of their histories of success with each operator. This is achieved by presenting each problem in terms of stick lengths with no numerical representations (i.e., three building sticks that can be added and subtracted to attain a desired stick length). Note that combining and comparing stick lengths in one's head is relatively inaccurate and error-prone; this deters subjects from doing much "mental" look ahead in the BST beyond the next step. In addition, because of our computer interface, BST solvers must make explicit every step they want to take by clicking on a certain area of the screen with the computer mouse. This allows us to record every intermediate step solvers make in a non-intrusive way. Another advantage of the BST for current purposes is that it makes success and failure salient and objectively distinct. For example, success is indicated when the current stick and the desired stick are equal (a readily apparent state), and subjects must explicitly acknowledge their success or failure by clicking on a "DONE" button or a "RESET" button, respectively.

The top rectangle in Fig. 1 presents a BST problem as solvers would see it in its initial state. Solvers must add and subtract the lengths of the building sticks to create a (current) stick that is equal in length to the desired stick. For instance, the bottom of Fig. 1 displays three paths emanating from the initial state, each one representing the state that results from applying a single operator. The left and right branches represent application of what we label *undershoot*, the operator for building up to the desired stick by initially undershooting the goal. Specifically, undershoot is instantiated here by selecting either of the two building sticks shorter than the desired stick and adding it to the current stick of length zero. The middle branch in Fig. 1 represents

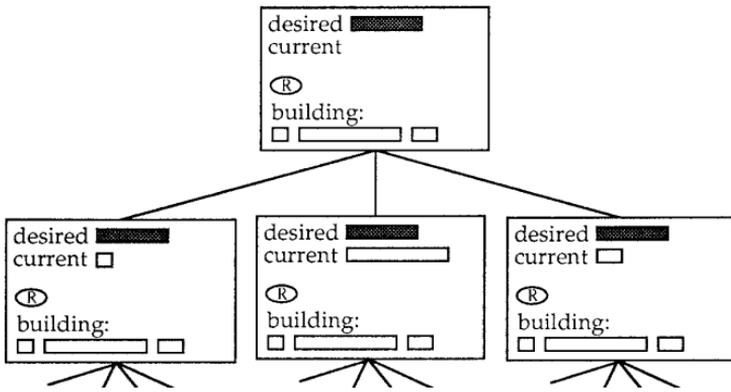


FIG. 1. Initial and successor states in the Building Sticks Task. The circled R represents reset.

an application of *overshoot*, the operator for initially overshooting the goal with a building stick that is longer than the desired stick. Overshoot is instantiated here by selecting the building stick that is longer than the desired stick and adding it to the current stick of length zero. Following an application of undershoot, solvers must select among moves that add to the current stick, whereas following an application of overshoot, solvers must select among operators that subtract from the current stick.²

Our model of operator selection in the BST makes several predictions regarding subjects' selections between overshoot and undershoot. These predictions are not limited to the problems designed to test for *Einstellung* but apply throughout the experiment. With respect to history of success, the prediction is that subjects' use of the overshoot and undershoot operators should mirror the relative success of those operators. If subjects experience a change in the success rates of the two operators, their selections should change in the same way. Also, if different subjects' histories of success are substantially different, then their operator selections should be different. Our model predicts a specific functional relationship between past history of success and subsequent selections. To test this relationship in multiple ways, we went beyond Luchins's design and tracked three different groups of subjects who experienced different histories of success (after a common training phase). By gathering data on each problem (e.g., which operator subjects used and whether it led to success or failure), we were able to record a constantly changing history of success for the subjects in each group. Given

² We define undershoot and overshoot as the two competing operators because (1) pilot subjects who provided talk-aloud protocols while solving the BST often made reference to the plans that undershoot and overshoot represent and (2) having two operators compete at the beginning of each problem instead of three (e.g., an operator for each building stick) has the advantage of parsimony.

our model, these different history-of-success profiles lead to predictions of different operator-selection tendencies that can be tested. In addition, this experiment allows us to verify that subjects' operator selections are also sensitive to distance-to-goal information as our model predicts.

Method

Subjects. Subjects in the experiment were 45 Carnegie Mellon University undergraduates, participating for credit or for \$5. The former were recruited from the Psychology Department subject pool, and the latter responded to an advertisement posted on a University electronic bulletin board.

Design. In this experiment, subjects were randomly assigned to one of three conditions. The conditions differed only in the type of BST problems presented on the ninth through eleventh trials (the variable phase). The first eight problems were the same for all subjects (the training phase); all had one building stick longer and two building sticks shorter than the desired stick, and all were solved by overshoot. We label these "o" problems. Subjects could attempt either overshoot or undershoot on "o" problems, but only overshoot would lead to a solution. Thus, the training phase was designed to lead to a strictly increasing number of successes for overshoot and failures for undershoot.

The variable phase followed the training phase and included one of three problem types: "o" problems that had one building stick longer and two building sticks shorter than the desired stick and could only be solved by overshoot, "u" problems that looked similar to "o" problems but could only be solved by undershoot (and in fewer steps than the "o" problems), and "U" problems that had all three building sticks shorter than the desired stick and could only be solved by undershoot. Note that overshooting on the first move was impossible for "U" problems, whereas both operators could be attempted on the first move for "o" and "u" problems.

These three problem types in the variable phase were designed to provide subjects in the different conditions with different histories of success. Subjects receiving "o" problems were supposed to experience overshoot successes; subjects receiving "u" problems were supposed to experience overshoot failures and undershoot successes; and subjects receiving "U" problems were supposed to experience no overshoot failures (because that operator was not applicable) but undershoot successes. After the variable phase, all subjects solved the same four "o" test problems. Three of the test problems were repeats from the first half of the training phase.

Apparatus. Subjects worked individually on a Macintosh IIfx computer. The BST interface was designed in cT 2.01 (Physics Academic Software, 1992) such that subjects only needed to use the mouse. The cT program that ran the experiment also collected information on each mouse click into a data file. Each rectangle in Fig. 1 is a sketch of the interface subjects saw.

Procedure. Before beginning the experimental trials, subjects read instruc-

tions for the BST on the computer screen and practiced each of the following actions: selecting a building stick, placing the selected stick in the building area, adding to the current stick, subtracting from the current stick, and restarting the problem. The instructions encouraged subjects to use the reset button whenever they wanted to restart the problem instead of trying to undo each move they had taken.

When subjects finished reading the instructions, the experimenter asked them to watch carefully as she solved two sample problems. Both sample problems had one building stick longer than the desired stick. The first sample problem was solved by an undershoot solution and the second by an overshoot solution so all subjects would have exposure to both operators. The experimenter demonstrated these solutions by making the correct moves in the same order for all subjects, pausing slightly between each move. (In subsequent experiments, these demonstration problems were presented automatically by the computer.) Subjects' remaining questions were answered before beginning the experimental trials.

The BST problems were divided into the three phases described above, although the transitions between phases were not indicated to subjects. After subjects completed all 15 BST problems, they were asked if the experiment had reminded them of anything they learned about in their psychology classes. Of the 45 subjects, only three mentioned the water jars task, and none remembered the relevant *Einstellung* result.³ Our analyses include all subjects' data.

Data analysis. Two types of dependent variables were used to analyze subjects' operator selections in this experiment. The first is a binary variable representing which operator a subject selected as the first move for a given problem—overshoot or undershoot. By design of the task and interface, it was easy to distinguish between the selection of these two operators by observing which stick a subject selected for the first move. The overshoot operator must be initiated with the selection of a building stick that is longer than the desired stick, and the undershoot operator must be initiated with the selection of a building stick that is shorter than the desired stick. This dependent measure is usually presented as a percentage of subjects whose first move was overshoot for a given problem. The other type of dependent variable used in this experiment is latency to make the first move. This latency is measured from the time the problem appeared on the screen until the subject selected an initial building stick.

Data on moves other than the first move of each problem are not presented below, but they are well accounted for by a simple hill-climbing metric. Over

³ Perhaps subjects mentioned Luchins' results so rarely because the experiment was conducted at the beginning of the semester before the *Einstellung* effect was taught in their courses. Subjects who mentioned the WJT but not its related empirical result were among those students in a general course in cognitive processes, in which the WJT was used to demonstrate verbal protocol analysis.

90% of solvers' nonfirst moves in an initial solution attempt took the current stick's length as close as possible to the desired stick's. When the initial solution attempt did not lead to a solution, subjects would restart the problem. Then, when they were back at the problem's initial state, they would usually choose the operator (overshoot or undershoot) that they had not already attempted on that problem and, again, generally follow a simple hill-climbing metric from then on.

Results and Discussion

Selection results. The percentage of subjects selecting the overshoot operator as their first step, for each of the eight training problems, is presented in the left portion of the top panel of Fig. 2. The first eight data points represent the training data for all three conditions. They are quite variable, presumably because of large differences in the training problems and because all subjects solved these training problems in the same order. As mentioned above, the important current problem feature predicted to affect operator selections is distance to the goal. To test whether the variability in these training data was due to this variable, we defined a measure of problem *bias*: A problem's bias is the difference in the hill-climbing distances for undershoot (selecting the medium-sized stick) and overshoot (selecting the big stick). The more positive a problem's bias, the closer overshoot takes the initial state to the goal; the more negative a problem's bias, the closer undershoot takes the initial state to the goal. (For undershoot, we only consider the medium-sized stick because that move always dominates the undershoot move based on the smallest stick.)

We used two distance metrics for the following analysis: one was based on differences and the other on ratios. They lead to very similar results, so we present the "difference" metric here and the "ratio" metric with the details of our model. According to the "difference" distance metric, the hill-climbing distance between the goal state and the state following an undershoot move is $|\text{desired} - \text{medium}|$, and the distance between the goal state and the state following an overshoot move is $|\text{desired} - \text{big}|$. Problem bias is the difference of these two differences, or $|\text{desired} - \text{medium}| - |\text{desired} - \text{big}|$.

A multiple regression with percentage of subjects using overshoot as the dependent variable and with problem bias and problem number as the independent variables shows that both variables have a significant, positive influence on selection: problem bias $t(5) = 4.21, p < .01$; problem number $t(5) = 3.62, p < .05$. The effect of problem bias suggests that subjects are indeed influenced by the distance-to-goal information: they are more likely to select overshoot the greater a problem's bias toward overshoot. The effect of problem number provides our first evidence of subjects' sensitivity to history-of-success information in the BST. This is because problem number is a proxy for the success of overshoot during the training phase. Figure 3 shows the percentage of subjects using overshoot, adjusted for the effect of problem

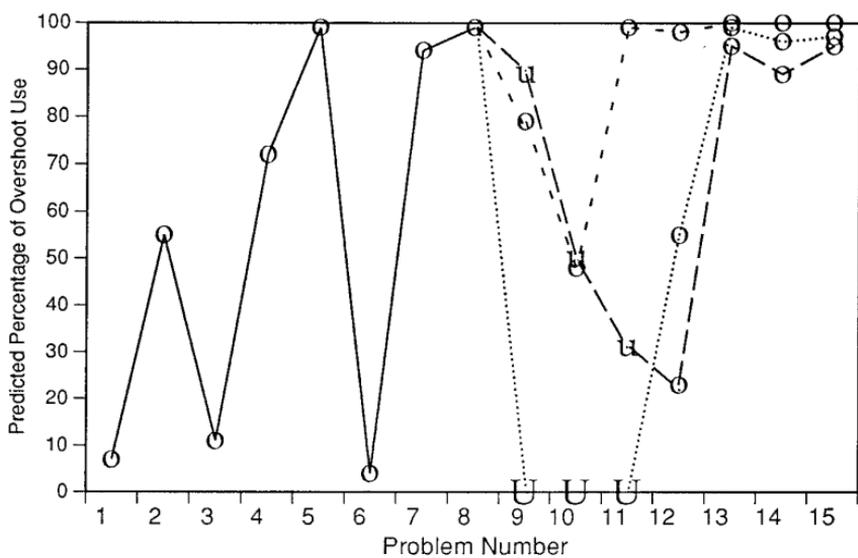
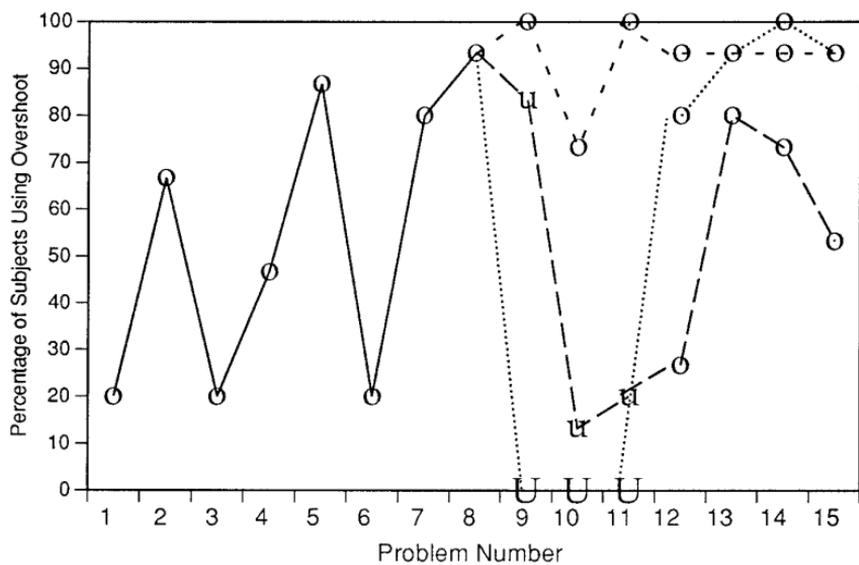


FIG. 2. Percentage of subjects using overshoot (top panel) and predicted percentages (bottom panel) for Experiment 1. Each data point is labeled by its problem type: "o" problems have one stick longer than the desired stick and are solved by overshoot; "u" problems have one stick longer than the desired stick and are solved by undershoot; "U" problems have no sticks longer than the desired stick and are solved by undershoot.

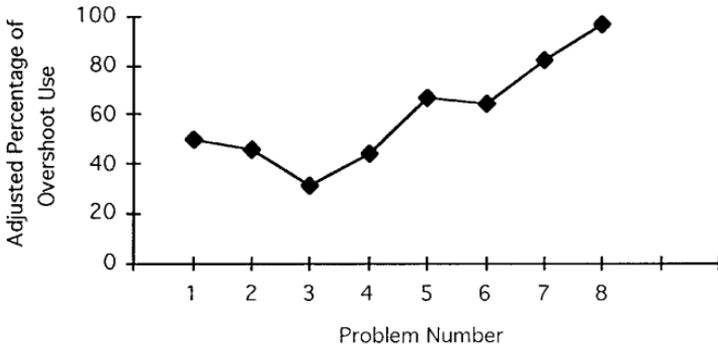


FIG. 3. Adjusted percentage of subjects using overshoot during training phase of Experiment 1.

bias; note the improved monotonicity (relative to Figure 2) across the training phase. Finally, note that these two effects together account for more than 86% of the variance in subjects' selection data.

After the training phase, the next effect to look for was an *Einstellung* effect. In this experiment, *Einstellung* would be exhibited if subjects receiving "u" problems in the variable phase (i.e., problems solvable only by undershoot) nonetheless continued to select overshoot. Indeed, on the first problem of the variable phase, 73% of the subjects receiving "u" problems initiated their solution with overshoot. These subjects were significantly more likely to use overshoot, the operator that had previously been successful, than a comparable group of 15 subjects (from a separate experiment) who were given the same problem with no previous BST experience, $z = 8.3$, $p < .01$. Thus, we reproduced an *Einstellung* effect in the BST.

One difference between Luchins's original experiment and this one is that, on his *Einstellung* problems, the previously successful solution method *and* a simpler method could lead to success, whereas, in our variable "u" problems, only a simpler method could lead to success. Thus, when solving our variable "u" problems, subjects who exhibited the *Einstellung* effect (i.e., selected overshoot) would experience overshoot failures and, hence, would be expected to select overshoot less often subsequently. Thus, we used only the first variable "u" problem as our test of *Einstellung*; averaging over all variable "u" problems would have attenuated the effect. Nevertheless, since our model predicts a decrease in *Einstellung* as a result of the undershoot successes and overshoot failures during the variable "u" problems, we test for it: Subjects in the "u" group selected overshoot significantly less often on the first problem of the test phase than on the first problem of the variable phase, $z = 2.9$, $p < .01$.⁴ This confirms the prediction that "u" subjects

⁴ One caveat to this comparison is that it is based on operator selections across two different problems, and we know from the training phase data that different problems can have a big impact on operator selections. So, we compared the two problems involved in this comparison,

should exhibit a weaker Einstellung effect after the variable phase than they did at its beginning. Indeed, it suggests a reversal of Einstellung similar to that observed in Luchins's (1942) data after the "failure" problem.

The main reason for including the variable phase in this experiment was to provide subjects in the three conditions with different histories of success with overshoot and undershoot so that we could compare selection tendencies between groups. The middle portion of the top panel of Fig. 2 presents the percentage of subjects in each condition who used overshoot on each of the variable problems (problems 9 through 11). Since subjects in the three conditions solved different problems in the variable phase, we do not make comparisons of these data. Rather, we present the variable phase data to demonstrate that subjects in the different conditions generally experienced the profile of successes and failures that the variable phase was designed to create: Subjects receiving "o" problems were very likely to use overshoot, leading them to experience three additional, guaranteed overshoot successes and few undershoot failures; subjects receiving "u" problems were also quite likely to use overshoot (at least on problem 9), leading them to experience many overshoot failures along with three guaranteed undershoot successes; and subjects receiving "U" problems could not use overshoot, leading them to experience no overshoot failures along with the three guaranteed undershoot successes.

With these different histories of success, our model predicts differences between particular conditions in operator-selection tendencies during the test phase (right portion of the upper panel of Fig. 2). Subjects who received "o" or "U" problems during the variable phase experienced similarly few overshoot failures. Hence, subjects in both of these conditions were expected to be very likely to use overshoot during the test phase. In contrast, subjects who received "u" problems during the variable phase experienced many overshoot failures, so they were expected to be less likely to use overshoot during the test phase. Consistent with these predictions, we found that subjects receiving "o" or "U" problems were not significantly different from each other in their high propensities to use overshoot during the test phase, $z = 0.4$, n.s., and yet subjects receiving "u" problems were significantly less likely to use overshoot during the test phase than subjects in the other two groups, $z = 3.95$, $p < .01$. In other words, the extra overshoot failures experienced by subjects receiving "u" problems led to a significant reduction in their selections of overshoot, relative to the other groups.

Finally, we can compare subjects' tendencies to select overshoot in the test phase to their tendencies to select overshoot on the same problems when they appeared during the training phase. Since all subjects experienced more suc-

and, although their stick lengths are different, their problem biases (according to both distance metrics) are virtually indistinguishable. This provides some evidence that the significant difference in selections is not due to differences in problem bias.

Overshoot:	Undershoot:
IF goal is to solve a BST problem & current stick is length 0 & desired stick is length <i>desired</i> & a bldg stick has length <i>stick</i> & <i>desired</i> < <i>stick</i>	IF goal is to solve a BST problem & current stick is length 0 & desired stick is length <i>desired</i> & a bldg stick has length <i>stick</i> & <i>desired</i> > <i>stick</i>
THEN add <i>stick</i> to the current stick & set a subgoal to subtract	THEN add <i>stick</i> to the current stick & set a subgoal to add

FIG. 4. Overshoot and undershoot productions.

cess with overshoot than undershoot, we combined the three conditions and predicted greater tendencies to select overshoot in the test phase than in the training phase. (Note that combining the three groups should only attenuate any increase in overshoot use since “u” subjects experienced some undershoot successes and overshoot failures in the variable phase.) This comparison confirmed our prediction by revealing that subjects were much more likely to select overshoot in the test phase than on the same problems in the training phase, 86% vs 29%, $z = 7.73$, $p < .01$.

The general, and quite surprising, conclusion from this experiment is that subjects’ operator selection tendencies could be distinguished after the experience of relatively minor changes in their history of success: We found differences between conditions that varied by as few as three problems as well as differences across time when the same problems were repeated only nine problems apart. Our model predicted all of these qualitative differences based on history-of-success information, but it remains to be seen how well the model provides quantitative fits to the selection data.

Description of the Model

Up to this point, we have only given a high-level description of our model and its predictions for operator selections in the BST. This section provides a detailed description of our model and presents quantitative fits to the selection data from Experiment 1. Appendices A and B provide a derivation of the model and more details of the fitting procedure.

The BST model is a production system within the ACT-R framework.⁵ Figure 4 provides an English version of the productions that correspond to overshoot and undershoot. Only when a production’s conditions (listed after the “IF”) match the current situation may its actions (listed after the “THEN”) be applied. Note that the overshoot and undershoot productions will both match a given BST problem when the current stick has length zero

⁵ The BST model may be obtained via the World Wide Web from <http://act.psy.cmu.edu/act/papers/BSTms.html>.

(i.e., the problem is in its initial state) and there is one building stick longer and another shorter than the desired stick. Other productions in the model include subtracting from the current stick's length, adding to the current stick's length, and resetting the problem to its initial state.

These productions allow the model to solve the same BST problems subjects do; in fact, these productions lead to the same operator-selection decisions that subjects face. For example, consider a problem with a desired stick of length 100 and building sticks of length 30, 170, and 40. The undershoot production matches with either the building stick of length 30 or 40, leading to two undershoot instantiations or moves (e.g., `UNDERSHOOT(30)` and `UNDERSHOOT(40)`). The overshoot production also matches, and it can be instantiated with the building stick that is longer than the goal stick (e.g., `OVERSHOOT(170)`). Therefore, three production instantiations are available to the model, just as three operator applications, or possible moves, are to solvers. Since the model (and solvers) can only execute one move at a time, we must specify a method for selecting among the production instantiations.

As mentioned earlier, an ACT-R model of operator selection is designed to select the production instantiation that yields the highest probability of success and lowest cost. Without loss of generality, we simplify our model to take into account only probability of success. (This is possible because in the BST probability of success plays a major role in differentiating among potential moves whereas cost does not.) In our model, as in the more general case, the actual probability of success associated with each possible move cannot be known in advance, so the model instead bases its selection on estimates of the predicted probability of success (PPS) of each move. To account for assumed inexactness in the internal computation of these estimates, independent Gaussian noise is added to each, and the move with the largest (noisy) PPS is selected. If two moves have very different estimated PPS values, then the noise is unlikely to affect the model's selection of the move with the larger estimate; however, if the estimated PPS values are very close, then the model will choose between the two moves essentially at random. Note that the added noise has mean zero and a variance that is constant across time; this variance is a parameter of the model.

Now, it only remains to describe how the model estimates PPS. Each move's PPS is a function of two quantities associated with the move—distance to the goal after the move is taken and history of success of the corresponding production. Below, we describe how our model computes these two values individually, and then we discuss how it combines them.

Recall that distance to the goal is measured in terms of the problem state that would occur after the potential move is taken. After an overshoot or undershoot move, the problem state's current stick will be equal to the building stick used in the move. Thus, distance-to-goal is a function of the desired stick's length and the length of the stick involved in the move. It is possible that there is some variability among solvers in the particular metrics used to

evaluate distance-to-goal. As such, we have considered two distance metrics for the computation of distance-to-goal in our model. The first of these, based on differences, was described above. The second, based on ratios, computes distance-to-goal (*dtg*) as:

$$dtg = \begin{cases} \frac{\text{desired}}{\text{medium}} & \text{for undershoot} \\ \frac{\text{big}}{\text{desired}} & \text{for overshoot} \end{cases}$$

(We normalize these ratios so that both undershoot's and overshoot's distance-to-goal will increase along the same scale: the closer the value is to 1, the closer the distance.)

Note that the ratio and difference metrics are strongly correlated and that they lead to very similar results. Thus, the choice of metric does not significantly affect the model fits reported in this paper. We use the ratio metric in our description and analysis of the model because we found it to be slightly more consistent with the retrospective reports provided by pilot subjects. More experimentation would be required to resolve the question of distance metric for situations in which it did make a difference.

The other value contributing to the PPS of a move is an estimate of the history of success of the production involved in the move; it is defined as $(\alpha + s)/(\alpha + \beta + s + f)$, where s and f are the number of past successes and failures associated with the production and α and β are parameters of the model associated with the production. This history-of-success estimate for each production is a Bayesian posterior probability of success, given the production's totals for success and failure and a Beta(α , β) prior distribution for the probability of success (see Berger, 1985). Note that a production's history-of-success value will asymptote to its actual proportion of successes as experience with the production accumulates.

According to a rational analysis of problem solving (Anderson, 1990), each of these two quantities plays a special role in estimating PPS. Distance-to-goal information is important because, in general, the distribution of such distances is different for states that will lead to eventual success versus eventual failure (i.e., distance to the goal tends to be greater in the failure case). And, since the ACT-R theory claims that problem solving should be tuned to this kind of statistical regularity, our model takes distance to the goal as inversely related to PPS. In addition, history-of-success information is important because it represents the expectation that a problem will be solvable by a particular production, averaged over all problems for which the production is applicable. Thus, history of success represents the statistical regularity associated with a particular production in the BST environment, and so our model takes history of success as directly related to PPS.

The key feature of our model's computation of PPS, then, is that it integrates

distance-to-goal and history-of-success information in a way that is consistent with a rational analysis of problem solving. Appendix A provides a mathematical specification of this analysis for the BST. The basic result is that the model's estimate of the log odds of success for a particular move (i.e., $\log(\text{PPS}/(1 - \text{PPS}))$) is related to distance-to-goal and history-of-success by the formula

$$\text{Log odds} = A - b * \text{Distance} + \text{History} \quad (\text{Central Equation}),$$

where Distance is the hill-climbing distance after the move is taken, History is the log odds version of history-of-success for the production involved in the move which is calculated as $\log((\alpha + s)/(\beta + f))$, and A and b are constants related to a_0 and b_0 discussed in Appendix A.

The Central Equation shows that moves with smaller distances to the goal and with productions having greater histories of success will have higher odds of success and, hence, will be more likely to be selected. Moreover, this equation shows that the history-of-success component acts as a scaling parameter, generally shifting the effect of distance-to-goal up or down. Critically, the model predicts that the history effect will be experienced throughout the distance spectrum, as long as the production still matches and will not depend on which particular problems the operator succeeded or failed (e.g., the history component of the Central Equation does not involve any variables specific to the problem context in which past successes and failures occurred). These are distinctive claims of our model that we will call the "generalization" or "independence" predictions.

The Central Equation and the definitions of History and Distance are presented in their simpler "log odds" format above. In this equation, the effects of History and Distance add together to determine the model's estimate of the log odds of success of a move. When the model's estimate of log odds of success is transformed to PPS, its range is limited to 0–1 and additivity will fail near the boundaries. This occurs when distance-to-goal is very large or very small or when history-of-success is extreme. Such ceiling and floor effects may arise in the data but do not argue against the independence prediction of our model.

Finally, with respect to latency data, our model makes the prediction that latencies will increase as a function of practice—but not as a function of overall practice with the task. Because of its production-based representation, our model predicts that a power-law speedup in latencies will occur as a function of number of practice opportunities with each production. (See Anderson, 1982, for more details.)

Fit of the Model to the Data

With an exact specification of our model given above, we can now assess how well the model reproduces subjects' selection behavior under particular

settings of the parameters. The model parameters required for fitting the full data are the noise variance σ^2 , and α_p and β_p for each production p . Since we are focusing on subjects' selections between overshoot and undershoot, we need only consider α_p and β_p for these two productions. Hence, the parameters are σ^2 , α_o , β_o , α_u , and β_u . We constrain these parameters such that $\alpha_o + \beta_o = \alpha_u + \beta_u$, $\alpha_u/(\alpha + \beta) = 1 - \alpha_o/(\alpha + \beta)$, and $\sigma = 0.05$. The first constraint is motivated by the fact that the model has no more prior history-of-success information about overshoot than undershoot, and, under the Bayesian framework, $\alpha + \beta$ indicates the equivalent number of data that the prior information represents. The second constraint is motivated by the complementarity of the two solution methods in the BST: Subjects can see only one approach, undershoot or overshoot, succeed on a given problem, and they must see at least one succeed before going on to the next problem. The third constraint merely places a restriction on the amount of noise added to the PPS estimates. This is a small amount of noise relative to the estimated PPS values for this experiment which are generally in the range of 0.3–0.9, with mean 0.62 and standard deviation .0985. Together, these constraints imply that $\alpha_o = \beta_u$ and $\alpha_u = \beta_o$, and they reduce the free parameters in the model to α_u and α_o .

We searched this two-dimensional parameter space using Powell's method (Press, 1992) to fit the percentages of subjects in each condition using overshoot on each problem. The numbers of overshoot and undershoot successes and failures preceding each problem (averaged by condition) were used to compute the model's predicted percentages (see Appendix B). The parameter values that minimized the SSE between observed and predicted percentages were $\alpha_u = 4.02$ and $\alpha_o = 3.21$. Thus, $\alpha_u/(\alpha + \beta) \approx .56$ and $\alpha_o/(\alpha + \beta) \approx .44$. These proportions indicate that the model's initial estimate for the success of undershoot is slightly higher than that of overshoot and suggest that subjects might have come into this experiment with a slight bias towards undershoot. Also, the absolute magnitude of the estimate for $\alpha + \beta$ is low, indicating that each new experience of success or failure with a production leads to a big change in the history-of-success estimate for that production.

We compared the model's predicted percentages of overshoot use based on these parameter estimates with the observed percentages for all problems and all conditions in the experiment. The predicted percentages are displayed in the bottom panel of Fig. 2. Regressing the observed values on the predicted values suggests that the two match fairly well: predicted = 4.18 + 0.981 * observed, $R^2 = .83$. Note that this best-fitting line has a y-intercept that is not different from 0, $t(27) = .69$, and a slope that is not different from 1, $t(27) = .22$.

Latency results. Yet another measure of solvers' performance that our model can predict is the time to make the first move of each problem (i.e., to apply overshoot or undershoot). The top panel of Fig. 5 plots subjects' average latencies to make the first move for each problem in the training

phase. The best-fitting power function, also plotted in the top panel of Fig. 5, has $R^2 = .89$. A χ^2 goodness-of-fit test suggests that the data deviate significantly from this curve, $\chi^2_6 = 65.12$, $p < .01$.⁶

This organization of the data, however, does not necessarily represent solvers' speedup in the most theoretically appropriate way. In our model, overshoot and undershoot are represented as two separate productions, and the ACT-R theory posits that productions' latencies should speed up as a power function of their practice. Since undershoot and overshoot are being selected differentially across these problems, problem number is not a good approximation for subjects' amount of practice with each production. The bottom panel of Fig. 5, then, presents subjects' latencies to make an undershoot or overshoot move plotted against the number of times the corresponding production had previously been used. For example, all subjects' latencies to apply overshoot when they were using overshoot for the third time are combined with all the latencies to apply undershoot when it was being used for the third time, and the average is plotted against the value "3". The abscissa on this graph only goes to seven "uses" because the number of data points contributing to each average became too small to be reliable after that point. The best-fitting power function, also plotted in the bottom panel of Fig. 5, has $R^2 = .99$ and does not deviate significantly from the data according to a χ^2 goodness-of-fit test calculated as above, $\chi^2_5 = 3.86$, n.s. This suggests that the systematic speedup in solvers' latencies is, as our model predicts, a function of amount of practice with each production and not necessarily a function of overall practice with the task.

Another interesting feature of these latency data is that solvers' speedup as a function of production practice is approximately equal for the two productions, even though solvers' success with overshoot and undershoot was very different. In particular, when we fit subjects' overshoot and undershoot latencies separately as a function of how much each had been practiced, the decay parameters for the two power functions were very similar, -0.59 and -0.42 . That is, during the training phase, all subjects were experiencing 100% failures with undershoot and 100% successes with overshoot, and yet they were becoming equally fast at applying both operators, given equal practice. These results suggest a dissociation between the time to apply a production and its likelihood of being selected: The speedup in production application can be described as a function of the total number of uses of the production, whereas the likelihood to select a production can be described as a function of its history of success. Separate effects of success and practice is a central feature

⁶ For all goodness-of-fit tests on latencies, we take the squared deviations between the observed means and the power curve's predictions, divide each by a measure of the variability of the data points contributing to the mean, and sum over trials (e.g., problem number or practice number). Due to the repeated measures in this analysis, we use the standard mean squared error of the subjects \times trials interaction as our measure of variability.

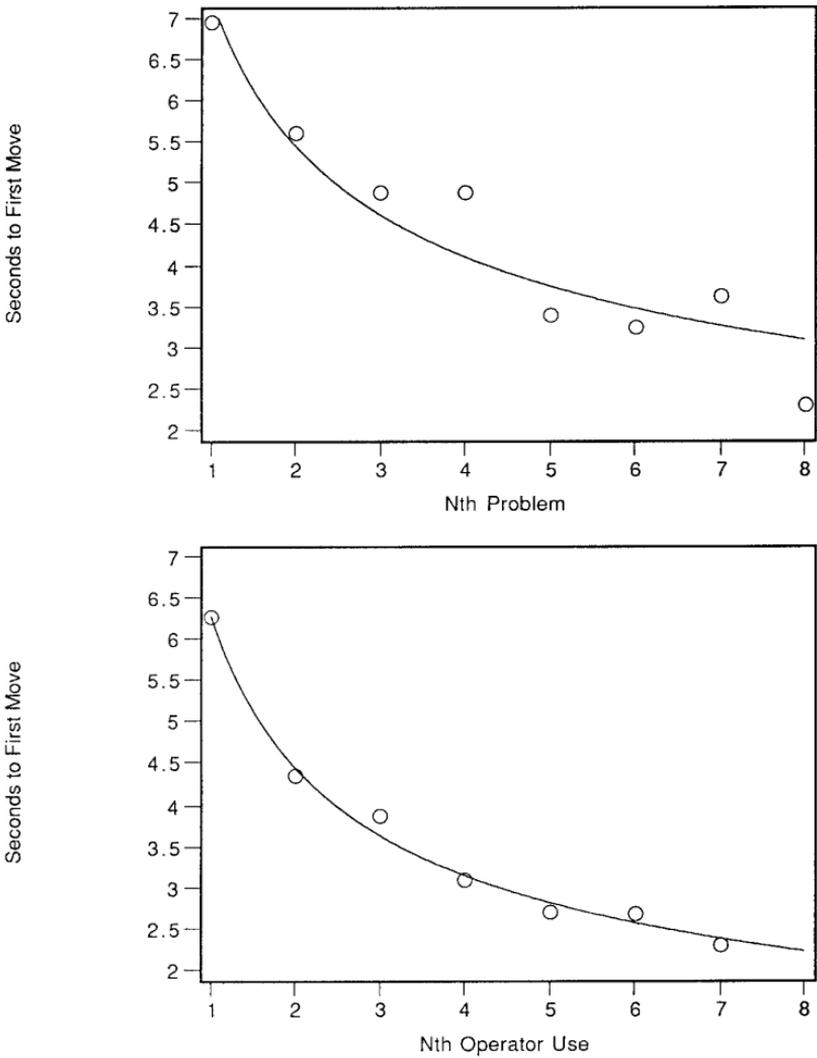


FIG. 5. Mean latency to make first move by number of training problems (top panel) and by number of uses of the production (bottom panel) in Experiment 1. For the top panel, the curve fitted to the data is $y = 7.2x^{-.41}$, $R^2 = .89$, and for the bottom panel, $y = 6.27x^{-.49}$, $R^2 = .99$.

of the current model and of the ACT-R theory. This contrasts with previous ACT theories such as ACT* (Anderson, 1983), which held that information regarding a production's success and overall use could be collapsed into a single index that affects both selection and latency.

Summary

The results of Experiment 1 provide preliminary support for five important predictions of our model of operator selection for the BST. First, the model

predicts that because of the gradually increasing history of success of overshoot relative to undershoot across the training phase, subjects' likelihood of using overshoot will also gradually increase. This trend was exhibited in subjects' selection data during the training phase, most clearly when the percentage of subjects using overshoot was adjusted to account for the different biases of the problems in the training phase. Second, subjects demonstrated an Einstellung effect for overshoot by their greater likelihood of using it on the first "u" problem of the variable phase as compared to subjects who had not received the training problems. In the model, this effect is predicted because of the increase in the history of success of overshoot that occurs during the training phase. Third, current problem features exerted a significant impact on subjects' operator selections because of the distance-to-goal information they convey. With problem bias defined as the difference between the hill-climbing distances for overshoot and undershoot, it accounted for almost 50% of the variance in subjects' selection data. Fourth, subjects in the three conditions exhibited different tendencies to use overshoot during the test phase that were consistent with the histories of success they had experienced during the variable phase. Finally, the model predicts that latencies should decrease as a power law of practice with each production. We found that analyzing latencies by the amount of practice subjects had with each production (instead of practice with the task overall) resulted in excellent power curve fits. Similar results have been obtained by Delaney, Reder, Ritter, and Staszewski (1994) and by Rickard (1994). In general, evidence from this experiment, that the absolute number of uses of a production and the proportion of successes with that production affect different performance measures, supports the ACT-R dissociation between production practice and history of success.

EXPERIMENT 2

The results of Experiment 1 support our model by showing that operator selection in the BST is influenced by the distance to the goal (measured from the resulting state of a potential move) and by the history of success (of the production involved in the potential move). In particular, our model's history-of-success component accounted for the Einstellung effect found in these data. However, some of the results from Experiment 1 are open to another explanation. For example, the finding that subjects receiving "u" problems were subsequently less likely to use overshoot than those receiving "U" problems could be explained if subjects were merely using the operator that worked on the most recent, similar problem. Recall that "U" problems were very different from the rest because all three building sticks were shorter than the desired stick. This could have led "U" subjects to treat their variable problems very differently or to note the transitions between the phases. Our model, in contrast, predicted differences between the conditions because of their different *overall* histories-of-success on all past problems, regardless of

the similarity to the current problem. The simplicity of Experiment 1 did not permit a distinction between these two explanations. More generally, however, it did not allow for a thorough investigation of how solvers use and combine distance-to-goal and history-of-success information.

The goal of the next two experiments, then, was to generalize the results of Experiment 1 to a situation where (i) both the history-of-success and distance-to-goal variables were systematically manipulated (without the confound of major differences in problem similarity) and (ii) the artificial "blocks" of overshoot and undershoot success in different phases were not required. These differences serve both to rule out the alternative explanation described above and to test the finer predictions of our model regarding information combination in a more natural problem-solving sequence.

One of those predictions, based on the Central Equation, is that subjects will generalize from the experience of an operator's success or failure on a particular problem (or problem type) to any new problem (of different type)—as long as the production corresponding to that operator is still applicable. That is, our model not only predicts that operator selections will be influenced by solvers' histories of success with each production, but that the influence will apply to new problems regardless of their similarity to the problem(s) on which past successes or failures occurred. This prediction stems from the fact that, in our model, information from a solver's history of experience is maintained at the level of productions, and so it will have an effect at the level of productions. For example, consider the following scenario: Assume that the model initially estimates that the history-of-success for overshoot and undershoot are approximately equal. At this point, neither production has an overall advantage. Next, suppose that the model solves several problems that are all biased towards undershoot but for which overshoot leads to more successes than undershoot. The histories-of-success of overshoot and undershoot will then change according to the number of successes and failures that occurred, leaving overshoot at an advantage in subsequent operator selections. The important point to note is that these new history-of-success values do not represent the fact that the successes and failures occurred on a particular type of problem. This implies that the model's shift towards selecting overshoot will occur for all new problems, regardless of whether they are biased towards undershoot (same type) or biased towards overshoot (new type). In other words, even when a production's successes and failures are focused on a problem with a particular bias, they will have a broad spectrum effect on all subsequent problems for which the production is applicable.

This generalized effect of history-of-success is a fairly nonintuitive prediction. It is quite possible that subjects who experience an operator's success or failure on a problem with a certain distance-to-goal would instead change their tendency to use that operator only on problems with similar distances. For example, models that perform problem-specific learning (e.g., Hammond, 1986; Logan, 1988; Nosofsky, 1984, 1986) might predict that experiences

associated with a certain type of problem will have a larger impact on similar future problems than on dissimilar future problems. This is because instance-based models, which store an almost complete representation for every problem encountered (e.g., its features, the operators used, the outcomes), solve new problems by reference to previous problems, and the more similar the previous problem, the greater its influence (e.g., Logan, 1988). In contrast, our model records history-of-success information independent of the particulars of the problems on which past success or failure occurred, so it does not predict an interaction between history-of-success and problem-specific information like distance-to-goal. The current experiment was designed to test our model's strong predictions about the generalized, independent influence of history-of-success information on operator selection.

Method

Subjects. Subjects in this experiment were 80 Carnegie Mellon University undergraduates. Fifty-four were recruited from the department subject pool and received course credit for participating, and 26 were recruited from an advertisement posted on a computer bulletin board and were paid \$5 for participating. Subjects from the two populations were approximately equally distributed across the various conditions of the experiment.

Design. The two main conditions in this experiment (biased and neutral) differ according to the type of problem on which failures were designed to occur. In the biased condition, subjects were led to experience failures on problems that were biased toward one operator (according to both hill-climbing metrics defined above) but solved by the other. We call these "false" problems, and they occurred once out of every three problems the biased subjects had to solve. The "false" problems were expected to lead subjects to experience a failure of the operator towards which they were biased and a success of the operator by which they were solved. In the neutral condition, subjects were led to experience failures on problems that were neutral (i.e., undershoot and overshoot resulted in problem states equally close to the goal) but could only be solved by one of the operators. These neutral problems occurred once out of every three problems the neutral subjects had to solve. They were expected to lead to failure of the operator that would not solve them 50% of the time. The remaining two-thirds of the problems were the same for all subjects. We call these "true" problems because they were solved by the operator toward which they were biased. Half of these were solved by overshoot and half by undershoot. Subjects were expected to experience a success (but no failure) on each "true" problem.

For each subject (biased or neutral), a single operator (overshoot or undershoot) was designated as the "failing" operator—the operator that would not solve the "false" or "neutral" problems. The other operator, then, was designated the "nonfailing" operator. The assignment of overshoot and undershoot to the "failing" and "nonfailing" operators was counterbalanced

across subjects. Note that subjects were not forced or guaranteed to experience failure of the failing operator. This label merely indicates our intended manipulation and expectation of what they would experience.

Procedure. This experiment was divided into several phases. First, subjects were introduced to the task and interface in an instructional phase. The experiment itself was composed of nine alternating “test” and “solve” phases: test 0, solve 1, test 1, solve 2, test 2, solve 3, test 3, solve 4, test 4. After these phases, subjects were asked to complete a postexperimental questionnaire concerning their self-reported strategies.

The instructional phase began with a description of the BST. Subjects read some text on the computer screen, practiced each of the mouse-clicking actions for building sticks, and watched as the computer solved two sample problems (one by undershoot and one by overshoot). Subjects were prompted to ask the experimenter any questions they might have at the end of the instructional phase before continuing.

The first experimental phase was test 0. All the test phases were the same for all conditions. (The one exception to this was that half of the subjects in each condition did not receive test 1 so we could compare performance with and without this test. Analyses revealed that the presence or absence of test 1 did not make a difference in subjects’ later performance, suggesting that the test phases were only a measure of and not an influence on operator selections. All comparisons presented below collapse over this factor.) Each was preceded by a screenful of text that instructed subjects to click on the stick that they would choose first if they were solving the problem. Subjects were then presented with ten BST problems in a random order. All test problems had one building stick longer than the desired stick and two building sticks shorter than the desired stick, so overshoot and undershoot were always both applicable. Two test problems were from each of the following categories: high undershoot bias, low undershoot bias, no bias (neutral), low overshoot bias, and high overshoot bias. (Recall that problem bias is a measure of the distance-to-goal after an overshoot move relative to the distance-to-goal after an undershoot move.) Subjects clicked on one of the three building sticks at the bottom of the screen to indicate their first move. After one such (legal) mouse click, the screen was immediately erased, and a new test problem was presented. Thus, subjects had to indicate which operator they would use to solve the problem but were never allowed to solve to completion. The stick selected and the time to click were recorded for each test problem.

There were also four solve phases in the experiment, with each preceded by a screenful of text that stated, “For the next set of problems, you should *solve* all the way . . .” Each solve phase consisted of five problem triplets, where a triplet includes one true undershoot problem, one true overshoot problem, and one “false” or “neutral” problem. The order of the five triplets within a solve phase was randomized as was the order of the problems within a triplet. Subjects were required to work on each problem until they solved

TABLE 1
Mean Number of Problems on Which Failures Occurred in Experiment 2

Condition	Solve phase			
	1	2	3	4
Biased				
Failing operator	4.5	4.0	3.3	3.7
Non-failing operator	1.3	1.5	1.4	2.0
Neutral				
Failing operator	2.9	2.5	2.5	2.5
Non-failing operator	1.0	1.3	1.3	1.3

Note. Each solve phase in Experiment 2 included a total of 15 problems.

it or until they had taken 50 steps, whichever came first. (In this experiment, subjects always managed to solve the problems before the 50-step limit.) Each click and its latency were recorded.

In the final post-experimental phase, subjects completed a questionnaire regarding their self-reported strategies. In addition, subjects were asked if the experiment had reminded them of any experiments they had read or learned about in class. None of the subjects were reminded of Luchins's water jars experiments.

Results and Discussion

Solve phase results. During the solve phases of the experiment, subjects solved different problems that were designed to lead them to experience certain profiles of success with the two operators, overshoot and undershoot. Since subjects were required to solve each problem to completion before they could go on, we had control over their proportions of successes. It is necessary, however, to verify that subjects were experiencing operator failures in accordance with the intended manipulations. Table 1 presents the average number of failures subjects experienced with both the operators for the biased and neutral conditions separately. These data confirm that subjects experienced more failures with the designated "failing" operator than with the "non-failing" operator, $F(1, 39) = 4.65$, $p < .05$, $MSe = 15.9$. In addition, they suggest that subjects in the biased condition experienced more failures with the failing operator than did subjects in the neutral condition, $F(1, 78) = 18.9$, $p < .01$, $MSe = 13.3$. This is expected since biased subjects were expected to experience a failure on almost every "false" problem, whereas neutral subjects were expected to experience a failure on approximately half of the neutral problems.

Although the test phases were specifically designed to test the effects of these different histories of success on operator selection, we can also compare

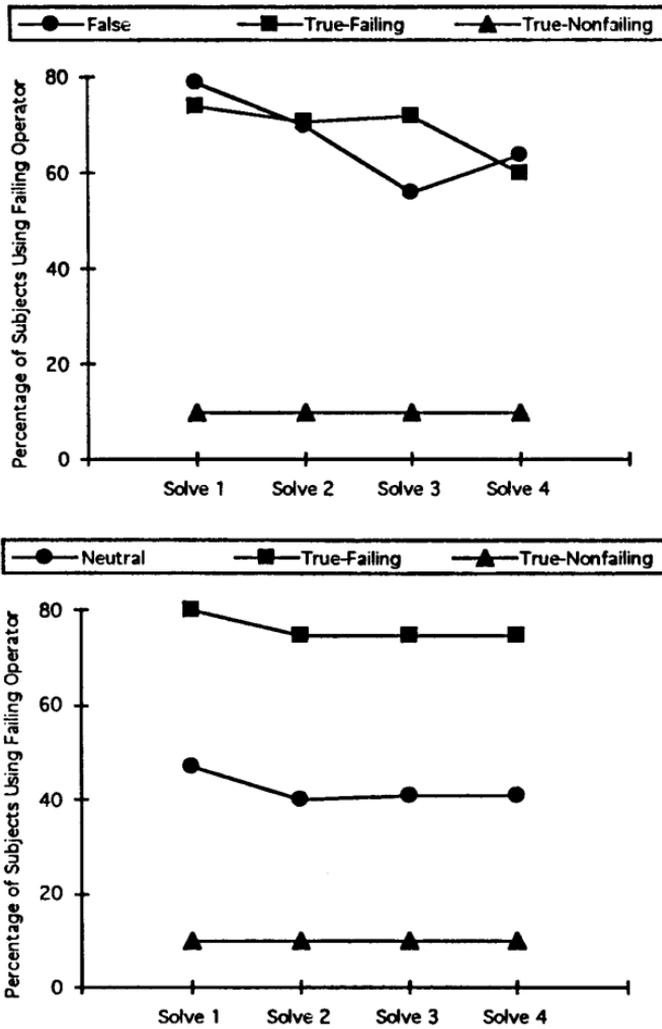


FIG. 6. Percentage of subjects in the biased condition (top panel) and neutral condition (bottom panel) using the failing operator for each problem type and solve phase in Experiment 2.

subjects' selection data across the four solve phases to get an initial idea of the effects. Recall that there were two types of "true" problems in each solve phase: true overshoot and true undershoot. These two true- x types can be relabeled as "true-failing" when " x " was the operator designed to fail and "true-nonfailing" when " x " was the operator designed not to fail. Figure 6 presents the percentage of subjects in the biased condition (top) and neutral condition (bottom) using the "failing" operator on each problem type.

There are several interesting features in these data. Across the four solve phases, the biased subjects exhibited a decrease in their likelihood to use the failing operator on both the true-failing and the false problems, $F(3, 117) =$

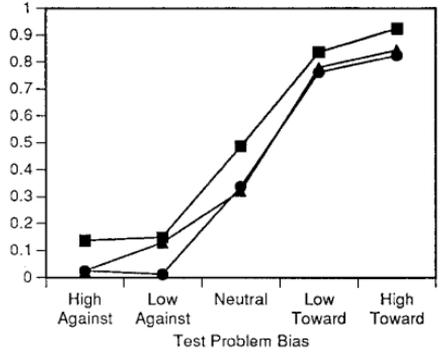
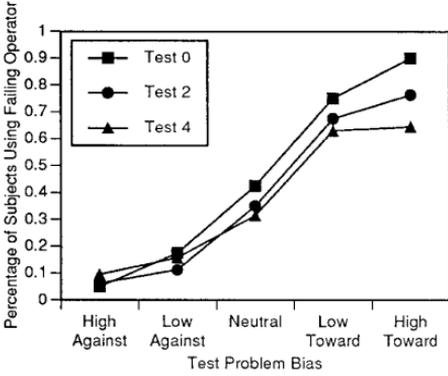
11.9, $p < .01$, $MSe = 0.03$. For these subjects, both problem types are biased toward the failing operator, and yet the true-failing problems are actually solved by the failing operator, whereas the false problems are not. As long as subjects cannot tell the two problem types apart, the model predicts that subjects' past experience of more failures with the failing operator, which occurred mainly on the false problems, will lead to a lower likelihood of using the failing operator on both problem types across solve phases. On the true-nonfailing problems, the biased subjects' use of the failing operator is low and steady across the four solve phases, $F(3, 117) < 1$, n.s., $MSe = 0.013$. Since these problems were biased against the failing operator, this result may represent a floor effect; subjects started with a low likelihood of using the failing operator on these problems, and it could not get lower.

For neutral subjects, use of the failing operator on neutral problems is approximately 50%, with no reliable change across the solve phases, $F(3, 117) < 1$, n.s., $MSe = .05$. Neutral subjects' selection of the failing operator on the true-failing problems did not exhibit a significant decrease either, $F(3, 117) < 1$, n.s., $MSe = .13$. However, when these data are combined with those of the biased subjects, a significant difference across solve phases is found for both of these problem types (true failing and neutral/false), $F(3, 237) = 10.3$, $p < .01$, $MSe = .04$, with no interaction between condition (biased vs neutral) and solve phase, $F(3, 237) = 2.00$, n.s., $MSe = .04$. This suggests that the change in biased and neutral subjects' behavior across the four solve phases could not be distinguished on these two problem types even though the biased data alone exhibited a significant change and the neutral data did not. Also, for the true-nonfailing problem type, the neutral subjects displayed the same behavior as the biased subjects: Their use of the failing operator on problems biased against that operator started out low and remained low across all four solve phases.

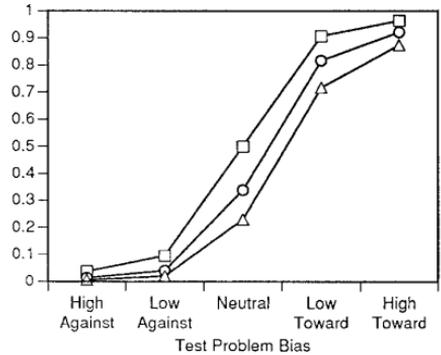
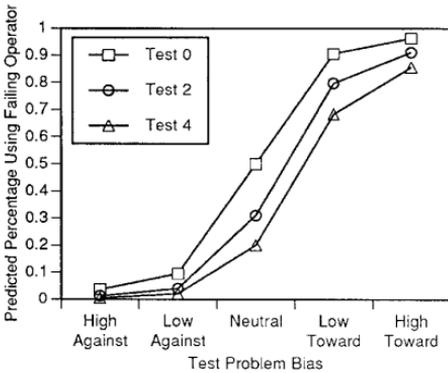
Test phase results. Figure 7 presents subjects' likelihood to use the failing operator at test 0, test 2, and test 4 for each of the five test problem types for the biased conditions (top left panel) and the neutral conditions (top right panel). (Test 1 was excluded from the figure because not all subjects were given this test phase, and test 3 was excluded so that the number of solved problems between the displayed test phases would be equal. The analyses below, however, include subjects' selection data from all the test phases except test 1.) Note that the test problems with bias "High Toward [the failing operator]" represent the type of problem on which the biased subjects experienced many of their failures, and test problems with "Neutral" bias represent the type of problem on which the neutral subjects experienced many of their failures.

A repeated measures ANOVA with the factors of test number (tests 0, 2, 3, 4), test problem bias (high-toward, low-toward, none, low-against, high-against), and condition (biased, neutral) was performed on these data. The analysis revealed that use of the failing operator decreased significantly across

Observed Percentages



ACT-R Predictions



Example-Based Predictions

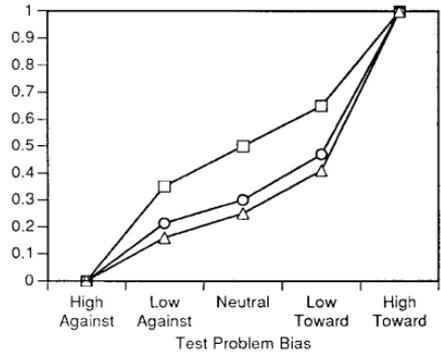
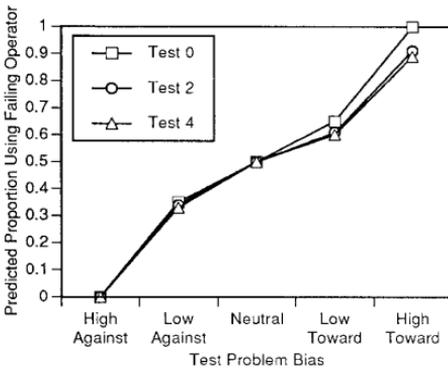


FIG. 7. Observed and predicted percentages of subjects using the failing operator on tests 0, 2, and 4 in the biased condition (left panels) and neutral condition (right panels) in Experiment 2.

tests, $F(3,225) = 12.263$, $p < .01$, $MSe = .076$, and increased significantly with increasing bias toward that operator, $F(4,300) = 281.276$, $p < .01$, $MSe = .129$. These two main effects reflect the history-of-success and distance-

to-goal predictions of our model, respectively. Our model also predicts a larger decrease across test phases for the biased condition than the neutral condition (due to the former's experience of more failures), but this interaction was not significant, $F(3,225) = 1.113$, n.s., $MSe = .076$.

The major prediction of our model (see Central Equation) is that the effects of history-of-success and distance-to-goal will be independent—that subjects will generalize their experiences of success and failure on one problem type to all problem types. This prediction stems from the model's representation of history-of-success information without the specific context of each success; it can be evaluated by testing the interaction between test number and test problem bias. This interaction is not significant in the complete data set, $F(12, 900) = 1.236$, n.s., $MSe = .070$, nor for the biased subjects alone, $F(12, 444) = 1.346$, n.s., $MSe = .08$, nor for the neutral subjects alone, $F(12, 456) < 1$, n.s., $MSe = .06$. (The powers of these two tests are .88 and .94, respectively, against an alternative generated by an alternate model described in the General Discussion. These values suggest that the tests would likely have revealed an interaction effect of that size if one existed.) However, when the analysis is restricted to the biased subjects' Test 0 and Test 4 data, this interaction reaches significance, $F(4, 148) = 3.0$, $p < .05$, $MSe = .08$. This is the only example of lack of independence in this experiment, and we will come back to it later.

Latency results. As with Experiment 1, we can examine subjects' latencies to make the first move on each problem. Grouping these latencies according to how many times the corresponding production had previously been practiced, we see a decrease in latencies that conforms fairly well to a power law (see Fig. 8). Recall that our model predicts a power-law decrease in latencies with number of production uses. The sum of the weighted deviations between these mean latencies and the best-fitting power curve is not significant, $\chi^2_{28} = 32.8$, n.s. (see Footnote 6).

Note that latencies for both the failing and nonfailing operators are included in Fig. 8 and that, when fit to separate power curves, the parameters for the exponent are almost the same, $-.123$ and $-.155$, respectively. As in Experiment 1, we have evidence of a dissociation between changes in the likelihood and speed of use of productions because one variable, history of success, influences how likely a production is to be selected while another variable, amount of use, influences how quickly it is to be executed.

Modeling Results

In this section, we present the operator-selection tendencies of the model when it is fit to the test phase data presented above. Because the next experiment has the same structure as this one, we combine the two data sets to estimate the model parameters. In this way, we show that our model, with a single parameter setting, can accommodate subjects' selection data from both experiments.

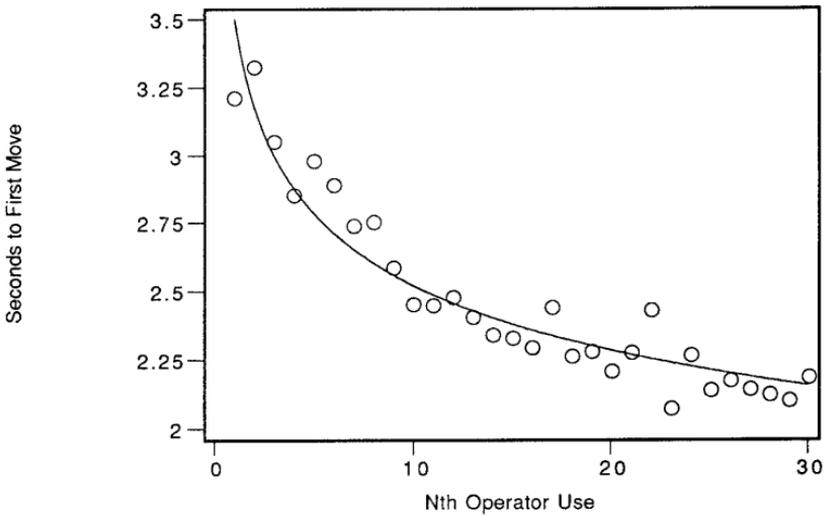


FIG. 8. Mean latency to make first move by number of production uses in Experiment 2. The curve fitted to the data is $y = 3.5x^{-1.4}$, $R^2 = .91$.

We obtain estimates for the parameters of the model by taking the average number of successes and failures experienced by subjects in each condition (from Experiments 2 and 3) as the input and the percentages of subjects in each condition using the failing operator on each test problem for each test phase as the output. The fitting procedure is the same as used in Experiment 1 (see Appendix B). With the same constraints employed there, our model requires two free parameters. The best-fitting values are $\alpha_o = 78.33$ and $\alpha_u = 80.46$. These values indicate that the prior estimates of probability of success for overshoot and undershoot are essentially equivalent: $\alpha_o/(\alpha + \beta) \approx .50 \approx \alpha_u/(\alpha + \beta)$. The larger magnitude of these parameters compared to those for Experiment 1 implies each experience of success and failure in this experiment had a smaller effect. Since Experiments 2 and 3 have many more trials than Experiment 1, this difference in parameter values may reflect a rapidly decaying effect of past experience, such that only recent successes and failures impact the model's history-of-success component. (See Elliott & Anderson's (1995) treatment of the accumulation of past experience.)

When these parameter estimates are used to calculate the model's predicted selections for Experiment 2, we achieve a good fit to the observed data. (See the middle panels of Fig. 7.) When predicted percentages are regressed on observed percentages, the best-fitting line is: predicted = $-7.41 + 1.18 \times$ observed, $R^2 = .97$. (Note that the predicted and observed values are on a 0–100 scale.) This level of fit is attained because the predicted percentages exhibit both the main effects of test number and test problem type with very little evidence of an interaction. Although the y-intercept is significantly

different from 0, $t(28) = 3.54, p < .01$, and the slope is significantly different from 1, $t(28) = 4.61, p < .01$, both values are relatively close to these ‘‘exact-fit’’ values. The fact that the slope is greater than 1 suggests that the predicted percentages are slightly more extreme than the observed percentages. This could represent solvers’ (but not the model’s) avoidance of extreme operator-selection tendencies.

The one place where the model does not fit well is on the ‘‘high-toward’’ test 4 data for the biased condition. Since this is the problem type on which biased subjects had seen the failing operator fail, it seems that subjects may be exhibiting some problem-specific learning. That is, a more drastic shift in operator selection on this problem type is consistent with subjects avoiding the failing operator more on test problems that were similar to the problems on which they experienced the failures. Upon further exploration of the problems solved by the biased subjects, we found that the ‘‘false’’ problems that half of them solved had the idiosyncratic feature that the desired stick’s length was always rather long. This was true only for the half of the biased subjects for whom overshoot was the designated failing operator. It is possible that these subjects were learning that the long desired stick was a cue to solve the overshoot-biased problems by undershoot. By coincidence, the ‘‘high-toward’’ test problems had that same feature, whereas the other test problem types did not. Thus, if some of these subjects were responding to this problem feature, they would be differentially biased against the failing operator on the ‘‘high-toward’’ test problems.

Selection data from the solve phases provide evidence that these subjects were exhibiting some unpredicted learning. Their tendency to select overshoot on true-overshoot problems stayed constant (and high), whereas it decreased with experience for the false-overshoot problems. These two types of problems were designed to be equivalent in bias and thus are indistinguishable to our model. And yet, the interaction between solve phase and problem type is significant in the data, $F(3, 57) = 4.09, p < .05, MSe = 0.3$, suggesting that these subjects were learning to discriminate between the false-overshoot problems and the true-overshoot problems. This is suggestive of some feature-specific or problem-specific learning. Indeed, since this unexpected learning was limited to a particular subset of the problems that had an unusual feature, it may be well described as the learning of exceptional cases, which occurs in addition to more general production-based learning. Our model explains the majority of subjects’ selection data but would require some modification to account for this additional kind of learning. For example, if we allowed our model to learn extra productions after experiencing repeated failures on especially salient problems, it could exhibit some problem-specific learning. That modification might be similar to the RULEX model (Nosofsky *et al.*, 1994), which learns general rules for categorization and then specifies exceptions to those rules, or to the ASC-M model (Siegler & Shipley, 1994), which incorporates both global and problem-specific learning.

EXPERIMENT 3

All but one of the results in Experiments 1 and 2 supported the basic claims of the current model. Solvers' operator selections were found to be influenced by both distance-to-goal and history-of-success information, and, most importantly, an independent combination of these two variables matched subjects' selection tendencies well, even when their failures had been focused on one problem type. Experiment 3 provides a replication of Experiment 2 with several modifications. First, we substituted an extreme-biased condition for the neutral condition. The extreme-biased condition was similar to the biased condition in that subjects were led to experience failures of one operator on problems biased toward the opposite operator, but the proportion of failures was much higher. The extra failures in the extreme-biased condition created a more extreme history of success than had previously been tested. Our model predicts that such a history of success will lead to even greater shifts in operator selections. Second, we modified some of the problems so that they were less likely to introduce exceptional cases like those we believe led to some feature-specific learning for half of the biased subjects in Experiment 2. Even with this change in stimuli, however, it is important to note that problem-based learning could still occur and, as in Experiment 2, was promoted by the fact that subjects experienced their failures on repeated problems of the same type. If subjects were storing past problems and referring to them exclusively to solve new problems, they should exhibit more of a change in selection tendencies on problems of the same type. On the other hand, if subjects were learning the overall success of the two operators, represented in separate units of procedural knowledge as in our model, we would expect more global shifts in their operator-selection tendencies.

Method

Subjects. Subjects were 38 Carnegie Mellon University undergraduates participating for credit or \$5. Subjects receiving credit versus payment were approximately equally balanced across the conditions.

Design. The design of this experiment was like that of Experiment 2, except that an "extreme-biased" condition was substituted for the neutral condition. Subjects' conditions determined the problems they were given to solve. In the biased condition, subjects received problem triplets each consisting of one true-undershoot, one true-overshoot, and one "false" problem. In the extreme-biased condition, subjects received problem sextuplets that each consisted of three true problems (two biased toward the "nonfailing" operator and one biased toward the "failing" operator), and three "false" problems. In both conditions, the "false" problems were biased toward the "failing" operator but only solved by the "nonfailing" operator. Thus, the expected proportion of failures of the "failing" operator was higher in the extreme-biased condition (3/4) than in the biased condition (1/2), leading to an even larger disparity in the history of success of the two operators.

TABLE 2
Mean Number of Problems on Which Failures Occurred in Experiment 3

Condition	Solve phase		
	1	2	3
Biased			
Failing operator	8.0	7.3	7.2
Non-failing operator	3.0	4.3	4.4
Extreme biased			
Failing operator	10.0	7.8	6.3
Non-failing operator	2.4	2.9	3.2

Note. Each solve phase in Experiment 3 included a total of 30 problems.

Procedure. The same instructional phase was included as in Experiment 2. After that phase was completed at the subjects' own pace, the experimental trials began. The sequencing of this experiment followed the alternating test-phase, solve-phase pattern of Experiment 2. However, there were 30 problems in each solve phase and only three solve phases in total. The sequence of phases, then, was test 0, solve 1, test 1, solve 2, test 2, solve 3, test 3. The procedure for the test phases was identical to Experiment 2.

Results and Discussion

As in Experiment 2, the main question to be addressed in this experiment was: Does the effect of success or failure on one problem type generalize across all problem types? This experiment also provided an additional test of the history-of-success component of our model by comparing operator selections of biased and extreme-biased subjects who experienced different proportions of successes. The same dependent measures will be analyzed here as were in Experiment 2: solve-phase selections, test-phase selections, and latencies.

Solve phase results. During the solve phases of the experiment, subjects in the different conditions solved different problems that were designed to lead them to experience certain profiles of success with overshoot and undershoot. Since subjects were required to solve each problem to completion before they could go on, we had control over their proportion of successes. It is still necessary, however, to verify that subjects were experiencing failures in accordance with their condition. Table 2 presents the average number of failures subjects experienced with the designated "failing" operator and with the designated "nonfailing" operator, for the biased and extreme-biased conditions. These data confirm that subjects experienced more failures of the "failing" operator than the "non-failing" operator, $F(1, 36) = 490, p < .01, MSE = 8.9$. In addition, although the extreme-biased subjects' total number

of failures was not significantly greater than the biased subjects', $F(1, 36) < 1$, n.s., $MSe = 8.7$, the former group did exhibit more of a change in the number of failures of the failing operator as evidenced by the interaction of condition and solve phase for the failing operator data only, $F(2, 72) = 4.2$, $p < .05$, $MSe = 5.3$. This is consistent with the intended manipulation because subjects in the extreme-biased conditions experienced a greater disparity in the success rates of the failing and nonfailing operators, leading them to change their selection tendencies more sharply. Note that at the first solve phase, subjects in the extreme-biased conditions experienced more failures of the failing operator than did the subjects in the biased conditions, $t(36) = 1.80$, one-sided $p < .05$.

With the failure manipulation verified, we can compare subjects' selection data across the four solve phases to get an initial idea of the effects. The top panel of Fig. 9 presents the percentage of subjects in the biased conditions using the failing operator on each of the three problem types across the three solve phases, and the bottom panel of Fig. 9 presents the same data for the extreme-biased subjects. Note that in both conditions, subjects exhibited a decrease in their likelihood to use the failing operator across solve phases—on both false and true-failing problems, $F(2, 38) = 5.5$, $p < .01$, $MSe = .03$ for biased subjects and $F(2, 34) = 7.9$, $p < .01$, $MSe = .04$ for extreme-biased subjects. The lack of an interaction between solve phase and these two problem types suggests that subjects' decrease in use of the failing operator was similar for the two problem types, $F(2, 38) < 1$, n.s., $MSe = .03$ for biased subjects and $F(2,34) = 1.3$, n.s., $MSe = .02$ for extreme-biased subjects. Finally, subjects' selection of the failing operator on the true-nonfailing problem types remained low and relatively constant across solve phases for both groups, $F(2, 38) < 1$, n.s., $MSe = .02$ for biased subjects and $F(2, 34) = 1.1$, n.s., $MSe = .01$ for extreme-biased subjects.

One difference between the biased and extreme-biased conditions is the amount of change in operator selections during the solve phase. The extreme-biased subjects exhibited a sharper decrease than did the biased subjects in their use of the failing operator on the true-failing and false problems across the experiment, $t(36) = 2.32$, $p < .05$. This difference is predicted by the model because the extreme-biased subjects experienced a greater proportion of failures of the failing operator, leading to a greater disparity in the two operators' histories of success. Thus, the solve phase data provide additional support for the relationship between history of success and operator selections that is specified by the model.

Test phase results. The top two panels of Fig. 10 present the biased subjects' and extreme-biased subjects' likelihood to use the failing operator during test 0, test 2, and test 3. (Data from test 1 are excluded from these graphs to make the presentation of the data more clear, but the results below include all the test data.) A repeated measures ANOVA with the factors test number (test 0, test 1, test 2, and test 3), problem bias (high-toward, low-toward,

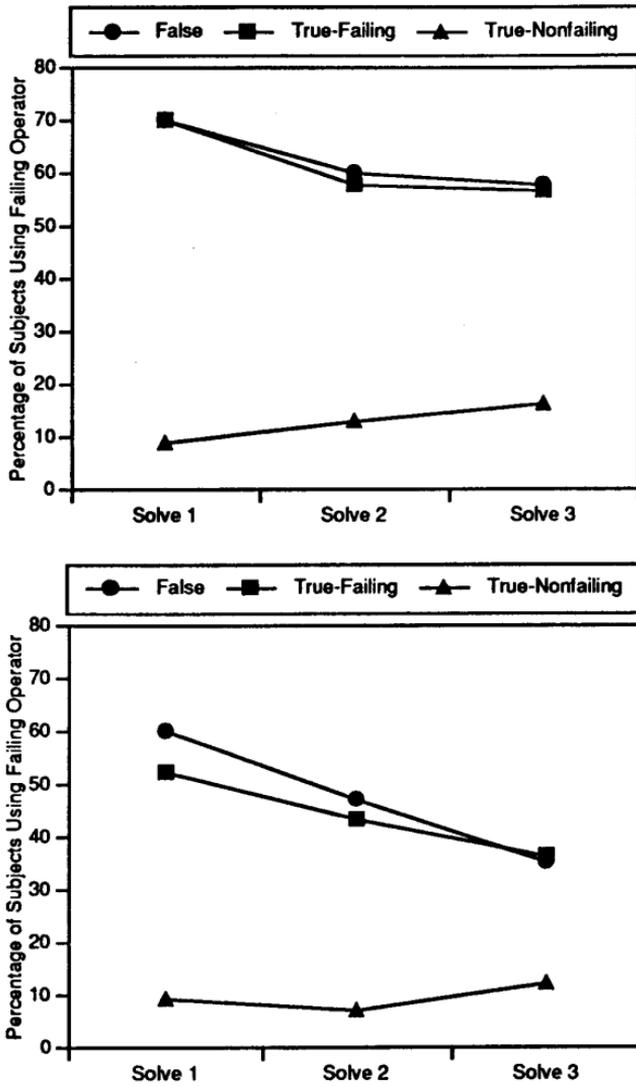
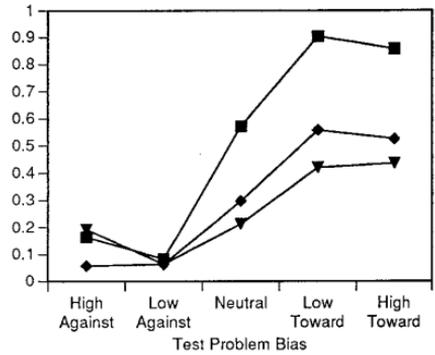
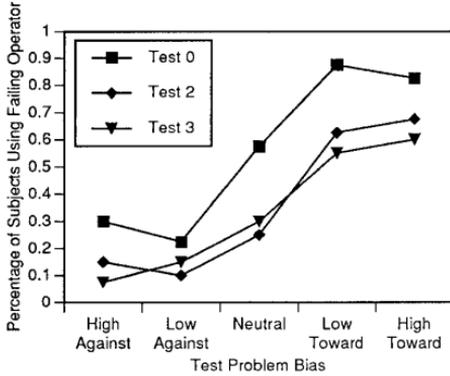


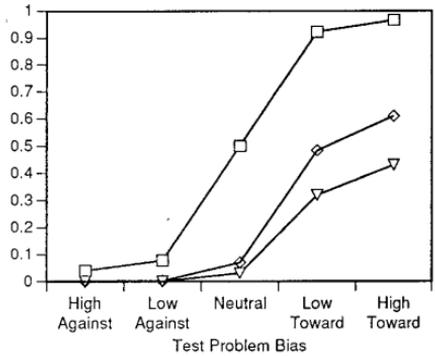
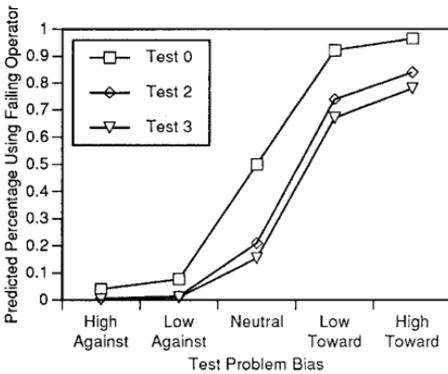
FIG. 9. Percentage of subjects in the biased condition (top panel) and extreme-biased condition (bottom panel) using the failing operator for each problem type and solve phase in Experiment 3.

none, high-against, low-against), and condition (biased, extreme-biased) was performed on these data. The analysis revealed that use of the failing operator decreased significantly across tests, $F(3, 105) = 25.380$, $p < .001$, $MSe = 0.082$. Also, use of the failing operator increased significantly with problem bias, $F(4, 140) = 69.894$, $p < .001$, $MSe = 0.136$. The interaction between these two factors was significant, $F(12, 420) = 3.032$, $p < .001$, $MSe = 0.077$. This interaction, however, is likely due to a floor effect at the high-

Observed Percentages



ACT-R Predictions



Example-Based Predictions

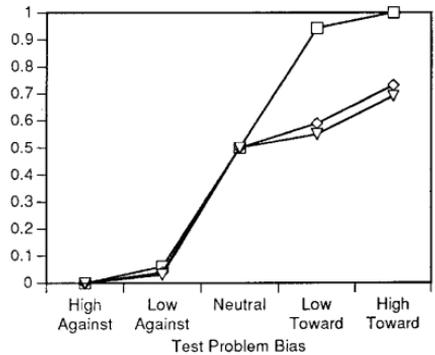
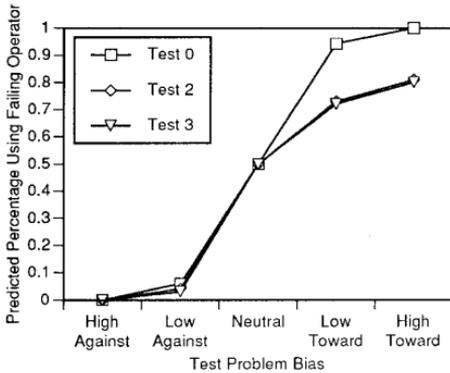


FIG. 10. Observed and predicted percentages of subjects using the failing operator on tests 0, 2, and 3 in the biased (left panels) and extreme-biased (right panels) conditions in Experiment 3.

against and low-against problem bias levels in the extreme-biased condition. If we remove the two “against” types from the analysis, the interaction is no longer reliable, $F(6, 210) < 1$, n.s., $MSE = .094$. Moreover, even without

removing these two problem types from the analysis, the interaction is not significant for the biased subjects alone, $F(12, 228) < 1$, n.s., $MSe = .073$.

The test phase data provide further support for our model's history-of-success component which specifies the quantitative effects of a solver's history of success. In particular, the model predicts that the decrease in use of the failing operator will be greater for the extreme-biased subjects than for the biased subjects because the extreme-biased subjects experienced a greater proportion of failures. Although the interaction of condition and test number was not significant in these data (even after removing the two "against" problem types to remove floor effects), $F(3, 105) = 1.366$, n.s., $MSe = .100$, a specific test that the shift in selection tendencies between Test 0 and Test 3 was greater for extreme-biased than for biased subjects was reliable, $t(36) = 1.78$, $p < .05$, one-sided. Combining this experiment with the results of Experiment 2, we have shown that subjects experiencing failure of an operator on approximately three-quarters of their attempts with it (extreme-biased subjects) decrease their use of that operator more than subjects experiencing its failure on approximately half of their attempts (biased subjects) who, in turn, decrease their use of that operator more than subjects experiencing its failure on approximately one-third of their attempts (neutral subjects). These results hold independent of the type of problem on which the failures occurred (neutral or biased) and independent of the distance-to-goal of the test problem on which selection performance was tested.

Latency results. As with Experiments 1 and 2, we can examine subjects' latencies to make the first move on each problem. Grouping these latencies according to how many times each operator had previously been practiced, we see a decrease in latencies that conforms fairly well to a power law (see Fig. 12). This power function curve does not deviate significantly from the observed latency means, $\chi^2_{38} = 34.32$, $p > .35$ (see Footnote 6).

Again, the model and data both distinguish between the effects of number of uses and proportion of successes: amount of use predicts speed of use and proportion of success predicts likelihood of use. Here, both operators are speeding up equally even though one is becoming more likely to be used and the other less. The power function exponents were $-.173$ for the failing operator and $-.170$ for the non-failing operator when latencies were fit according to number of uses with each. As in Experiments 1 and 2, this provides evidence of a dissociation between changes in likelihood of use and speed of use.

Modeling Results

As with Experiments 1 and 2, we can also evaluate the fit of our model to these selection data. Using the parameters obtained by fitting the combined data sets from Experiments 2 and 3, we calculated the model's predicted percentages of using the failing operator for the biased and extreme-biased conditions of Experiment 3. These predictions are presented in the middle

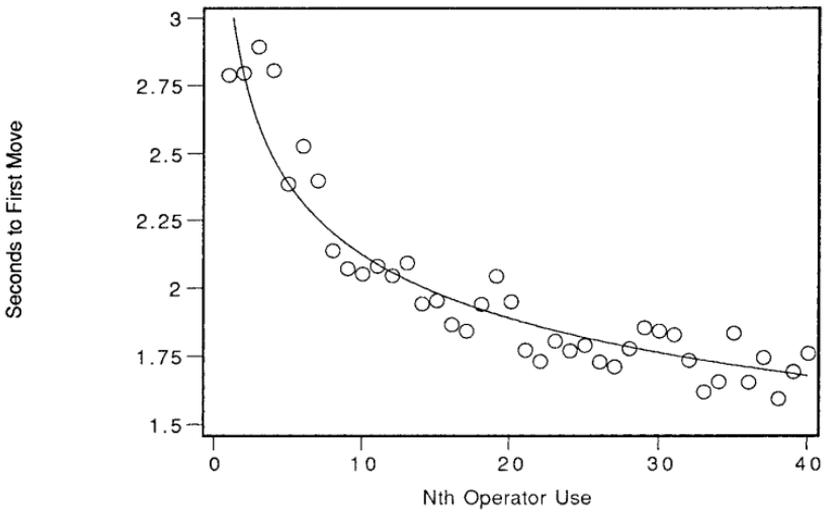


FIG. 11. Mean latency to make first move by number of production uses in Experiment 3. The curve fitted to the data is $y = 3.1x^{-.17}$, $R^2 = .88$.

panels of Fig. 11. Note that the extreme-biased predictions display an apparent interaction due to a floor effect for the high-against and low-against test problem types just as our data did. When the observed percentages are regressed against the corresponding predicted percentages, the best-fitting line is: $\text{predicted} = 6.00 + 0.85 * \text{observed}$, $R^2 = .97$. As was the case in Experiment 2, this line's slope and y-intercept values suggest that the model is not systematically overpredicting or underpredicting and provides a fairly close one-to-one match with the data. Nevertheless, the intercept is significantly different from 0, $t(28) = -5.12$, and the slope is significantly different from 1, $t(28) = 4.47$.

GENERAL DISCUSSION

In the experiments presented above, we found several converging lines of evidence to support the four major claims of our ACT-R model of operator selection: (1) Information about the solver's history of success is stored with each production, and more successful productions are more likely to be selected. Subjects demonstrated their sensitivity to each production's history of success by reproducing the Einstellung effect in Experiment 1 and by changing their operator-selection tendencies in Experiments 2 and 3. (2) Current problem information that measures the distance to the goal has a privileged status in our model: A simple hill-climbing distance metric is computed for each applicable production, and moves that minimize this distance are more likely to be selected. Subjects' sensitivity to this metric was also exhibited in their selection data. (3) The above two sources of information are evaluated and combined independently. In general, when solvers experienced failures

of a particular operator on one problem type, they decreased their use of that operator across all problem types. According to our model (in which each operator has its own production), this is due to the production-based storage of history-of-success information. Our Central Equation that specifies this independent contribution of current and experience-based information is similar to Massaro's fuzzy logical model of perception (Massaro & Cohen, 1991; Massaro & Friedman, 1990; Massaro & Hary, 1986) in that it claims independence both in the evaluation of the separate dimensions and in their integration. Almost all of the tests of this independence prediction supported our model. (4) Different measures of solvers' past use of an operator influence their likelihood to use it vs their speed at using it. The relative proportion of success with an operator influences the likelihood of selecting that operator, but the number of times a solver uses a particular operator determines the speedup in the application of that operator, regardless of its proportion of successes. In all three experiments, subjects exhibited a speedup in their use of the operators being studied, even for the operator that was less successful and less likely to be selected.

These results provide new specifications for the combined influence of history of success and current problem features on solvers' operator selections. More generally, they provide insight into how current and experience-based information are integrated in problem solving. Our basic conclusion is that the processing of both types of information plays an important and quantifiable role in problem solving. More specifically, however, we have shown that these two types of information are combined in an additive fashion when they are quantified in terms of what each predicts for the likelihood of success of different actions. That is, likelihood of success is the single dimension into which both past and current information can be "transformed" and on which solvers seem to be basing their selections. This is a practical (and, dare we say, rational) metric for solvers to use because it should, on average, lead to more frequent success.

Our model of the influence of current and experience-based information on operator selection also uses likelihood of success as the global metric by which it selects among productions. Along with this global metric, the fundamental idea behind our model is that each production represents a separate unit of procedural knowledge. Together, these ideas define the main features of how our model makes operator selections. For example, the production-based representation in our model means that history-of-success information is stored separately for each production and aggregated over all trials on which the same production was attempted. This led our model to select among operators in a way very similar to the way people do. Assuming that a solver's production set captures his or her knowledge about which procedural units are distinct and which are not, it makes sense to aggregate success information at this level. This representation also implies that solvers with a more refined set of productions will

be able to represent history-of-success information at a finer grained level. Finding this kind of correspondence would provide additional support for a production-like representation of success knowledge. Indeed, one could even reverse the logic of this paper and use the independence prediction of the Central Equation as a way of detecting how procedural knowledge was divided into productions. That is, the degree of generalization solvers exhibited after experiencing success or failure of a particular operator in a particular situation would indicate whether that operator was represented as a separate unit of procedural knowledge or not: less generalization of the history-of-success effect implies more separate productions.

As noted above, in almost all cases, solvers' operator selections suggested that they were storing and using history-of-success information at the level of productions in our model and generalizing over individual problems. There was one case (half of the biased subjects at test 4 in Experiment 2), however, in which some problem-specific or feature-specific learning was also exhibited. This raises a general question: When will our independence prediction be maintained and when will more context-specific forms of learning tend to dominate?

The experiments in this paper cannot answer that question, but they do provide some clues to guide our speculation. Independence between history-of-success and distance-to-goal seems to occur when solvers have effective operators (i.e., productions) that are applicable in a fairly wide range of situations. In our experiments, undershoot and overshoot fulfilled these criteria because they both could apply to almost all of the problems. But the BST is not the only task in which this is the case. For example, Reder's (1987) work suggests that the general question-answering strategies of retrieval and plausibility judgment may also lead to an additive combination of history-of-success information with other information that is specific to the features of the question. Problem-specific learning, in contrast, may tend to occur when there is something memorable or discriminable about a particular problem or set of problems. We believe this was the case for a subset of the problems in Experiment 2, and it fits with Logan's (1988) suggestion that more memorable stimuli will require fewer repetitions for instance-based retrieval.

Example-Based Models

Another question that relates to problem-specific learning is: If solvers seem to exhibit problem-specific learning in one case in Experiment 2, could this kind of learning explain our other results also? The class of models called example-based models, which learn exclusively by storing and referring to past instances, provide a useful contrast to the production-based learning exhibited by our model. Example-based models rely on a very different kind of processing that has been shown to be quite versatile and flexible. For example, Medin and Schaffer's (1978) paper showed that example-based models could exhibit prototype effects in categorization, originally a phenom-

enon thought to be a hallmark of abstract processing. So, the question of whether example-based models could account for our operator selection results is an intriguing one.

Unfortunately, no standard example-based model can be directly applied to our task, so answering this question is a complicated endeavor. What we present below is a selection of two example-based models, one by Logan (1988) and one by Nosofsky (1984). Our goal here is not to create a competition between our model and these example-based models. Since each could be extended or modified to fit the data better and better, any search for the “best” model would be fruitless. Rather, we wanted to get a sense of how easily two “vanilla” example-based models could capture the above results. If they could not do so with ease, it would help clarify the value added by our model.

Logan’s (1988) theory of automatization claims that learning occurs via acquisition of a knowledge base of previously processed instances. The theory claims that an algorithmic process and a memory-retrieval process are in competition or, more specifically, are in a race to propose an answer to the current problem. Initially, only the algorithm can be used to solve problems, but if the current problem has been processed before, there will be an instance of it stored in memory, so the memory-retrieval process has a chance to win. Note that in Logan’s model, retrieval is based on identity between the current problem and a past example. Under the assumptions of Logan’s theory, when the same problem has been stored multiple times, the minimum time for one of its instances to be retrieved decreases as a power function of the number of times the problem has been seen. This leads to a gradual increase in the probability that the memory retrieval process wins the race against the algorithm.

We can map Logan’s theory onto operator selection in the BST by defining the algorithmic process as selecting the closest operator according to our distance metric and the retrieval process as selecting, from a stored instance, the operator that solved the problem in the past. When a solver sees a problem, retrieval of previous instances of that problem races against the algorithm for the selection of an operator. As soon as one process wins (i.e., a single instance is retrieved or the algorithm completes), the appropriate operator will be applied at a certain latency. Taking Logan’s theory literally, the speedup in time should decrease as a power function of the number of occurrences of the *same* problem and the operator selected should depend on the proportion of past instances of the *same* problem that were solved by each operator.

The learning that subjects exhibited in our study, however, is certainly not limited to previous experience on identical problems. We found substantial decreases in subjects’ latency to select their first operator before they ever received a repeated problem. For example, in Experiment 1, subjects did not see an identical problem until the final (test) phase, and yet their latencies

decreased as a power function of operator uses across the first two phases. Similar across-problem learning also occurred with respect to operator selections. For example, in Experiments 2 and 3, subjects changed their operator-selection tendencies on (test) problems they had never solved before. In fact, their test-phase selections changed for all five problem types even though they only had solved problems from two or three of the types. Logan's (1988) model, when taken literally, cannot account for these findings of across-problem learning.

However, as Logan suggests, strict identity may not be required for the use of past instances. To get an idea of how example-based learning would perform without the strict identity requirement, we look to an example-based model which allows for generalization across similar instances. Such a model is Nosofsky's (1984) categorization model.

We can map Nosofsky's model of categorization onto the current task by viewing the selection of an operator on a BST problem as a categorization task (e.g., "Should I select undershoot vs overshoot on this problem?" becomes "Is this an 'undershoot problem' or an 'overshoot problem?'"). According to this mapping, past problems (called examples) are stored in memory along with their correct categorization (i.e., the operator that led to a solution). Then, the probability that a particular operator will be selected for a new problem is calculated as the proportion of all past examples on which that operator was used, with each example weighted according to its psychological similarity to the new problem. Thus, the probability of selecting undershoot (u) on test problem t , given a set of previous examples x , is:

$$P_u(t) = \frac{\sum_{x \in U} \eta(x, t)}{\sum_{x \in U} \eta(x, t) + \sum_{x \in O} \eta(x, t)},$$

where the function $\eta(x, t)$ defines psychological similarity between example x and problem t , U is the set of past examples solved by undershoot, and O is the set of past examples solved by overshoot. This kind of processing should lead to an interaction between past history of success and distance to the goal: past successes of an operator on one problem type (e.g., high overshoot bias) will exert more influence on the current problem if the two are similar than if they are dissimilar.

We ran simulations of such a model to test whether this interaction would indeed be produced and to compare the model's output with the observed data. The comparison is not meant as a full-fledged test of Nosofsky's model but rather an exploration into how a typical example-based model, specified enough to be fit to data, might account for our results. To provide the model with some prior information about our distance metric, we included 30 preexamples as its initial database. These preexamples served as example-based training of our hill-climbing metric and were created from the solve phase problems of Experiments 2 and 3. For all of the preexamples, the categorization stored as "correct" was always the operator toward which the problem

was biased. The features of each preexample were also stored. We defined the features x_f of problem x as the numerical lengths of its sticks (i.e., the three building sticks and the desired stick). We defined the similarity function $\eta(n, t)$ between problems x and t as:

$$\eta(x, t) = e^{-d(x, t)},$$

where

$$d(x, t) = \sqrt{\sum_f w_f (x_f - t_f)^2}$$

and w_f is the weighting associated with feature f . The weights w_f specify the distance function d .

The bottom panels of Figures 7 and 10 present the best-fitting predictions of this example-based model for Experiments 2 and 3, respectively, when the weights w_f were estimated from the combined data sets. With four features, there were four such weights or free parameters in the example-based model; their best-fitting values are 1.28, 0.48, 2.12, and 2.68 for the three building sticks (left to right) and the desired stick, respectively. Regressing these predicted percentages on the observed percentages, we obtained the following equations: predicted = 10.0 + 0.68*observed, $R^2 = .85$, for Experiment 2 and predicted = -1.00 + 0.82*observed, $R^2 = .73$, for Experiment 3. In both cases, the R^2 values are lower than those obtained with our model's predictions, and yet our model required only two parameters; our model seems to be providing a better fit with fewer free parameters.

By inspection, one can see that the example-based model fits produced the expected interaction between test problem type and test number such that greater shifts occurred for the test problems more similar to the problems on which failures occurred (i.e., the "high toward" type for the biased and extreme-biased conditions and the "neutral" type for the neutral condition). In contrast, our rule-based model produced more equal sized shifts from test phase to test phase for the different test problem types (with the exception of floor effects). This distinction between the models seems to explain the better fit of our model to the data and suggests that, as in our model, solvers' history of information about successes and failures may be stored and used at the production level, not at the problem level.

One other difference is that Nosofsky's (1984) categorization model was not designed to make predictions on the time course of categorization decisions, whereas our model made specific latency predictions that were confirmed in the data. A more recent model by Nosofsky and Palmeri (submitted) presents an exemplar-based random walk model for categorization that predicts response times in speeded, multidimensional perceptual classification. This model combines elements from Nosofsky's (1984, 1986) previous models and Logan's (1988) model. If this model were generalized to apply in the

domain of problem solving, it is conceivable that it could predict the pattern of latencies observed in our data. Nevertheless, it is unclear if it would simultaneously be able to overcome the difficulties (exhibited by its sister model above) in fitting the selection data.

Although they make different predictions in the particular cases explored above, example-based and production-based models can be viewed as lying on the same continuum. Both types of models are maintaining some record of a history of experience that can be used to adapt subsequent performance. It is just a matter of the range of problem situations over which those records are aggregated. In example-based models, experience obtained on one problem may be combined with other experiences on the same problem (e.g., Logan's theory) or with experiences on problems that have similar features. In essence, it is the sensitivity of the similarity metric that determines what range of past experiences will have an impact. In production-based models, there is a particular range over which past success information is combined as well. By definition, this range is limited to the set of problems for which the production is applicable. In our model, the history-of-success of a production is an estimate of the probability of success of that production, given that the production's conditions are matched. If a production's conditions were very specific, its history-of-success would be a conditional probability that is very similar to referring to past examples based on a strict similarity metric. If, however, a production's conditions were very general, its history-of-success would be averaged over and applied to a great variety of problems. This would correspond to an example-based model with a very loose or insensitive similarity metric. The similarity metric of the example-based model described above had difficulty reflecting subjects' problem-general sensitivity to history-of-success while maintaining a context-specific sensitivity to distance-to-goal information in the current problem. In contrast, our model was sensitive to both sources of information due to its problem-general (but production-specific) history-of-success calculations and its problem-specific distance-to-goal calculations.

Conclusion

The experiments presented in this paper replicate the common finding that solvers are very sensitive to how close a particular move will take the current state to the goal state, even when they have had almost no experience in the task. Moreover, solvers tend to combine this distance-to-goal information with their knowledge of the success of different operators in an independent fashion. Our results support the notion that when solvers experience a success (or failure), they are learning something more general than whether it was a good (or bad) idea to try that particular operator on that particular problem. Rather, they seem to be attributing each success to the overall successfulness of the operator, which in turn can influence their likelihood of using that operator on all sorts of new problems. This finding speaks to the generality

of what people can learn by solving problems. We did find, however, that this kind of processing is not universal. Under particular situations, solvers can also exhibit problem-specific learning. Certainly, it would be implausible to claim that problem-specific learning never occurs. The most likely scenario is that both types of learning are at work in some balance. An issue for further research is under what conditions does each type of learning dominate and by what mechanisms?

APPENDIX A: SIMPLE MATHEMATICAL VERSION OF THE BST MODEL'S OPERATOR SELECTION

The analysis below is a special case of the rational analysis provided in Anderson (1990, 203–214). It is a simplification in that only predicted probability of success (and not expected cost nor effort expended so far) is used to evaluate potential moves. An advantage of this case is that we can derive the predicted probability of success of each move in closed form. (For the related computation of the model's probability of selecting one move over another, see Appendix B.)

The analysis below specifies a statistical model for the probability that a particular production, applied to the current problem state, will lead to success. The model is developed using a Bayesian framework in which the solver is gathering data from the observation of a sequence of such situations. First, we specify a parameter $\theta = (a, b, r)$ that describes the population from which the problem states are drawn. This multidimensional parameter characterizes the statistical regularities in the sequence of problem states that the solver observes. We assume that the problem states are independent and identically distributed given the parameter θ .

Next, we describe the distribution of various observable quantities of the problem states, given θ . The two quantities we are interested in are *distance-to-goal* (a continuous variable, x , that measures the distance between the goal state and the successor to the current state after the production is applied) and *Success* (a binary variable, S , that describes whether applying the production to the current state leads to success). The likelihood function (A.1) describes the joint distribution of S and x for a given θ :

$$P(S, x|\theta) = P(x|S, \theta)P(S|\theta). \quad (\text{A.1})$$

The first component on the right-hand side specifies the distribution of distance-to-goal values, conditional on whether the problem is solvable or not. We use an exponential distribution of distances with mean a for solvable problems ($S = 1$) and another exponential distribution of distances with mean b for unsolvable problems ($S = 0$):

$$P(x|S = 1, \theta) = ae^{-ax} \quad \text{and} \quad P(x|S = 0, \theta) = be^{-bx}.$$

Note that a and b are two of the components of the parameter θ identified

above, and “solvable” here means “solvable by the particular production being considered.” We define the second component of the right-hand side of (A.1) in terms of r , the overall probability of the production leading to success:

$$P(S = 1 | \theta) = r \quad \text{and} \quad P(S = 0 | \theta) = 1 - r.$$

Note that r is the third component of the parameter θ .

Next, we specify a prior distribution for θ ; this represents the solver’s uncertainty about the population of problem states. The solver will use the observed data to update this distribution to a posterior distribution for θ . For simplicity, we choose the priors on a and b to be point masses at the same values used in Anderson (1990): $a_0 = 1$ and $b_0 = 5$; more generality is possible. The fact that $a_0 < b_0$ represents the fact that distance-to-goal values tend to be shorter in success states than in failure states. We take r to have a beta distribution with parameters α and β . Note that α and β are not components of θ ; rather, they are model-fitting parameters that describe the solver’s prior.

The above specification implies that the solver will learn from experience, where a particular history of experience is $H_n = ((x_1, s_1), (x_2, s_2), \dots, (x_n, s_n))$ and each pair (x_i, s_i) describes the distance-to-goal and success information of a previously observed problem state. Learning occurs in two ways. First, the distribution of the parameter θ is updated. The posterior for θ given a history H_n is computed according to Bayes Theorem, which states that the posterior is proportional to the likelihood function multiplied by the prior. In our case, the posterior for θ involves a change in r . Since r was taken to be from a beta distribution, its posterior is in the same family and depends only on the number of successes and failures experienced:

$$P(\theta | H_n) \sim \langle a_0, b_0, \text{beta}(\alpha + \sum s_i, \beta + \sum(1 - s_i)) \rangle.$$

Second, the posterior predictive distribution for new values of S and x given a history H_n can be derived from the posterior distribution for θ . This is computed as

$$P(S_{n+1}, x_{n+1} | H_n) = \int P(S_{n+1}, x_{n+1} | \theta) P(\theta | H_n) d\theta,$$

which in our case is

$$P(S_{n+1} = 1, x_{n+1} | H_n) = a_0 e^{-a_0 x_{n+1}} \frac{\alpha + \sum s_i}{\alpha + \beta + n}$$

$$P(S_{n+1} = 0, x_{n+1} | H_n) = b_0 e^{-b_0 x_{n+1}} \frac{\beta + \sum (1 - s_i)}{\alpha + \beta + n}.$$

Finally, since x_{n+1} , the distance-to-goal for the new problem state, is observed, we can condition on it and thus use the posterior predictive distribution

to compute the posterior predicted probability of success (PPS) for the production under consideration. In particular, the posterior log odds of success is:

$$\begin{aligned} \log\left(\frac{P(S_{n+1} = 1 | x_{n+1}, H_n)}{P(S_{n+1} = 0 | x_{n+1}, H_n)}\right) \\ = \log\left(\frac{a_0}{b_0}\right) - (a_0 - b_0)x_{n+1} + \log\left(\frac{\alpha + \sum s_i}{\beta + \sum (1 - s_i)}\right). \end{aligned}$$

This is referred to as the Central Equation in the text. Note that the last term is the posterior log odds of success of the production under consideration and the preceding term involves the distance-to-goal, x_{n+1} , after applying that production to the current problem.

APPENDIX B: DESCRIPTION OF THE MODEL-FITTING PROCEDURE

Given an initial parameter setting, the idea is to calculate noisy PPS values (z_i) based on the (prenoise) PPS values (v_i) for each move i . For simplicity, we consider two possible moves for each problem: the overshoot move (denoted by subscript o) and the closer of the two undershoot moves (denoted by subscript u). This simplification makes the computation below much easier and is consistent with both the data (subjects almost never select the smallest building stick for their first move) and the model's behavior (the closer undershoot move always dominates the other undershoot move by enough that the model would almost never select the smallest building stick).

The probability of selecting each move equals the probability that that move's noisy PPS is larger than its competitor's. This probability can be computed, given the estimated PPS values for each move (v_i) and the noise variance σ^2 . The derivation below presents the model's predicted probability of selecting undershoot, p_u , with Φ as the standard normal cumulative distribution function:

$$\begin{aligned} p_u &= P(z_u > z_o), & z_i &\sim N(v_i, \sigma^2) \\ &= P(z_u - z_o > 0) \\ &= P(v_u - v_o - \varepsilon > 0), & \varepsilon &\sim N(0, 2\sigma^2) \\ &= P\left(\frac{\varepsilon}{\sigma\sqrt{2}} < \frac{v_u - v_o}{\sigma\sqrt{2}}\right) \\ &= \Phi\left(\frac{v_u - v_o}{\sigma\sqrt{2}}\right). \end{aligned}$$

The probability of selecting overshoot, then, is $p_o = 1 - p_u$.

For each problem being fit, p_u from above is compared with the observed

percentage of subjects selecting undershoot. Using Powell's method (Press, 1992) to search the parameter space, we minimize the sum of the squared deviations between the predicted and observed probabilities and present the best-fitting parameters.

REFERENCES

- Anderson, J. R. (1982). Acquisition of cognitive skill. *Psychological Review*, **89**(4), 369–406.
- Anderson, J. R. (1983). *The Architecture of Cognition*. Cambridge, Massachusetts: Harvard University Press.
- Anderson, J. R. (1990). *The Adaptive Character of Thought*. Hillsdale, New Jersey: Lawrence Erlbaum Associates.
- Anderson, J. R. (1993). *Rules of the Mind*. Hillsdale, New Jersey: Lawrence Erlbaum Associates.
- Anzai, Y., & Simon, H. A. (1979). The theory of learning by doing. *Psychological Review*, **86**, 124–140.
- Atwood, M. E., Masson, M. E. J., & Polson, P. G. (1980). Further explorations with a process model for water jug problems. *Memory & Cognition*, **8**(2), 182–192.
- Atwood, M. E., & Polson, P. G. (1976). A process model for water jug problems. *Cognitive Psychology*, **8**, 191–216.
- Bareiss, E. R., Porter, B. W., & Wier, C. C. (1988). Protos: An exemplar-based learning apprentice. *International Journal of Man-Machine Studies*, **29**, 549–561.
- Berger, J. O. (1985). *Statistical Decision Theory and Bayesian Analyses*. New York: Springer-Verlag.
- Bhaskar, R., & Simon, H. A. (1977). Problem solving in semantically rich domains: An example from engineering thermodynamics. *Cognitive Science*, **1**(2), 193–215.
- Cooper, L. A., & Regan, D. (1980). Attention, perception, and intelligence. In R. Sternberg (Ed.), *Handbook of human intelligence*. New York: Cambridge University Press.
- Delaney, P. F., Reder, L. M., Ritter, F. R. & Staszewski, J. J. (1994). *The power law of practice applies by strategy within task*. Manuscript submitted for publication.
- Elliott, S. W. & Anderson, J. R. (1995). Effect of memory decay on predictions from changing categories. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, **21**, 815–836.
- Hammond, K. (1986). *Case-based Planning: An Integrated Theory of Planning, Learning, and Memory*. Ph.D., Yale University.
- Jeffries, R., Polson, P. G., Razran, L., & Atwood, M. E. (1977). A process model for missionaries-cannibals and other river-crossing problems. *Cognitive Psychology*, **9**, 412–440.
- Klahr, D. (1985). Solving problems with ambiguous subgoal ordering: Preschoolers' performance. *Child Development*, **56**(4), 940–952.
- Larkin, J. (1981). Enriching formal knowledge: A model for learning to solve textbook physics problems. In J. R. Anderson (Ed.), *Cognitive Skills and their Acquisition*. Hillsdale, New Jersey: Erlbaum Associates.
- Logan, G. D. (1988). Toward an instance theory of automatization. *Psychological Review*, **95**(4), 492–527.
- Luchins, A. S. (1942). Mechanization in problem solving. *Psychological Monographs*, **54**(Whole No. 248).
- Luchins, A. S., & Luchins, E. H. (1959). *Rigidity of Behavior*. Eugene, Oregon: University of Oregon Books.
- Massaro, D. W., & Cohen, M. M. (1991). Integration versus interactive activation: The joint influence of stimulus and context in perception. *Cognitive Psychology*, **23**, 558–614.
- Massaro, D. W., & Friedman, D. (1990). Models of integration given multiple sources of information. *Psychological Review*, **97**(2), 225–252.
- Massaro, D. W., & Hary, J. M. (1986). Addressing issues in letter recognition. *Psychological Research*, **48**, 123–132.

- Medin, D. L., & Schaffer, M. M. (1978). Context theory of classification learning. *Psychological Review*, **85**, 207–238.
- Newell, A. (1990). *Unified Theories of Cognition*. Cambridge, Massachusetts: Harvard University Press.
- Newell, A., & Simon, H. A. (1972). *Human Problem Solving*. Englewood Cliffs, New Jersey: Prentice-Hall, Inc.
- Nosofsky, R. M. (1984). Choice, similarity, and the context theory of classification. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, **10**, 104–114.
- Nosofsky, R. M. (1986). Attention, similarity, and the identification-categorization relationship. *Journal of Experimental Psychology: General*, **115**(1), 39–57.
- Nosofsky, R. M., & Palmeri, T. J. *An Exemplar-Based Random Walk Model of Speeded Classification*. Manuscript submitted for publication.
- Nosofsky, R. M., Palmeri, T. J., & McKinley, S. C. (1994). Rule-plus-exception model of classification learning. *Psychological Review*, **101**(1), 53–79.
- Reder, L. M. (1987). Strategy selection in question answering. *Cognitive Psychology*, **19**, 90–138.
- Reder, L. M. (1988). Strategic control of retrieval strategies. In *The Psychology of Learning and Motivation* (pp. 227–259). Academic Press.
- Rickard, T. C. (1994). *Bending the power law: The transition from algorithm-based to memory-based performance*. Unpublished doctoral dissertation, University of Colorado, Boulder.
- Siegler, R. S., & Shipley, C. (1994). Variation, selection, and cognitive change. In G. Halford & T. Simon (Eds.), *Developing cognitive competence: New approaches to process modeling*. Hillsdale, New Jersey: Lawrence Erlbaum Associates.
- Thorndike, E. L. (1932). *The fundamentals of learning*. New York: Columbia University Press.
- (Accepted August 18, 1995)