

# Modeling changes in strategy selections over time

Christian D. Schunn Lynne M. Reder

Department of Psychology  
Carnegie Mellon University  
Pittsburgh, Pennsylvania 15213  
{schunn,redner}@cmu.edu

## Abstract

In this paper we present a method for fitting strategy choice data at the individual subject/individual trial level, and demonstrate, using the SAC model, that good fits to data can be obtained at this level. We conclude by discussing the sources of power in using this method.

## Introduction

How do people decide which strategy to use in solving a problem or answering a question? It is generally accepted that superficial features of a problem influence the strategy that is chosen to solve the problem. On the other hand, the claim that people select a question-answering strategy prior to executing any strategy, specifically prior to searching for the answer (Reder, 1987) has not been generally accepted (e.g., LeFevre, Greenham, & Waheed, 1993; Siegler & Jenkins, 1988). Even more controversial has been the claim that one of the criteria used in the decision of whether to try to retrieve the answer (or compute the answer by some other means) is a rapid "feeling-of-knowing" and that this rapid feeling-of-knowing depends *only* on features of the question and not at all on partial retrievals of the answer (Reder, 1987; Reder & Ritter, 1992; Schunn, Reder, Nhoyvannisvong, Richards, & Stroffolino, in press).

In this paper we focus on one set of arithmetic gameshow experiments that illustrate these results. In that series of experiments, this phenomenon was explored, carefully controlling for prior knowledge and tracking how learning and retrieve/compute strategy selections changed as a function of exposure to problems (Reder & Ritter, 1992; Schunn et al, in press). Arithmetic problems that people were unlikely to know before the experiment (e.g.,  $37 * 23$ ) were used as stimuli. Subjects were exposed to problems over and over (up to 28 times in one experiment) and each time made a very rapid assessment of whether they would be able to quickly retrieve the answer. The subjects indicated their strategy selection by pressing one key if they thought they could retrieve the answer, and other key if they thought they had to calculate the answer. This quick assessment took place in about a half second and was too little time to actually retrieve the answer. Subjects were allowed to study the answer to the problem after each trial and were given incentives to learn the answers and to select retrieve. However, there were disincentives for selecting retrieve if the answer was not known.

Evidence that these rapid strategy selections were based on aspects of the question, and not a partial retrieval of the answer, came from several results of the experiments. First, time to make the retrieve/compute decision was affected by practice with the task, but not practice with a specific problem. In other words, the practice with a specific problem that led to faster answers did not lead to faster preliminary decisions. The second result came from operator switch problems—some of the problems were distorted such that the operator was switched so that the two operands had appeared together before but not with that operator. For example, if  $21+35$  was presented earlier, then  $21*35$  could appear as a special operator switch trial. For such a problem, the subject could not know the answer since it had not been previously presented. But, such a problem would still look familiar as a first impression. In fact, it was operand co-occurrence that predicted retrieve judgments, not how often the exact problem had been seen.

A third result that supports the view that a rapid feeling-of-knowing comes from exposure to the problem and not knowledge of the answer comes from Experiment 1 of Schunn et al. In this experiment sometimes subjects were exposed to problems without getting a chance to actually answer the problem (either by calculation or retrieval) and without having the opportunity to study the answer. This manipulation was done for only a subset of the problems, called infrequently-answered problems. As one would expect, speed and accuracy in producing the answers were affected by how often subjects studied the questions; however, tendency to select retrieve was only affected by exposure to the problem itself.

In the rest of this paper, we describe a model of this retrieve/compute strategy selection process. The SAC model was first proposed, implemented and fit to these data by Reder with the help of Stroffolino and Richards. The original simulation and model fits have been extended and will appear in detail in Reder and Schunn, (in press) and Schunn, Reder, Nhoyvannisvong, Richards, and Stroffolino (in press). We present the model and the fits to the data from Schunn et al. to demonstrate how a model can be applied to subject data at the individual subject/individual trial level and provide an excellent fit to the data. At the end of the paper we will discuss which features of the model and simulations were important sources of power in the model fits.

## The SAC Model

The model is called SAC, which stands for *Source of Activation Confusion*. The representation used by the SAC model consists of interassociated nodes representing concepts that vary in long term strength. Here, we apply the SAC model to several arithmetic experiments. For these simulations, nodes represent numbers, operators, and whole problems. The nodes representing whole problems connect the operands and operators to the answers.

Each node has a base-level or long-term strength. The strength of a node represents the history of exposure to that concept, with more exposure producing greater strengthening. Nodes that represent arithmetic problems such as  $27 * 34$  would start out weak at the beginning of the experiment, as these problems were initially unfamiliar to the subjects. By contrast, nodes for familiar problems would be strong even at the beginning of the experiment. However, the experiments did not use problems that were likely to have pre-experimental familiarity, and the simulations presented here assume that all problem nodes are created for the first time during the experiment.

Strength can also be thought of as the base-line or resting level of activation of a node. Increases and decreases in this base-line strength change according to a power function:

$$B = c \sum t_j^{-d} \quad (1)$$

where  $B$  is the base level activation,  $c$  and  $d$  are constants, and  $t_j$  is the time since the  $j^{\text{th}}$  presentation. This function captures both power law decay of memories with time, and power law learning of memories with practice.

In addition to the base or resting level of activation of a node, there is also the *current activation* level of a node. The current level of a node will be higher than its base-line whenever it receives stimulation from the environment (i.e., when the concept is mentioned or perceived, or when the concept receives activation from other nodes). While base-line strength decays according to a power-function (i.e., first quickly and then slowly), *current* activation decays rapidly and exponentially towards the base level. Let  $A$  represent the current level of activation and  $B$  represent the base level of activation. Then, the decrease in *current* activation will be:

$$\Delta A = -\rho (A - B) \quad (2)$$

such that, after each trial, the current activation will decrease for every node by the proportion  $\rho$  times that node's current distance from its base level activation.

Activation spreads between nodes via links. Links connect nodes that are associated through conceptual relations. For example, links connect nodes that represent the components of a problem—operands and operators—to the node that represents the entire problem. Links also connect the nodes representing the entire problems to the nodes representing the answers. These links will vary in strength depending on how often the two concepts have been thought of concurrently. Strength of links also

depends on the delay between exposures. Specifically, link strength is determined by a power function given by:

$$S_{s,r} = \sum t_i^{-d_L} \quad (3)$$

where  $S_{s,r}$  is the strength of the link from the node  $s$  to node  $r$ ,  $t_i$  is the time since the  $i^{\text{th}}$  co-exposure, and  $d_L$  is the decay constant for links.

The current activation level of a node can rise from environmental stimulation or from associated nodes that send activation to it. The amount of activation that is sent depends on the activation level of the source (sending) node and on the strength of the link from the source node to the receiving node, relative to the strength of all other links emanating from the same source node. The change in activation of some node  $r$  is computed by summing the spread of activation from all source nodes  $s$  directly connected to node  $r$  according to the equation:

$$\Delta A_r = \sum (A_s * S_{s,r} / \sum S_{s,i}) \quad (4)$$

where  $\Delta A_r$  is the change in activation of the receiving node  $r$ ,  $A_s$  is the activation of each source node  $s$ ,  $S_{s,r}$  is strength of the link between nodes  $s$  and  $r$ , and  $\sum S_{s,i}$  is sum of the strengths of all links emanating from node  $s$ . The effect of the ratio  $S_{s,r} / \sum S_{s,i}$  is to limit the total spread from a node  $s$  to all connected nodes to be equal to the node  $s$ 's current activation  $A_s$ .

In this spreading activation model, feeling-of-knowing judgments are based on the activation level of the node representing the problem. In essence, we assume that feeling-of-knowing monitors intersection of activation from two source nodes. Specifically, when two terms in a problem send out activation to associated concepts and an intersection of activation is detected by bringing an intermediate node over threshold, a person will have a feeling-of-knowing response.

In our current simulations, we assume that when a problem is presented, all the nodes representing the components are activated. For example, in the problem  $23 * 14$ , the nodes representing  $23$ ,  $*$ , and  $14$  are all activated. Then, activation spreads from the component nodes to all the connected problem nodes. Problem nodes connected to several of the components receive the greatest amount of activation (e.g.,  $23 * 14$ ). The extent of activation that accumulates at the problem node affects the likelihood of selecting retrieve as the strategy of choice. In a similar fashion, activation spreads from problems nodes to answer nodes. This is how answers are retrieved.

Because activation that spreads to a node is added to the base activation, the selection of which problem node will have the highest final activation will also depend on the relative base level activations. The current activation level of the most (currently) active problem node is used to determine feeling-of-knowing. Based on the feeling-of-knowing, a decision is then made to retrieve or compute. That is, if the problem node has a relatively high activation level, then retrieval will most likely be selected; and if the problem node has a relatively low activation level, then computation will most likely be selected.

## Model Details

In addition to predicting feeling-of-knowing decisions (i.e., decisions between retrieval and computation), this model can also predict which answers are retrieved from memory, and the speed with which the answers are retrieved. Here, however, we focus on the feeling-of-knowing, or retrieve/compute, decisions. The computer simulation is given as input the same problems presented to each subject. Since each subject received a different set of problems in random order, a separate simulation was conducted for each subject. This precise yoking of the simulation to subjects was important because on a given trial the expected activation level for a problem would vary depending on the exact sequence of trials: for any subject on a given trial, the number of links, the current activation, and strengths would be different from any other subject's values. The simulation output is a probability of selecting to retrieve on each trial. We will now step through the process by which that probability is determined.

At the start of the experiment, the representation of memory for the simulation is identical regardless of the experimental stimuli to be seen. Nodes for the operands are assumed to already exist, whereas nodes for the problem components are assumed not to exist (i.e., the problems are novel). For simplicity, the initial base level strengths of the operand and operator nodes are set to a constant amount. When problems are seen for the first time, a problem node is created, as are the links from the component operand and operator nodes to the novel problem node. The initial base level strengths of the problem nodes and of the links is simply determined by the equations determining power-law growth and decay.

On each trial, all the nodes representing the problem components are activated to the same constant amount. Activation then spreads along the links emanating from nodes representing each of the problem components to nodes representing the complete problems. Activation only spreads to directly connected nodes at this point.

Once the activation has spread across these links, activation of the problem nodes can be used to make a strategy selection between retrieve and calculate. The activation value of the most active node is used. The simulation predicts a probability of choosing retrieve based on this activation value. This probability of choosing retrieve is calculated by assuming a normal distribution of activation values with a fixed variance and activation threshold for selecting retrieve. This probability is reflected in the formula:

$$P = N[(A-T)/\sigma] \quad (5)$$

where  $A$  is the activation of the most-active problem node,  $T$  is the subject's threshold,  $\sigma$  is the standard deviation, and  $N[x]$  is the area under the normal curve to the left of  $x$  for a normal curve with mean=0, and standard deviation=1.

A single value for the standard deviation parameter was used for all simulations. However, we assume that subjects vary in their thresholds for choosing between retrieve and compute. That is, some subjects are conservative and have

high thresholds, whereas other subjects are optimistic and have lower thresholds. A value between 30 and 200 was selected for each subject to maximize fit to their data. This value reflects the subject's overall base-rate of selecting retrieve. This wide range of possible values mirrored the large between-subjects variance that was found within each of the experiments in the retrieval selection rates. While the subjects might have differed on other dimensions as well, there were no other obvious differences (with the exception of the one mentioned below), and so, for parsimony's sake, the other six parameters were held constant across subjects.

After each trial, all the node strengths and activations are updated using Equations 1 and 2. Link strengths are updated for each link, following the same kind of function used to determine changes to base level activation—all the links connecting the problem component nodes to the problem node in the just-presented problem are strengthened, whereas all other links in the network are weakened. It is at this point that if a new problem has been presented for the first time, that a new node representing that problem is created, and links are created connecting the component nodes to the problem node.

The simulation just described involves seven parameters, listed in Table 1 along with the values that we used.

Table 1. SAC Model Parameters and Values

parameter name	value
input activation	50
$\rho$	0.8
$c$	5
$d$	0.175
$d_L$	0.12
$T$	30-200
$\sigma$	45
never-retrieve	T/F

There is one final component of the SAC model that required an additional parameter. This eighth parameter was only used for simulating some of the subjects. The parameter was simply a binary value by subject reflecting whether the subject had a predilection not to choose retrieve for a particular operator. This parameter was added because we found that some subjects had a strong aversion to choosing retrieve for a particular operator. To model these subjects, the probability of selecting retrieve on that operator is set to zero. For those subjects, the probability of selecting retrieve for the other (non-meta-rule) operator was simply determined as for the regular subjects—by the equations given in the SAC model. A simple 5% cutoff is used to select which subjects to model with this never-retrieve rule: subjects have to select retrieve for less than five percent of the trials with a particular operator.

In sum, there are eight parameters for the simulations, six of which were held constant for all simulations.

## Model-fitting Methodology

To compare the SAC model's predictions to subjects' actual retrieve/compute decisions, we used an aggregation

procedure developed by Anderson (1990). For each trial, for each subject, the model produced a probability of choosing retrieve based on the calculated activation values resulting from the trial history for that subject. That is, the probability reflected the model's experience with the exact same problems given to the subject. Since subjects made binary decisions and the simulation produced probabilities, it was necessary to aggregate trials. That is, all trials for a given subject in which the simulation predicted that the probability of selecting retrieve would fall between 0% and 10% were grouped together; all trials where the probability fell between 10% and 20% were grouped together and so on. We tabulated the proportion of retrieval strategy selections that were made by that subject for the exact same trials in each probability range. This was done for all probability ranges. Note that each subject contributes data points to each (or at least many) of the ranges. The fit of the model was tested by plotting mean actual proportion of retrieval strategy selections against mean expected percent retrievals. A perfect fit would be a straight line with a slope of 1 and a y-intercept of 0. We plot this desired line to show where the fitted points should actually lie.

Rather than plot the full scatter plot of each subject's value in each probability range, which often contains too many points from which to abstract the central tendency accurately, we plot the mean subject value (i.e., mean of subject means) within each range. To present an estimate of the subject variance, we also plot standard error bars. Furthermore, we present the  $r^2$  between predicted and actual values based on the full scatter plot, not the mean responses across subjects. This value presents an estimate of the amount of variance that the model accounts for at the individual subject level, a fine-grained level of detail not typically presented in tests of computational models. To assess whether there are any systematic biases in the model's predictions, we also present the slope and y-intercept of the best fitting regression line.

Since the number of subjects and data points per subject varied for the various experiments and analyses, it was necessary to vary the size (and hence number) of the probability ranges. Values were selected for each analysis using the following rule: the ranges were made sufficiently large such that subject contributed at least 5 data points to most of the ranges, thereby ensuring stable proportions.

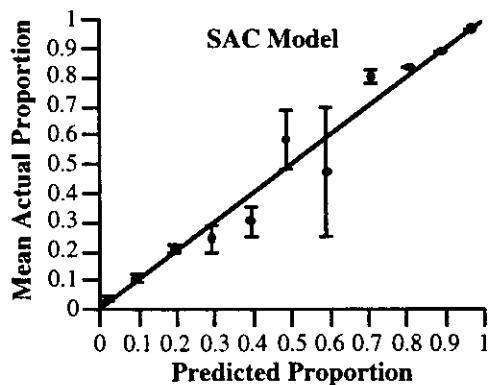


Figure 1. Model fit for all problems in Reder and Ritter.

## Simulation of Feeling of Knowing Data

As a first test of the SAC model, the strategy choice data from Reder and Ritter's Experiment 2 (described above) were compared with the SAC model's predictions. Using the aggregation procedure described earlier, subjects' actual retrieve/compute decisions were well predicted by SAC, with an  $r^2$  of .85 (see Figure 1). Note that the line drawn in the graph is the desired line actual=predicted. The slope of the best fitting line was not significantly different from 1 (slope=0.993,  $t(56)=0.125$ ,  $p>.9$ ), nor was the intercept significantly different from 0 (intercept=-0.001,  $t(56)=0.029$ ,  $p>.9$ ). In other words, the SAC model accounted for a large percent of the variance of the subject's strategy selections even at the individual subject level, and there were no systematic biases in the predictions.

A key result of Reder and Ritter was that subjects were as likely to select retrieve for operator-switch problems as for the training problems. The SAC model predicts this effect: Operators are associated with a large number of problems and the activation spread from a node along each link is inversely proportional to the total connection strength of the links emanating from that node. The fit of the SAC model to the operator-switch retrieve data is quite good ( $r^2=.82$ ). Figure 2 presents this fit. Fewer groupings were used in this analysis because there were relatively few operator-switch problems. Again, the slope of the best fitting line was not significantly different from 1 (slope=1.17,  $t(23)=1.42$ ,  $p>.15$ ), nor was the intercept significantly different from 0 (intercept= -0.009,  $t(23)=0.22$ ,  $p>.8$ ).

## Value of Each Parameter

One criticism of our model is that it contains many free parameters. This leads to the question: are all the parameters necessary? Rather than testing the value of all the parameters individually, we address this issue more globally by exploring one particular reduced alternative model. This alternative model might be called the everything-is-in-the-threshold-values account. Since each

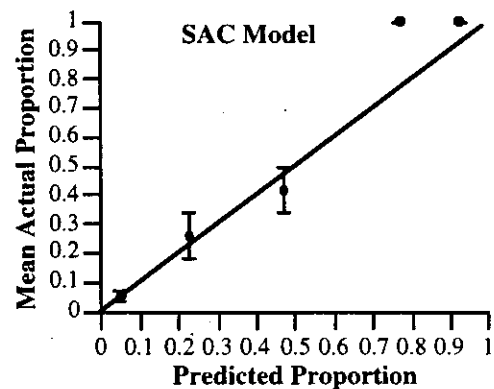


Figure 2. Model fit for the operator-switch problems only.

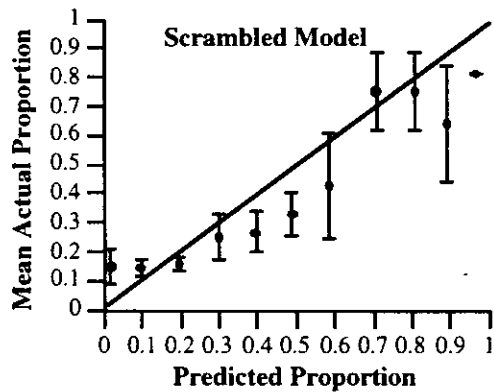


Figure 3. The Scrambled model fit for Reder and Ritter.

subject was given a different threshold value, and there are more subjects than probability ranges, one might argue that the good fits are due to having more free parameters than data points. To evaluate this alternative, a variant of the SAC model was created in which the model's predictions for each subject were scrambled. That is, the original model's predictions for each subject were kept, but the pairing with the subject's actual responses was randomly reorganizing. For example, rather than having the model's prediction for the first trial paired with the subject's response to the first trial, the model's prediction for the first trial might be paired with the subject's response to the tenth trial, or perhaps the 100th trial.

This scrambled model was able to account for 54% of the variance, suggesting that subject thresholds were an important part of the SAC model's good fit. However, this fit is much worse than the 85% of the variance for which the original model can account (see Figure 3). Furthermore, the scrambled model's best fitting regression line deviates significantly from the desired line: its slope differed significantly from 1 (slope=0.694,  $t(56)=3.56$ ,  $p<.001$ ), and the intercept differed significantly from 0 (intercept=0.068,  $t(56)=2.04$ ,  $p<.05$ ).

As another alternative to the SAC model, there is a class of strategy selection models which we call base rate models (e.g., Anderson, 1993; Siegler & Jenkins, 1988). Base rate models assume that strategies are selected according to the relative proportion of times each strategy has been successful. Such a model would correctly predict that subjects should initially select to calculate and gradually shift to selecting to retrieve because the experiment was designed such that subjects would initially know none of the answers and gradually know an increasingly larger percentage of the answers.

To evaluate such a base rate account could we tested the following model. We assumed that there was a linear increase over trials in the probability of selecting retrieval since analyses of the data had suggested that there were no significant curvilinear trends over time. Each subject was assigned two parameters: the initial retrieval rate, and the rate at which retrieval selections increased over time. Despite having many free parameters, the base rate model was only able to account for 71% of the variance in the

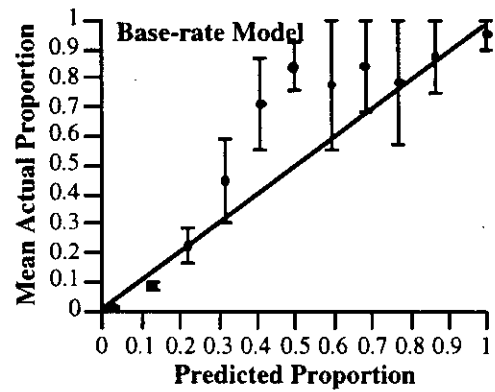


Figure 4. The Base-rate model fit for Reder and Ritter.

individual subject strategy selections (see Figure 4), significantly lower than the 85% produced by the SAC model. The slope of the best fitting line was not significantly different from 1 (slope=1.09,  $t(39)=0.75$ ,  $p>.45$ ), nor was the intercept significantly different from 0 (intercept= 0.032,  $t(39)=0.61$ ,  $p>.5$ ). Yet, Figure 4 shows that there were serious deviations between the predicted and actual strategy selections rates. Furthermore, the base rate model could not explain why subjects would be sensitive to the familiarity of operator-switch problems—it would predict that the current base rate would be used no matter what the familiarity of the operator-switch problem.

To test the model further, we applied the SAC model to the data from Schunn et al's Experiment 1 (described above). In order to provide a much stronger test of the SAC model, the model's parameters were set to the same values that were used in the simulation of Reder and Ritter. The only parameters that we did not take from the simulation of Reder and Ritter were the two subject-specific parameters: the subject's threshold, and whether they used the never-retrieve rule for an operator.

As with the Reder and Ritter data, the SAC model fit the new Experiment 1 data quite well, producing a  $r^2$  of .69 (see Figure 5). The slope of the best fitting line was not significantly different from 1 (slope=0.951,  $t(254)=1.22$ ,  $p>.2$ ), nor was the intercept significantly different from 0 (intercept=0.011,  $t(254)=0.54$ ,  $p>.5$ ). Thus, the SAC model generalizes very well to other data sets with all but the

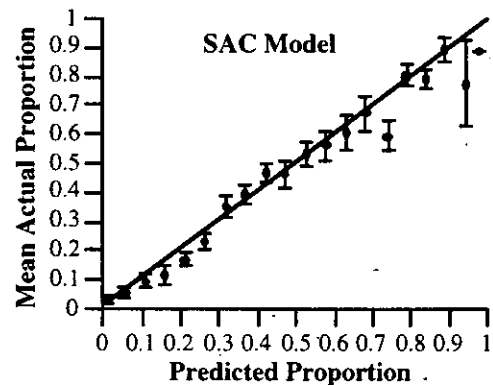


Figure 5. Model fit for all problems in Schunn et al. Exp. 1.

subject specific parameters held constant.

In sum, with the parameters set to the same values used for the simulations of the Reder and Ritter data (but still allowing subject-specific parameters), the SAC model produced a very good fit to the data. The SAC model, again holding all but the subject-specific parameters constant, has also been applied to the second experiment of Schunn et al. Unfortunately, for lack of space, we can only say that the model produced just as good a fit to that data set as was obtained for the model fits presented here.

## Discussion

In this paper, we described the application of a model to strategy selection data at the individual subject/individual trial level, and did so across several data sets holding the parameter values constant. Why were we able to provide such a good fit at this detailed level? In other words, what were the sources of power in the model? There are several factors that we believe were important. First, we employed a general class of model (spreading activation) that has been successfully employed in accounting for many aspects of human memory and behavior. By using a kind of model that has already demonstrated a lot of generality, we were more likely to find generality of our own model in applying it to different data sets using the same parameter values.

Second, the experiments that we modeled carefully controlled subjects' exposure to the problems (by using novel problems). In this way, we were able to capture individual subject performance without having to postulate individual differences in pre-existing knowledge, which normally would be very difficult to characterize with a small number of parameters. The remaining subject differences were general biases in strategy selection, which were fairly easy to capture with two simple parameters.

Third, we yoked our model to each individual subject's history of exposure to the various problems, rather than giving the model as input the average subject's experience. In this way, we were able to precisely capture the influence of different experiences both across subjects and across time. Not only is it likely that this feature contributed to the power of our model, but it also ruled out a particular excuse should we not have obtained good fits to the data—using this method, we could not hide behind vague claims about strange mixtures or orderings of problems or as a source of individual deviations from model predictions.

Finally, the experiments that we modeled provided a reasonably large number of data points per subject. The advantage of this is that we could posit two free parameters per subject without leading to problems of over-fitting the data, which in turn would have led to lack of generality across data sets. In general, whenever there are free parameters associated with each subject, the number of data points per subject rather than number of subjects seems more crucial to being able to accurately test the model.

In sum there are several features of the experiments that we selected to model and the way in which we conducted the model fitting process that are likely to have contributed to the strong fits that we obtained. There is one final feature of our model fitting process that deserves comment: the methods of presenting and evaluating the model fits. These methods of evaluating the fit of the models are not the most commonly used ones. For example, typically one computes some form of a goodness-of-fit test, which is typically a sum squared error between predicted and actual observations. We believe that this method has two serious drawbacks. First, it confounds relative ordering deviations and absolute magnitude deviations. That is, when there is a poor goodness-of-fit, it is unclear whether it is caused by a poor absolute magnitude of fit, or whether the relative ordering of predictions are in error. Our method of looking at both  $r^2$  (a measure of relative order relations) and the slope and intercept of the best fitting regression line (a measure of absolute magnitude fit) does distinguish these two aspects. Second, there is an unfortunate property of statistical goodness-of-fit tests: the model appears to do better the fewer the data points. In other words, one is punished for applying the model to large data sets. Third, the absolute magnitude of sum squared error is not very meaningful to the reader, whereas  $r^2$  is. This is why we believe that our evaluation methodology is better than the typically-used chi-square procedure.

## References

- Anderson, J. R. (1990). The adaptive character of thought. Hillsdale, NJ: Erlbaum.
- Anderson, J. R. (1993). Rules of mind. Hillsdale, NJ: Erlbaum.
- LeFevre, J., Greenham, S. L., & Waheed, N. (1993). The development of procedural and conceptual knowledge in computational estimation. Cognition and Instruction, *11*, 95-132.
- Reder, L. M. (1987). Strategy selection in question answering. Cognitive Psychology, *19*, 90-137.
- Reder, L., & Ritter, F., (1992). What determines initial feeling of knowing? Familiarity with questions terms, not with the answer. Journal of Experimental Psychology: Learning, Memory, & Cognition, *18*, 435-451.
- Reder, L. M & Schunn, C. D. (In press). Metacognition does not imply awareness: Strategy choice is governed by implicit learning and memory. To appear in: L. M. Reder, (Ed.), Implicit Memory and Metacognition. Hillsdale, N.J.: L. Erlbaum.
- Schunn, C., Reder, L., M., Nhouyvanisvong, A., Richards, D., & Stroffolino, P. (In press). To calculate or not calculate: A source activation confusion (SAC) model of problem-familiarity's role in strategy selection. Journal of Experimental Psychology: Learning, Memory & Cognition.

Siegler, R. S., & Jenkins, E. (1988). How children discover new strategies. Hillsdale, NJ: Erlbaum.