**Title:  Modeling the distinct phases of skill acquisition**

Caitlin Tenison[1]*, John R. Anderson[1]

[1] Department of Psychology, Carnegie Mellon University, Pittsburgh, PA. 15213

**Correspondence to:**

Caitlin Tenison

Department of Psychology ~ Baker Hall 342c

Pittsburgh, PA 15213

Email: ctenison@andrew.cmu.edu

**Abstract:**

A focus of early mathematics education is to build fluency through practice. Several models of skill acquisition have sought to explain the increase in fluency due to practice by modeling both the learning mechanisms driving this speedup and the changes in cognitive processes involved in executing the skill (such as transitioning from calculation to retrieval). In the current study, we use hidden Markov modeling to identify transitions in the learning process. This method accounts for the gradual speedup in problem solving and also uncovers abrupt changes in reaction time, which reflect changes in the cognitive processes that participants are using to solve math problems. We find that as participants practice solving math problems they transition through three distinct learning states. Each learning state shows some speedup with practice, but the major speedups are produced by transitions between learning states. In examining and comparing the behavioral and neurological profiles of each of these states, we find parallels with the three phases of skill acquisition proposed by Fitts and Posner (1967): a cognitive, an associative, and an autonomous phase.

**Introduction:**

Fluency, defined as the ability to quickly and accurately solve a problem, is a focus of early mathematics education (Kilpatrick, Swafford, & Findell, 2001). Studies have shown that students who score high on fluency measures maintain their skills over time (Singer-Dudek & Greer, 2005) and perform better on more complex math tasks than students with lower fluency scores (Skinner, Fletcher & Hennington, 1996). Furthermore, in a study of fluency-building interventions, Codding et al. (2007) found that low fluency students showed the most benefit from interventions that targeted accuracy, whereas students identified with intermediate levels of fluency responded best to interventions that emphasized speed. This differential response suggests that at different stages of practice, different skills are in play. To this end, we look at how different cognitive theories of fluency building view the decreases in latency due to practice. In defining the differences among these theories, an important distinction is made between the qualitative shifts in how students are solving problems and the gradual quantitative changes in proficiency of performance.

Within psychology, the speed at which a person performs a task can be used as a rough measure of the cognitive processes underlying the task. Almost universally, it has been observed that people speed up with practice; however, various theories have emerged that examine whether qualitative or quantitative changes underlie this speedup. In a classic reference paper, Newell and Rosenbloom (1981) observed that performance tends to speed up as a power function of the amount of practice. Dubbed, the 'Power Law of Practice', Newell and Rosenbloom

suggested that this speedup was due largely to a sequence of small qualitative changes in task execution. Subsequent attempts have been made to refine this characterization of practice. A major issue is whether the underlying learning function is best characterized as a single speedup function, suggesting continuous improvement in the same process, or if some of the speedup might be attributed to abrupt qualitative changes in cognitive processes (Anderson, 1982; Logan, 1988; Rickard, 1997). While not belaboring the question of the exact functional form of the quantitative change, the current study pursues the qualitative-quantitative debate within the context of a laboratory task that has clear analogs in the development of mathematical skill. In particular, we will argue that with practice participants move from calculation, to retrieval, to automatic recognition of the answer. In the present study, besides examining latency, we also consider the participants' self-reports and fMRI brain activation patterns. Using this multipronged approach, our aim is to gain insight into the cognitive changes associated with skill acquisition.

*Cognitive Models of Fluency Building*

Several existing cognitive models identify potential changes in cognitive processes that occur during practice and explain the mechanisms for how students become fluent. The Race model, which is part of the Instance Theory of Automization (Logan, 1988, 2002; Compton & Logan, 1991), proposes that each time a problem is practiced it becomes encoded in memory, then, when a participant sees the problem again, each of the previously encoded 'instances' independently race to generate the response and the fastest wins the race. The speedup results by

increasing the number of instances in memory, each of which has taken less time to execute than the previous instance. This theory predicts that the practice data would fit a single three-parameter power function regardless of what procedure a person is using to reach the answer (Logan, 1988).

$$\mu_{ret} = I + \beta n^{-\alpha}$$

Where $\mu_{ret}$ is the time it takes to retrieve the answer, I is the asymptotic latency (i.e., the fastest possible time), $\beta$ is the amount of latency that can be reduced with practice, $n$ is the number of practice opportunities, and $\alpha$ is the learning rate. This function describes the speedup as a single quantitative change. Later work suggested an improved fit using the exemplar based Random Walk model, a type of Instance model (Nosovsky & Palmeri, 1997; Palmeri, 1999). Whether or not the three-parameter power function is the most appropriate for modeling speedup due to practice is unclear (Heathcote, Brown & Mewhort, 2000; Myung, Kim & Pitt, 2000); however, the Race model, and other theories of practice related speedup, are concerned not only with which function best describes the speedup (e.g. three-parameter power function, four- parameter power function, three-exponential function), but also with the qualitative changes to the cognitive processing that occurs with practice. These models are distinguished by their theories of the type of cognitive processes employed to produce a response and how practice influences the use of these processes.

In contrast to the Race model, The Component Power Law Theory (CMPL) proposed by Rickard (1997) suggests that a single memory association is formed and then strengthened with practice. On seeing a problem, a participant must

choose between two strategies, either applying the problem solving algorithm or retrieving the answer. This qualitative distinction between calculation and retrieval is marked by a 'step-function' decrease in problem solving latency associated with the reduction in cognitive steps needed to implement the retrieval strategy. Within each strategy, however, CMPL supports quantitative changes as well. The CMPL model suggests that fluency building best fits two functions, one governing calculation and the other governing retrieval. In his 2004 paper, Rickard illustrated how a normal distribution around a fixed mean best fits the calculation stage of problem solving, and a power function best fits the problem-solving speedup of retrieval. Delaney et al. (1998) had earlier argued that learning consisted of two power functions, one for calculation strategies and one for retrieval strategies, and that fitting a separate power function to each strategy better fits the data than a single power function, as suggested by the Race model. Additional papers supporting CMPL, have noted that speedup in the execution of computation procedures during this first learning stage fits a power function (Bajic & Rickard, 2009, 2012; Rickard 1997). This quantitative speedup in both calculation and retrieval is explained by the mechanism of association and the strategy choice bottleneck.

Both the CMPL model and Delaney et al. (1998) share similarities with the Strategy Choice and Discovery (SCADS) model (Shrager & Siegler, 1998), which is an extension of the earlier Adaptive Strategy Choice (ASCM) model (Siegler & Shipley, 1995). This model suggests that the shift to retrieval strategies arises out of an increased association between the math problem and the solution (Shrager &

Siegler, 1998). When students are first introduced to problems, it takes all their attention to solve that problem; however, with practice, attention resources are freed up for strategy discovery. As students discover and practice new strategies, attention resources are then apportioned towards detecting and deleting redundancies in problem-solving strategies, leading to determining optimal order in which to execute strategy steps and eventually retrieval of the answer from memory.

In addition to associative speedup processes, strategy choice and use contribute to problem solving latency as well. Reder and Ritter (1992) found that participants presented with a problem assessed whether or not they knew the problem prior to employing a strategy to solve it. Bajic and Rickard (2009, 2011) suggested that choosing strategies produced a "pause-effect"; a pause before initiating the algorithm that increased in length prior to the first retrieval. Furthermore, they found effects of a strategy choice bottleneck after people began to retrieve that were then reduced with practice. Bajic and Rickard (2009, 2011) observed a *partial-retrieval effect* as people performed some mix of applying the algorithm interrupted by retrieving the answer from memory. In these studies, the frequency in which the *partial-retrieval effect* occurred was dependent on the type of task. In addition to these effects of practice, there were also individual differences in strategy choice. Across several different age groups researchers have identified individuals who are able to retrieve the answer, but choose instead to calculate (Hecht, 2006; Siegler, 1988; Touron & Hertzog, 2004). This work suggests that

factors outside of practice may also influence what strategies participants choose to use.

Another view goes back to the early proposal by Fitts and Posner (1967) that a skill progresses through three qualitative stages: the cognitive, the associative and the autonomous stage. Anderson operationalized these three stages using ACT theory and identified several learning mechanisms to simulate skill acquisition (Anderson, 1982). This early work has evolved, and the learning mechanisms have been changed in more recent years with ACT-R (Anderson et al. 2004). ACT-R is a cognitive architecture grounded in psychological theory that can model learning and performance. Knowledge in ACT-R is represented as either declarative knowledge, which consists of facts (i.e., 2 is a number), or procedural knowledge, which is composed of productions representing mappings of states to actions (i.e. if the goal is to solve 2 +2, then the answer is 4). ACT-R models the three Fitts & Posner stages as follows:

1. **Cognitive:**  The participant must perform a sequence of calculations to produce the answer. For instance, if someone knows the addition table, but not the multiplication table, they can calculate the answer to 3 x 5 by repeated addition.

2. **Associative:**  The participant learns the declarative fact that 3 x 5 = 15 and now can retrieve that fact from declarative memory when faced with the question.

3. **Autonomous:** With enough practice the learner can acquire a production that simply produces the answer 15 in response to the problem without

querying declarative memory. While strategy choice occurs in the two prior

phases, during the final autonomous phase, time is not spent choosing a

strategy; rather, answering the question effectively become a reflex.

The ACT-R model allows for speed up within the first two stages. In the cognitive

stage performance speedups are due to knowledge compilation. *Knowledge*

*compilation* combines two learning mechanisms: *proceduralization,* which is when

production rules take on a more task specific form and *composition,* which is when

multi-step procedures are collapsed into single procedures (Anderson, 1982;

Taatgen & Lee, 2003). This mechanism both captures the quantitative speedup due

to strengthening the retrieval of newly learned instructions, as well as,

improvements in the computation (such as when someone learns a partial result). In

the associative stage, *declarative strengthening* increases the speed of retrievals and

impacts performance (Anderson, 1982; Anderson, 2007). With subsequent

retrievals, the stimulus becomes more strongly associated with the correct response

causing faster and more accurate retrievals. Once the autonomous phase is reached,

ACT-R does not have any further basis for speedup, although one could imagine the

potential increase by acquisition of new motor programs (perhaps, for coding the

entry of the answer).

*The Current Study*

In the current study, we use a novel math task in order to study the effect of

practice on how participants solve problems. At the beginning of the task, students

are taught a novel math operation that requires counting a sequence of numbers,

summing those numbers, and reporting a total value. Over the course of the experiment participants practice solving three problems 36 times each. While we focus our modeling efforts on the problem solving latency data, we also collect fMRI data and strategy reports (both concurrent and retrospective) as a means of better understanding our learning states and as a check that our statistical model is in fact detecting distinct cognitive states and not arbitrary cut-offs in student efficiency. After half of the trials, participants report their problem-solving strategy by choosing from a list of commonly used strategies (calculation, retrieval, partial calculation, and other). Our goals in analyzing these data are (1) to test how many states best represents the learning data and (2) to verify if the states we detect are cognitively distinct.

Attempts to evaluate the different models of skill acquisition are limited by difficulties in identifying where strategy shifts may lie. Previous work has identified where the strategy shift occurs by designing tasks where participants must take distinct actions for distinct strategies (e.g., Bajic & Rickard, 2007, 2010). We are interested in where strategy shift occurs as well; however, our approach does not presuppose the strategies, but rather, uses an unsupervised modeling method that evaluates different strategy possibilities according to how well they fit the data. In particular, we will be focusing on latency and evaluating when changes in latency might reflect a change in strategy. Rickard (2004) suggested a method for identifying a single shift point (between calculation and retrieval) using only latency data. While his study was focused on comparing one- state to two- state models, we broadened our study to compare any number of states. In the face of uncertainty

about where a shift point lies, our method does not simply pick the best for that problem, but weighs the various possibilities by how probable they are given data from the whole problem set.

A clear signature of a change in strategy would be a major drop in latency from one trial to the next, but given the noise in latency data one can never know for sure whether or not there is a shift. We need a way of quantifying this uncertainty and evaluating different models in the face of such uncertainty. In this study, we use hidden Markov models (HMMs) to allow us to efficiently calculate all the possible trajectories participants may take among states, as well as, the relative likelihood of each. An HMM treats the actual state of the system as hidden, but whose probability can be estimated. An HMM is based on the Markov assumption that future states of a process depend on the present state. This is critical for the dynamic programming solutions that accompany HMMs. Beyond being well suited theoretically for modeling skill acquisition, HMMs avoid the problem of fitting average data. Fitting learning to problem-solving latencies averaged across all subjects, rather than for each practiced item separately, has been shown to be problematic (Clauset, Shalizi & Newman, 2009; Haider & Frensch, 2002; Heathcote, Brown & Mewhort, 2000; Myung, Kim & Pitt, 2000). Using HMMs allows us to assess each problem for each participant individually, yet still combine the evidence from all the problems and participants when judging among competing models.

The two-part structure of the paper first addresses our work modeling skill acquisition, and second, explores the results of our modeling work. We will apply HMMs to latency data in our task to decide both the number of learning states and

whether or not there is within-state speedup. We will use both behavioral and brain data to assess whether or not these latency-identified learning states are qualitatively distinct. Specifically, we will assess if one state involves computation, another just retrieval, and a third state, a stimulus-response link deleting all cognitive involvement, including retrieval.

**Method**

*Participants*

Twenty-three university students (9 women; mean age 22/ SD 2.3) participated in the study. Two participants were excluded because of excessive movement that caused the fMRI signal to fluctuate. A third excluded participant did not understand the instructions. Participants gave informed written consent and received monetary compensation for their participation. All participants were right handed and had normal or corrected-to-normal vision. The university ethics board approved the study.

*Materials*

To investigate how students become fluent, we taught all participants a type of arithmetic operation they had never seen before and asked them to practice solving problems using this operation. This novel operation, called a 'Pyramid problem', uses the same algorithm as in the experiment by Anderson et al. (2011). Pyramid problems use '$' as an operator between two values, the base (B) and the height (H) (i.e., B$H, where B is the base and H the height). The base is the starting

number, and the height indicates the total number of terms to be added where each

term is one less than the previous. For example, 8$4 expands to 8+7+6+5. To further

explain, starting with the base number (8), participants identify the number of

terms (4) in a descending sequence (8,7,6, and 5). Then, by adding the terms

correctly, the participants find the value for the operation. In the example above, the

value is 26. To solve these problems, participants kept a running total in their head

as they added several integers. We created our problem sets using three heights, 3, 4

and 5. The maximum base was 11. The minimum base for problems of a given height

was the height plus 1. This produced a sample size of 21 unique problems from

which we created our practiced and novel problems sets (See Appendix D for

complete list of problems). Over the course of the arithmetic training task, we

selected three problems, one from each height, for each participant to practice (we

refer to these as 'practiced problems') and 18 novel problem repeated twice, but

only after each participant had solved more than 50 problems. This allowed us to

compare two problem sets to distinguish between the effects of practice and

becoming fluent with a specific problem. To minimize overlap between problems,

we selected problems where the base was not equal to, or within 1, of the other

practiced problems. For example, if we selected 5$3 (5+4+3), we would not select

5$4 (5+4+3+2), because practice solving one problem could easily transfer the

ability to solve another.

      We used two methods to assess strategy use. During the arithmetic training

task, participants completed a concurrent multiple-choice assessment after solving

each problem within the selected block, and a retrospective assessment at the end of

the arithmetic training task. The concurrent assessment presented a list of

strategies from which participants were encouraged to choose the strategy that best

matched the one used to solve the previous problem. We compiled the strategy

options by considering the frequency of reported use from a previous experiment

using the same problems (Tenison et al. 2014). At the start of the blocks during

which the concurrent assessment was collected, the instructions included

definitions of the available strategy choices. "Retrieve" was defined as remembering

the answer; "calculate" was defined as using arithmetic to find the answer; "partial"

was described as partially calculating and partially remembering the solution.

Participants were instructed to indicate "other" for any strategy that did not fit the

three categories. The concurrent assessment was presented on the screen after

participants had entered an answer to the Pyramid problem. Participants were

prompted, "How did you solve the problem?" and were given four choices: 1)

Retrieve 2) Calculate 3) Partial 4) Other. Only one participant reported 'other'

strategy. On investigating the participant's retrospective strategy report, it was

found that, even though a calculation strategy was used, the participant considered

it distinct from the calculate category. For the purposes of this paper we collapsed

the report of 'other' strategy used into the calculate category. A retrospective

strategy assessment was given at the completion of the arithmetic training task,

where each participant solved 15 paper-based problems that included all the

practiced and some of the novel problems and wrote detailed explanations for how

each problem was solved. We used these reports to establish what participants were

reporting as the 'partial' and 'calculate' strategies during the concurrent assessment.

Participants used a numeric keypad with a standard keyboard number arrangement to type the solution to the math problems and to indicate the problem solving strategies. The backspace key was disabled so that each participant was unable to delete numbers once imputed.

*Warm-up Task Procedure*

The experiment consisted of a warm-up task and an arithmetic training task. The aim of the warm-up task was to familiarize participants with what they would be doing in the arithmetic training task. Following the warm-up, the arithmetic training task consisted of 6 blocks during which participants practiced a set of Pyramid problems. During the arithmetic training task, participants were inside a fMRI scanner; during the warm-up task they were not.

During the warm-up task, the participants solved Pyramid problems, reported strategy use, and practiced using a numeric keypad. The Pyramid problems were taught by explaining the role of base and height in constructing a string of addition facts. The investigator instructed the participants to start with the base and then count down until the number of terms listed was equal to the height. The investigator next instructed the participants to add these values together. After each problem was solved, the participant received correctness feedback and instruction on how to solve the problem. This feedback was visually represented as the chain of addition productions required to calculate the answer. Following this feedback, participants reported the strategy employed to solve that problem. None of the problems used in the warm-up task were used in the arithmetic training task. In

addition, the participants practiced keying their responses on a numeric keypad.

Since participants would not be able to see the keypad during the arithmetic task, a

box cover was placed over the keypad during the warm up task so that participants

could practice touch only entry. Participants had eight opportunities to press the

corresponding key while viewing a screen presentation of random numbers from 0-

9. The warm-up session required approximately 15 minutes to complete.

*Arithmetic training procedure*

The arithmetic training task was divided into 6 blocks during which we

recorded fMRI data and behavioral responses. The 6 blocks comprised of 3 practiced

problems, with each one solved 36 times, 18 novel problems, with each one solved

twice, and 6 warm up problems (See Appendix D for set of problems used). Each of

the 6 blocks contained 25 items; a warm-up problem for which the response was

discarded, 18 practiced items (3 problems with each problem repeated 6 times) and

6 novel problems. Novel problems matched practiced problems in difficulty and

repeated once, but only after the participant had solved more than 50 problems and

would be unlikely to remember the problem. Participants also completed

concurrent strategy assessments during the 2nd, 4th, and 6th blocks (the

alternating blocks allowed the experimenter to check the reactivity of the

assessment).

Pyramid problems were presented on the screen following a 2-second

fixation period. Once the problem appeared on the screen, the participant was

allowed a maximum of 30 seconds to indicate knowledge of a solution by pressing

the return key on the numeric keypad. After pressing 'return', participants had 5

seconds to input a solution using the keypad and then press the return key. Problem

solving time was defined as the time between the appearance of the math problem

and the point at which the participant indicated that they were ready to input the

answer. Limiting the time available for the participants to input the solution forced

them to finish solving the problem prior to beginning input phase. After answering

the problem, the participant was given correctness feedback and information about

how the problem should have been solved. During the 2nd, 4th, or 6th blocks, a

screen appeared that prompted participants, "How did you solve the problem?".

Participants were given 5 seconds in which to select the number that corresponded

to the strategy used. At the end of each problem-solving trial, a 1-back task was

presented onscreen for a randomly selected time (between 6 to 12 seconds). This

was to prevent metacognitive reflection on the previous problem and allow the

hemodynamic response of the brain to return to baseline. In the 1-back task,

participants were asked to judge if the letter in the center of the screen was the

same as the previous letter seen.

*fMRI Data Acquisition and Analysis*

   Images were acquired using gradient echo-echo planar image acquisition on

a Siemens 3T Verio Scanner using a 32 channel RF head coil, with 2 s. repetition

time (TR), 30 ms. echo time, 79° flip angle, and 20 cm. field of view. The experiment

acquired 34 axial slices on each TR using a 3.2 mm thick, 64×64 matrix. This

produces voxels that are 3.2 mm high and 3.125 x 3.125 mm2. The anterior

commissure-posterior commissure line was on the 11th slice from the bottom scan slice. Acquired images were pre-processed and analyzed using AFNI (Cox, 1996). Functional images were motion-corrected using 6-parameter 3D registration. All images were then slice-time centered at 1 sec and co-registered to a common reference structural MRI by means of a 12-parameter 3D registration and smoothed with an 6 mm full-width-half-maximum 3D Gaussian filter to accommodate individual differences in anatomy.

Our ROI analysis focused on six predefined regions based on the predictions from the skill acquisition framework. The regions and the reasons for their selection are discussed in context of our behavioral results. These data were analyzed using a general linear model (GLM). Our design matrix consisted of a regressor for each problem (a total of 2880), one across all problems for response activity, and one across all problems for feedback. The design matrix was constructed by convolving the boxcar functions of these variables with a hemodynamic function (using the standard hemodynamic response – Friston et al, 1998). This normalizes the MRI signal so that the mean value of any region across an entire block was 100. Thus, deviations of 1 unit from 100 can be interpreted as a 1% change in activation.

**Results**

*Practice Effects on Fluency Building*

We first examine the effect of practice on speed and accuracy of problem solving to verify that the participant was able to become fluent in solving the practiced problems by the end of the experiment.
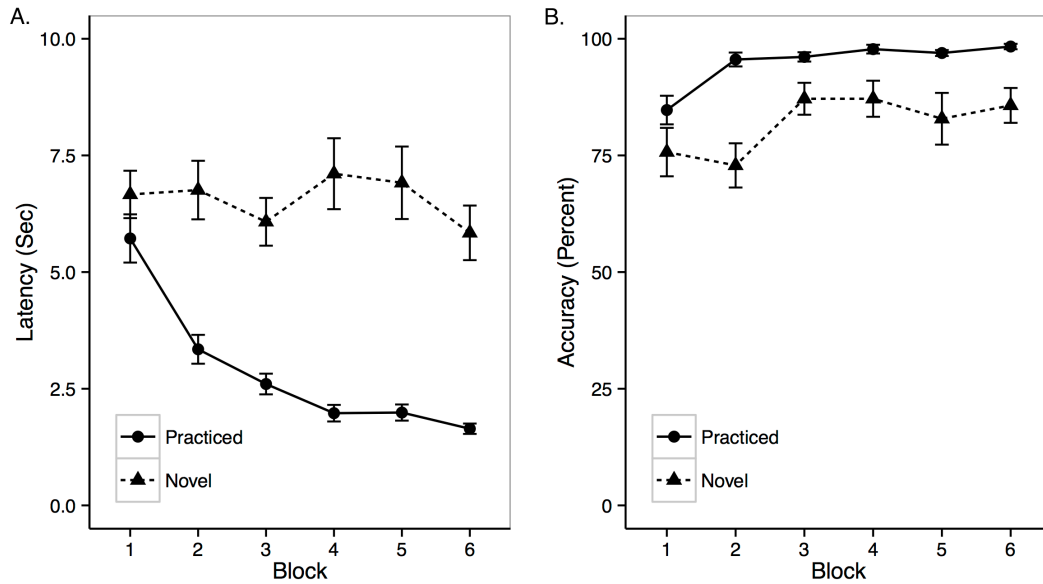
*Figure 1*. Mean latencies (a) and percent accuracy (b) on the six blocks of the experiment for both practiced and novel problems. Error bars represent standard errors.

A repeated measures ANOVA run on the latency data revealed a significant main effect of problem group (practiced vs. novel), $F(1,19)=69.26$, MSE= 296.8, $p<.001$, block, $F(5,95)=23.59$, MSE=9.4, $p<.001$, and a significant problem by block interaction, $F(5,95)=14.24$, MSE=4.5, $p<.001$. The time it took participants to solve practiced problems decreased, but the time to solve novel problems remained fairly constant (Figure 1.A). A repeated-measures ANOVA run on accuracy data revealed similar effects to the latency data. Participants' problem solving accuracy increased over time, but their accuracy solving novel problems remained relatively constant (Effect of practice: $F(1,19)=18.2$, MSE=1.0, $p<.001$, effect of block: $F(5,95)=8.7$, MSE=0.1, $p<.001$ , interaction: $F(5,95)=2.8$, MSE=0.3, $p<.05$) (Figure 1.B). We fit a regression to the problem solving latencies for novel problems for each participant.

We found that for each novel problem solved the problem solving latency was reduced an average of 0.03 seconds (SD= .01 seconds).

*Strategy Assessments*

We next examined whether the inclusion of a concurrent assessment altered the participants' problem solving latencies. Because concurrent assessments appeared on the 2nd, 4th, and 6th problem solving block for each participant, it is difficult to disassociate the effects of practice from the effects of the assessment; however, we can look for this effect amongst novel problems. To test this, we combined the first, middle, and last two scan blocks to create 3 time blocks. We then conducted a repeated-measures ANOVA on the novel problems investigating the relationship between time block and assessment. The analysis revealed a marginal main effect of time block, $F(1,19)=3.8$, MSE=2.0, $p=0.06$, but no significant effect of assessment, $F(1,19)=1.3$, MSE=1.3, $p=0.27$), nor a significant interaction between the variables, $F(1,59)=1.1$, MSE=0.6, $p=0.29$.

We collected concurrent reports during block 2, 4 and 6 (illustrated in Table 1). Table 1 reports the percentage of problems solved with a given strategy for practiced and novel problems separately. This shows an increase in reports of retrieval and a decrease in partial and calculate strategy for trained problems over the course of the experiment. On the other hand, novel problems show relatively little change in reports. Retrospective reports of partial and calculation strategies gave us some insight into the range of procedures participants were describing in their concurrent assessments. For the most part, participants reported a partial

Table 1
Mean percent strategy use reported within the concurrent assessments collected at three
time points during the task for practiced and novel problems

|  |  | Block 2 (%) | Block 4 (%) | Block 6 (%) |
|---|---|---|---|---|
| Practiced Problems | Calculate | 15 | 4 | 3 |
|  | Partial | 22 | 10 | 3 |
|  | Retrieval | 63 | 86 | 94 |
| Novel Problems | Calculate | 83 | 83 | 83 |
|  | Partial | 13 | 15 | 14 |
|  | Retrieval | 4 | 3 | 3 |

strategy as using previously solved problems to help solve new ones. If given 5$3 to

solve, a participant might report remembering the answer to 5$4 and then

subtracting 2 from that sum. In these cases, the problems recalled were either one of

the practiced problems or the previously solved problem. This partial strategy is

distinct from the strategy Bajic and Rickard (2009, 2011) identify as partial-

algorithm retrieval (interrupting calculation with retrieval) . We were unable to

identify cases where the participant interrupted calculation by retrieving the

answer; however, it is possible by the end of the experiment when the retrospective

reports for this strategy were collected, it was no longer used. In addition,

participants reported some minor variants of calculation strategies. Of these, the

majority of participants reported calculating the solution by starting with the largest

number and adding the numbers sequentially. Some participants used a calculation

strategy that involved adding all number groups that summed to ten prior to adding

other numbers, for example, in 7$5, one might add 7+3 and 6+4, prior to adding 5.

*Hidden Markov Model Fitting*

 Our goal was to determine how many distinct cognitive states participants went

through as they acquired fluency, and whether or not all the improvement in solving

problems was in the transition between states or during within-state speedup. By fitting the data using hidden Markov models (HMMs), we were able to address some of the challenges to modeling skill acquisition. First, to identify when a latency drop for an item reflects a state change versus some combination of general speedup plus noise in the data, we tested different HMMs that reflected different assumptions about the number of states and whether or not there was learning within a state. This allowed us to test several possibilities: from a situation where participants were just continuously speeding up within a single state (as in Logan's Race model), to one in which all speedup comes from the transitioning from states of constant rates of response. Furthermore, simulation analysis has indicated the importance of accounting for individual and item-level differences in the transition between states (Haider & Frensch, 2002). In this experiment, we expect that participants will make these transitions after different amounts of practice for different items, so to accommodate for such differences, we modeled the latency history of each practiced problem for each participant individually. To account for such variation in transitions, we fit the HMMs to every practiced problem. We refer to the specific practice opportunity of a problem as a 'trial'.

To be fairly exhaustive in examining the possible numbers of states, we fit HMMs that had one- to five- learning states and also varied whether there was within-state speedup or not. Historically, several different speedup functions have been used to model skill acquisition. Although the focus of the current paper is on identifying the phases of skill acquisition, we ran our analysis using several different speedup functions to determine whether or not function type significantly alters the number

of states our model fits. Besides a no-learning model (that had no within-state speedup), we tested three different functions to describe within-state speedup: a three-parameter power function, a three-parameter exponential function, and the four-parameter APEX function. This created a total of 20 HMMs (one to five states each with four options for within-state speedup). The HMMs fitting one- through three- learning states model hypotheses of the three theories of skill acquisition previously discussed: the Instance Theory; CMPL and ACT-R. An HMM with one-learning state represents the view that learning is best captured by a continuous speedup function (Logan, 1988, 2002). The HMMs with two- learning states fits the CMPL theory, and the HMMs with three-learning states models the three phases of skill acquisition (Anderson, 1982).

   We do not expect problem solving latencies to behave like perfect flat functions or smooth decreasing functions within a learning state, as there will be variability around the predicted problem solving times for that state, so we modeled the variability with a gamma distribution. If the predicted response time for an item in learning state $i$ practiced $j$ times was $T_{ij}$ then the shape parameter of the gamma distribution was set to 3[1] and a scale parameter to 1/3 $T_{ij}$. This predicts that as participants' response latencies decreases so does the variance of the response latencies. The error bars around the curves in 2b illustrates this variance.

---

[1] To minimize the number of parameters, we assumed a fixed shape of 3. Results do not seem to change for other values.
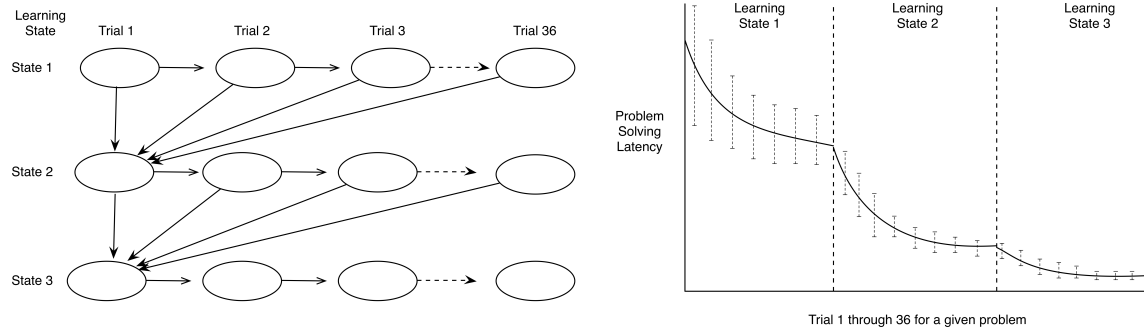
*Figure 2.* The right diagram (a) illustrates the structure of the three-State HMM we fit to the data. Each latent state of the HMM corresponds to an observation in a learning state. Students have some probability of transitioning between observations within a single learning state, or into the first node of the next learning state. The graph on the left (b), illustrates how the information from the HMM translates into latency predictions for a practiced problem. We estimated the variance around the function according to a gamma distribution with a scale factor that decreased with practice.

The models with within-state speedup would violate the Markov assumption (that behavior only depends on current state and not past history) if we had only one Markov state for each learning state. Therefore, we used a more complex Markov model that had a distinct state for each number of trials in a learning state. Figures 2a and 2b show the distinction between HMM trial-states and the learning states that we are considering (in this figure we show 3 learning states). In order to model within-state speedup, we need to have a different HMM state for each number of trials in a learning state. Thus, if there were $i$ learning states, there were $36i$ HMM states. Because the HMMs reflect the uncertainty in determining where the learning of a problem occurs, we can characterize the probability that a participant is in any of the $36i$ HMM states on any trial. Note that many trajectories can lead to the same HMM state. For instance, a participant may have reached the fourth trial in the third learning state after many different combinations of trials in the first and second state. We used the standard expectation maximization algorithm (Rabiner,

1989) to estimate the parameters associated with the model. The number of parameters we estimated was dependent on the function we used to describe behavior within a learning state. Some parameters were estimated for each learning state, whereas others were estimated once across all learning states. Appendix B provided a mathematical expression for the probability that the HMM is calculating and maximizing.

In fitting an $i$-state HMM, we estimated $i$-1 transition probabilities between the learning states. In the case of a no-within state speedup, the expected response latency within state $i$ was $T_i$ for all trials in that state.

$$T_{ij} = T_i$$

Therefore, for an $i$ state no-learning model we estimated $2i$-1 parameters in total for both the transition probabilities and the mean within-state latencies. In both non-learning and learning models, because we had a fixed shape parameter, no further parameters beyond mean latency were required to model the variability in latency.

In the case of within-state speedup, a trial's latency was not only a function of the state $i$, but also the number of the trials $j$ within that state, which we will denote as $T_{ij}$. We tested models of several different speedup functions. In the case of a power function the expected time $T_{ij}$ on trial $j$ in State $i$ was calculated as:

$$T_{ij} = \text{I} + \beta_i j^{-\alpha}$$

The parameter I is the intercept, the parameter $\beta_i$ was the amount of time that could speedup, and the exponent $\alpha$ controlled the rate of speedup. We assume intercept and exponent were constant across learning states and just estimated a

different $\beta_i$ or each state. In the case of an exponential function the expected time $T_{ij}$ was:

$$T_{ij} = I + \beta_i e^{-\delta j}$$

As with the power function, $\beta_i$ indicates the total amount that the latency can be reduced by. The learning rate is captured by the exponential function by the exponent $\delta$. Both the power and exponential functions required estimating $i+2$ latency parameters for $i$ states, which in addition to the transition parameters gave these models a total of $2i+1$ parameters, two more than the no-learning model. The APEX function, proposed by Heathcote and colleagues (2000) calculated the expected time $T_{ij}$ as a combination of the previously described exponential and power functions.

$$T_{ij} = I + \beta_i j^{-\alpha} e^{-\delta j}$$

As with the previous functions for HMMs using the APEX function, we estimated a single $\alpha$, $\delta$, $\beta_i$, and I for each model. For each state within a model, a separate $\beta_i$ was estimated. Thus, producing a total of $2i+2$ parameters for an $i$-State HMM.

We assessed the 20 models through a comparison of BIC values. The Bayesian Information Criterion (BIC) metric uses the likelihood of a model fitting the data, but also includes penalties for added parameters (Schwarz, 1978). This criterion is regarded as a parsimonious means of determining nonlinear model fit (Spiess & Neumeyer, 2010). The BIC is calculated as:

$$BIC = -2 * LL + K * \ln(N_{Observe})$$

Table 2
BIC scores for No-Learning and Learning models

|  |  | 1 State | 2 States | 3 States | 4 States | 5 States |
|---|---|---|---|---|---|---|
| No Within State Speed Up | Parameters | 1 | 3 | 5 | 7 | 9 |
|  | BIC | 8911.9 | 6981.0 | 6625.0 | 6601.2* | 6610.8 |
| 3 Parameter Power Function | Parameters | 3 | 5 | 7 | 9 | 11 |
|  | BIC | 7507.9 | 6713.1 | 6590.3 | 6590.0* | 6604.4 |
| Exponential | Parameters | 3 | 5 | 7 | 9 | 11 |
|  | BIC | 7516.4 | 6739.7 | 6590.6* | 6591.6 | 6606.5 |
| Apex | Parameters | 4 | 6 | 8 | 10 | 12 |
|  | BIC | 7510.5 | 6715.8 | 6595.4* | 6596.9 | 6612.0 |

*Note.* Asterisks indicates the best number of states within a type of model according to BIC

Where $LL$ is the log likelihood of the data, $K$ is the number of parameters, and $N_{Observe}$ is the number of observations (60 items x 36 trials = 2160 latencies). BIC scores for the different models are shown in Table 2. Looking at these BIC measures it would appear that within-state speedup always leads to better performance than no-within state speedup. While the four-learning states power function HMM has the best BIC scores, the improvement from three-to four-learning states is negligible (Kass & Raftery, 1995). Thus, we conclude that the three-learning state, three-parameter power function provides the best fit to our data. Also the difference between the power and exponential functions are negligible. The fits of the APEX model are substantially worse indicating that its additional parameter is not needed.

We ran a simulation analysis (described in Appendix A) to explore the sensitivity of our HMM modeling technique for discovering the true number of states underlying a dataset and identifying the shape of the learning within the states. We simulated data using a three-parameter power function, a three-parameter exponential function and no speed up For each function we generated

samples reflecting one-, two- and three- state models of skill acquisition. Without

exception, when the data was best fit by a three-state speedup function it was

generated by a speedup function. The ability, however, to discriminate between the

power function and exponential functions was poor. These simulations suggest that

our data represent three states with within-state speedup. We, however, cannot

discriminate between power and exponential forms of the speedup function.

*Examining the three- learning state model*

In the following examination we will used the 3-state power function. The

best fitting parameters for this model were

I -- Latency Intercept (shared by all states):   0 sec.

$\alpha$ -- Power Exponent (shared by all states):  .118

$\beta_1$ – Latency that speeds up in State 1: 9.09 sec.

$\beta_2$ – Latency that speeds up in State 2: 3.72 sec.

$\beta_3$ – Latency that speeds up in State 3: 1.86 sec.

$\pi_{12}$ -- Probability of transitioning from state 1 to 2 on a trial:  .174

$\pi_{23}$ -- Probability of transitioning from state 2 to 3 on a trial: .069

Note that the intercept I estimates to be zero, which in a sense means there is one

parameter less. Since the intercept is zero all the latency can eventually speed up,

but since the exponent $\alpha$ is quite small that speed up is slow. Most of the latency

speed up is produced by the sharp drop in the $\beta_i$ parameters between states.
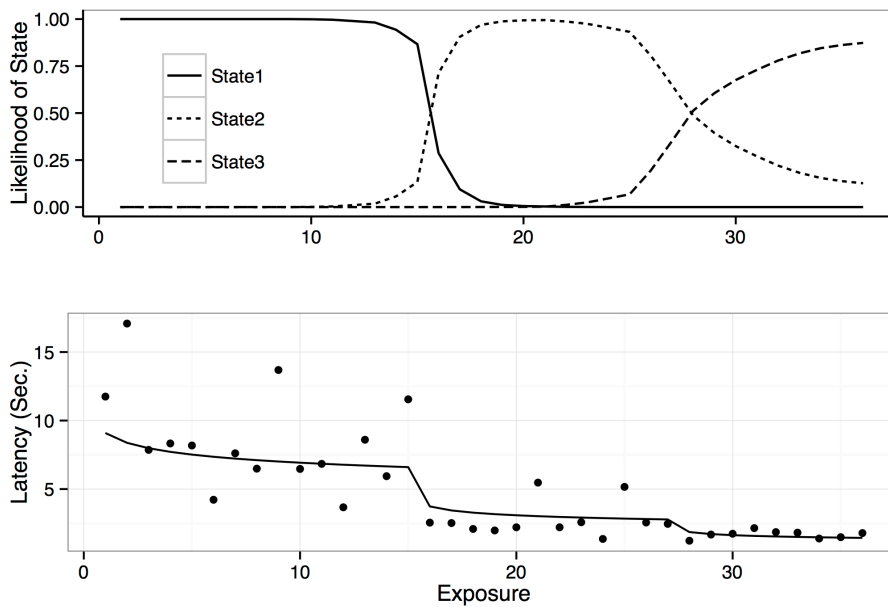
*Figure 3. (a)* The likelihood of each state for a single item practiced 36 times. (b) The latency scores for each practice opportunity plotted are plotted alongside the predicted latency assuming each trial is in the majority state.

The HMM we fit to the data estimates a probability that, on a given trial, a problem is in a certain state. Figure 3 shows one practiced problem with its estimated probabilities of being in one of three learning states, As seen in most, but not all trials, there is a fairly high probability of being in one of the three states; each of the 2160 trials has a probability of being in one of the three states greater than .5, and 83% have a probability greater than .8. The variability in the times is quite representative of the observed variability. Finally, Figure 3 shows the likelihood of state belonging alongside the predicted latency, assuming each trial is in its majority state. While Figure 3, illustrates data of a single individual practicing a single item, more individual's data are displayed in the lattice plot in Appendix C. This lattice

plot provides an additional comparison of the latencies predicted by the one-, two-

and three-state power function HMMs.

Table 3
Descriptive statistics indicated average properties of the three states

|  |  | Latency (sec) | Accuracy (%) | Calculate (%) | Partial (%) | Retrieval (%) | Total Count |
|---|---|---|---|---|---|---|---|
| Three State Model | State1 | 7.98 | 84.0 | 28 | 55 | 17 | 307 |
|  | State2 | 2.95 | 94.7 | 10 | 20 | 70 | 759 |
|  | State3 | 1.41 | 98.1 | 3 | 2 | 95 | 1094 |
| Novel Problems |  | 7.34 | 80.8 | 83 | 14 | 3 | 720 |

*Assessing the Three States*

To assess whether state transitions (such as those illustrated in Figure 3) are

qualitatively distinct, we turn to the properties of the different trials with the

assumption that they are in their most probable state. Table 3 shows some basic

statistics associated with the three states. Not surprising, the states show strong

differences in mean latency, given that the HMM used latency to form these learning

states. Additionally, they show substantial differences in accuracy and reported

strategy use. The concurrent reports show that State 1 is dominated by procedural

computation (either calculation or partial), whereas, problems labeled as State 2

and 3 are generally reported as solved by retrieving the answer.

We next investigated if there was an effect of difficulty for problems

categorized in the three-learning states. Several models of skill acquisition would

suggest that problem difficulty would impact latency when participants are

performing computations to produce the answer, but not when retrieving the

answer (Compton and Logan, 1991; Logan, 1988; Rickard, 1997). The problems that

the current study participants practiced differed in difficulty by the number of terms

added to find the solution. While height 3 problems only required two additions,

height 4 problems required three additions, and height 5 problems increased to four

additions. Because each type of problem for each participant was not necessarily

assigned to each state by the HMM, the number of observations we had for each

difficulty level varied considerably and a repeated measures ANOVA could not be

run. Instead, we used a linear mixed-effect to model the latency data. It is common

in such models to take the log of the latency scores as the dependent variable in

order to meet the normality assumptions of the model (Kliegl, Masson, and Richter,

2010). Our model included three fixed effects and an interaction: problem difficulty,

state assignment, number of times a problem had been seen (we will refer to this as

*exposure*), and the interaction between problem difficulty and state assignment. We

also included a random intercept for participants to account for individual

differences. We initially fit a maximal model that included a random slope and each

of our fixed effects (Barr, Levy, Scheepers, and Tily, 2013). Including problem

exposure as a random effect led to our model not converging. We also found that the

level of variability in state assignment was not sufficient to be included as a random

effect. Removing this random effect from the model decreased the BIC score (2300

to 2259), indicating improved model fit. We calculated t statistics and p-values for

our fixed effect using the Kenward-Roger approximation for degrees of freedom

(Kenward and Roger, 1997). We found a significant effect of both problem difficulty

($F_{2,24.3}$ =.8.87, p<=0.001), state assignment ($F_{2,927.1}$ =584.0, p<0.001) and problem

exposure ($F_{1,1650}$=79.3, p<0.001), as well as a significant interaction between
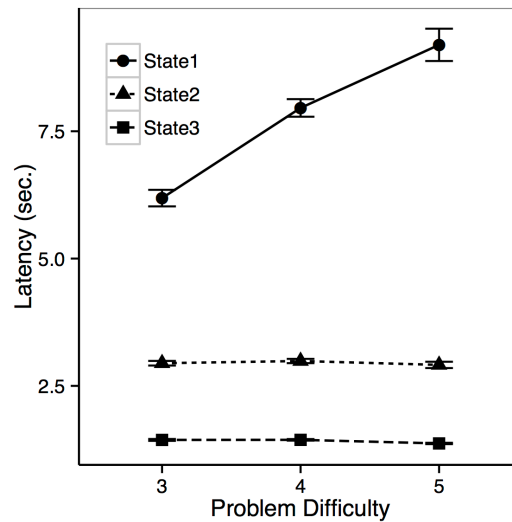
*Figure 4.* Mean latencies for problems of different 'Heights' (height provides an indicator of how many calculations the participant must do) in each of the three states. Error bars represent standard error.

difficulty and state assignment ($F_{4,419}$ =7.6, p<0.001). Figure 4 indicates the meaning

of our significant interaction, showing a problem difficulty effect for items in State 1,

but none apparent in State 2 or State 3.

Given the reports and the effects of problem difficulty, we can be fairly

confident that participants were calculating an answer according to the assigned

algorithm in State 1, but not doing so in the other two states. This indicates that a

state difference, identified by looking at latency, maps onto a qualitative difference

in cognitive processing. The difference between States 2 and 3, on the other hand,

seems less dramatic. Both states show a majority of reported retrievals and no effect

of problem size. There were differences in accuracy but this was smaller than the

difference with State 1. To see if we can find any stronger evidence for the

distinction between the second and third states, we now turn to our fMRI data.

*Using neural indicators to understand the three-state model*

We collected fMRI data while participants were trained on these problems so that we could later use differences in brain responses to better understand the results of our three-state model. We ran an ROI analysis to consider the distinction our model made between State 1 and State 2, as well as between State 2 and State 3. We chose to examine only regions related to the cognitive processes we hypothesized to be involved in the three states. First, the mean latencies, the reported strategy use, and the effect of problem difficulty suggest that State 1 involved a sequence of calculations while State 2 may have involved a simple retrieval. Second, State 2 may have featured an effortful retrieval that, with practice, transitioned into State 3, a more automatic process, as in the Fitts and Posner's *automization* phase of skill acquisition. We ran a region of interest (ROI) analysis to see if regions implicated by these hypothesis would differ significantly in percent activation between State 1 and State 2, and between State 2 and State 3. To test the hypothesis that State 1 contains more calculation than State 2, we tested the left horizontal intraparietal sulcus (HIPS), an area commonly associated with numerical processing (Dehaene, Piazza, Pinel, and Cohen, 2003). To test the theory that State 3 is a highly automatized retrieval as described by Fitts and Posner's automization phase of learning, we chose two ROIs: the left fusiform gyrus, which involves visual recognition of familiar stimuli (McCandliss, Cohen and Dehaene, 2003), and the left motor cortex, which is associated with responding. Finally, to test the hypothesis that State 2 involves a more effortful retrieval than State 3, we chose the left lateral

inferior prefrontal cortex (LIPFC), which is associated with effortful retrieval (Anderson and Fincham, 2014; Badre and Wagner, 2007).

To adjust for participants who did not reach all three learning states, and account for within subject variance, we ran a linear mixed-effects model (Chen et al, 2013; Judd, Westfall and Kenny, 2012) using the *lmer* function in R (Bates, Maechler & Bolker, 2012). We wanted to contrast two possible explanations as to why activations may be changing over the course of the experiment. The first hypothesis is that there is some continuous change with experience of an item. The second is that there is a discontinuous change associated with a change of state. Therefore, we included trial of exposure to represent continuous change and two variables to distinguish the three states. The first variable, S1-S2, was the difference between the probability that the item was in State 1 minus the probability that the item was in State 2. The second, S2-S3, was the difference between the probabilities of State 2 and State 3. For example, the S1-S2 variable would be near 1 if there were a high probability of being in State 1, and near -1 if there was a high probability of being in State 2. The S1-S2 variable would be near zero either when it was uncertain, or when there was a high probability of being in State 3.

For each ROI, we fit the data using a linear mixed-effects model to the percent activation in that region with fixed effects for S1-S2, for S2-S3, and for number of times that they had been exposed to that problem before. We also included a random intercept for each item for each participant to account for item-level differences. Finally, we included the average activation of the left and right auditory region as both a fixed and random effect. With scanner noise remaining

constant throughout the task, including this region captures general brain variance (such as effects of breathing) and regresses it out from the percent activation of a specific region (Birn, Smith, Jones and Bandettini, 2008).

For each ROI, we initially fit a maximal model, which included random effects for S1-S2, S2-S3, exposure, and average percent activation of the left and right auditory region. In addition, we included fixed effects for all four of these random effects. We then used the procedure outlined by Barr and colleagues (2013) for selecting the most appropriate model[2]. For each ROI, we report the F statistics and p-values for our fixed effect using the Kenward-Roger approximation for degrees of freedom from the best fitting model (Table 4). Positive coefficients for S1-S2 indicate greater brain activity associated with greater likelihood of State 1 and negative coefficients indicate greater brain activity associated with higher probability of State 2. In the S2-S3 contrast variable, positive coefficients indicate greater brain activation associated with higher probability of State 2, whereas negative coefficients are associated with greater probability of State 3. Figure 5 shows the coefficients for the S1-S2 and S2-S3 contrast for the three regions with significant coefficients. For our exposure variable, positive coefficients indicate increases in activation associated with increases in problem exposure (the number

---

[2] After fitting the maximal model, we removed any random effects that did not improve model fit (reducing BIC by at least 2), and re-ran the model. At this point we considered the fixed effects and removed any which did not show a significant contribution to the model (p>.05). If model fit worsened with the removal of non-significant variables then that meant that including the variable helped account for general variance in the model that was unrelated to the percent activation. In testing both random and fixed effects we used BIC as a means of comparing these different models and judging whether removing either of these effects improved model fit.

Table 4
Summary of Mixed-Effects Models which Best Fit the Percent Activation of Four Regions of Interest

| Region | Model | Intercept (SE) | Fixed Effect | Slope (SE) | NDF,DDF | F | P | Model BIC |
|---|---|---|---|---|---|---|---|---|
| Left HIPS | F: S1-S2 S2-S3 Exposure Auditory RS: Exposure Auditory RI: Subject | .25 (.04) | S1-S2 | .06 (.03) | 1,2102.3 | 5.86 | .02 | 3609.7 |
| | | | S2-S3 | .08 (.03) | 1,2042.4 | 13.8** | <.001 | |
| | | | Exposure | -.003 (.003) | 1,20.3 | .74 | .4 | |
| | | | Auditory | .05 (.06) | 1,18.5 | 65.5** | <.001 | |
| | | | RS: Exposure | .0001 (.01) | | | | |
| | | | RS: Auditory | .06 (.25) | | | | |
| Left Motor | F: S1-S2 S2-S3 Exposure Auditory RS: Auditory RI: Subject | .28(.05) | S1-S2 | -.19 (.03) | 1,573.3 | 36.1** | <.001 | 3832.5 |
| | | | S2-S3 | -.2 (.02) | 1,966.6 | 56.5** | <.001 | |
| | | | Exposure | .007 (.002) | 28.3 | 6.3 | .02 | |
| | | | Auditory | .77 (.08) | 18.7 | 90.3** | <.001 | |
| | | | RS: Auditory | .01 (.32) | | | | |
| Left Fusiform | F: Auditory RS: Auditory RI: Subject | -.02 (.04) | Auditory | .71(.1) | 1, 18.84 | 44.4** | <.001 | 4013.5 |
| | | | RS: Auditory | .20(.44) | | | | |
| Left PFC | F: S1-S2 S2-S3 Auditory Exposure RS: Auditory RI: Subject | .17 (.04) | S1-S2 | .11 (.04) | 1,2135.9 | 14.4** | <.001 | 3968.9 |
| | | | S2-S3 | .18 (.03) | 1,1954.8 | 92.0** | <.001 | |
| | | | Exposure | -.005 (.001) | 1,2143 | 7.8* | <.01 | |
| | | | Auditory | .58 (.07) | 1,18.6 | 18.6** | <.001 | |
| | | | RS: Auditory | .09 (.30) | | | | |

*Note.* F = Fixed Effects, RS = Random Slope, RI = Random Intercept. Numerator (NDF) and denominator degrees of freedom (DDF), F, and P value calculated using the Kenward-Rodgers approximation for degrees of freedom.

*p<.01. **p<.001

of times a problem is solved). Finally, random slopes within the model fits indicate

that the coefficients representing these relationships differ between participants.

We find that increased activity in the left HIPS is associated with increased

activity of being State 1 rather than State 2, and of being in State 2 rather than State

3 (Figure 5). Because inclusion of random slopes does not improve the model for

either S1-S2 or S2-S3, this suggests that the relationships between these factors and

percent activation are invariant across subjects. Furthermore, while we do not find a
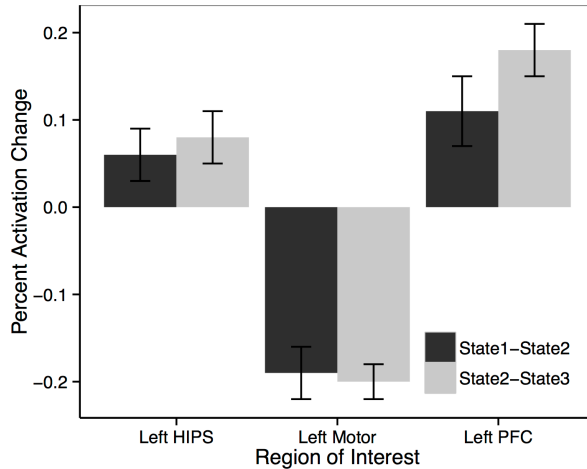
*Figure 5.* The difference in activation between State 1 to State2 and between state State 2 and State 3. These values are from our mixed-effect models and account for both subject level differences and the overall effect of practice.

significant effect of exposure on percent activation, including it as both a fixed effect and random slope accounted for variance within the sample and improved the model fit. These results indicate that as participants progress through the three states, HIPS activation reduces with the change in states rather than gradually with practice, supporting the hypothesis that the large changes in latency which cause our model to identify state changes reflects a change in cognitive processing.

We find that increased activity in the left inferior prefrontal cortex (LIPFC) is associated with increased likelihood of being in State 1 over State 2 and being in State 2 over State 3. Additionally, a significant negative relationship exists between activity in the LIPFC and exposure, indicating that as participants practiced problems less activation in the LIPFC is observed. These effects do not vary across individuals. The significance of both the change in state and the effect of practice

suggest that this reduction in the LIPFC reflects a process which contributes to both within and between state speedup.

We next consider left fusiform gyrus and the left motor cortex. We selected these regions to test the hypothesis that the distinction between State 2 and 3 reflects a shift from stimulus-retrieval-response process to a stimulus-response process. We find no effect in the left fusiform gyrus in either the S1-S2 or S2-S3 contrasts or exposure. Neither State assignment, nor exposure, accounts for significant variance in the left fusiform gyrus, suggesting that the visual recognition processes are not impacted by the removal of the retrieval process. We did find the increased likelihood of being in State 2 rather 1, and State 3 rather than 2 associated with increased activity in the left motor region. Additionally, we find increased activation in this region associated with increases in the number of times the participant solved a problem (exposure). This effect varies between individuals, whereas the S1-S2 and S2-S3 effects do not. These findings suggests that as the response latency decreases, the initiation of the motor response to report the solution becomes the dominant process occurring during the trial.

In conclusion, the brain activity indicates that state-change marks a shift in the amount of cognitive activity (LIPFC and HIPS) relative to motor activity. Both LIPFC and HIPS are more active in the earlier states, and the left motor cortex is more active for latter states. The results from this analysis provide some insight into the roles each of these regions play in accounting for both the gradual speedup due to practice, as well as, the cognitive processes underlying the changes between states.

**General Discussion**

*Model Building*

The use of HMMs allows us to avoid some of the problems of other approaches. First, it allows us to fit individual items rather than the much criticized (.e.g. Clauset, Shalizi & Newman, 2009; Haider & Frensch, 2002; Heathcote, Brown & Mewhort, 2000; Myung, Kim & Pitt, 2000) practice of fitting average data. Fitting individual items is relatively straightforward if a single speedup function is being fit to all the trials for a practiced problem, but gets more challenging when one models where the item transitions among states. Variability in latency leaves considerable uncertainty about when these transitions take place. Unless one has a highly reliable external indicator of strategy[3], one should deal with that uncertainty rather than take a best guess. HMMs allow one to calculate the probability of the data over all the possible transitions trajectories an item may have taken.

In the current study, we fit each practiced item for each participant to an HMM which captured learning-related speedup. More complex models with more parameters will always have a higher likelihood of fitting the data. To correct for model complexity, we used BIC to evaluate the fit of the models. Using BIC we found that regardless of how many states our models fit, the within-state learning models better fit the data than the no-speedup models. This suggests that even accounting for abrupt decreases in problem solving latency, that with practice there is some speedup of execution of cognitive tasks. Among the different learning functions tested, the three-parameter power function model fit the data better than other

---

[3] We do not think our reports have that kind of reliability although a method like that of Bajic & Rickard may have.

learning functions although its advantage over an exponential function was

marginal. Across all the learning functions tested, we found that while including

more than three states sometimes produced better BIC scores, these improvements

were not significant. We ran a simulation analysis to explore the reliability of our

method when determining the type and number of functions underlying our dataset.

This analysis indicated that it is unlikely that we fit more learning states than are

actually present in the dataset. Furthermore, the simulation analysis showed that

while our HMM method can discriminate between data generated with and without

in-state learning, it is unable to discriminate the type of learning function generating

the data in our experiment.

Discriminating between power and exponential functions is difficult and it

may be that our experiment did not have enough trials to discern the shape of the

function. Therefore, we looked at the data from Rickard's (1997) Experiment 2

(available at http://www.newcl.org/node/6) that involved many more trials per

item. We chose this dataset because it involved alphabet arithmetic, which is similar

to our task in that the problems are to be solved by a repetitive procedure. There

are 440 items with at least 72 observations and we chose to fit the models to these.

The models chosen were no within-state speedup, three-parameter power function,

and three parameter exponential function. We investigated 1 to 3 States. The results

are presented in Table 5. Consistent with the analysis of our data, the best fitting

model was the three-state power function. However, in this case the difference in

BIC scores between the three-state exponential function and the three-state power

Table 5
BIC scores for No-Learning and Learning models fit to Rickard (1997)

| | | 1 State | 2 States | 3 States |
|---|---|---|---|---|
| BIC Scores | Flat | 94,355 | 70,034 | 67,068 |
| | Power | 74,316 | 66,236 | 65,562 |
| | Exponential | 73,583 | 66,625 | 65,817 |

function is substantial; this probably reflected that the Rickard data involve many more items studied for twice as many trials.

The parameters of the three-state power model offer an interesting comparison with the estimates from our experiment.

I  -- Latency Intercept (shared by all states):  0 sec.

$\alpha$  -- Power Exponent (shared by all states):  .184

$\beta_1$ – Latency that speeds up in State 1: 5.65 sec.

$\beta_2$ – Latency that speeds up in State 2: 2.93 sec.

$\beta_3$ – Latency that speeds up in State 3: 1.76 sec.

$\pi_{12}$ -- Probability of transitioning from state 1 to 2 on a trial:  .060

$\pi_{23}$ -- Probability of transitioning from state 2 to 3 on a trial: .050

Again, the intercept I estimates to be 0 indicating that all the latency is subject to speedup, furthermore the exponent $\alpha$ is small indicating that most of the speedup is in the state transitions. The latencies $\beta_i$  are in similar ratio to our experiment but smaller, reflecting the somewhat simpler task. The probabilities of transitioning between states is considerably lower, probably reflecting the larger number of problems being repeated.

*Model Interpretation*

To what extent do our results support or challenge the models of skill acquisition? The three models of skill acquisition discussed in the Introduction differ not only in the number of abrupt changes between states due to the cognitive processes being used to solve the task, but also in an explanation of what learning processes underlie gradual within-state changes. The Instance Theory takes the perspective that skill acquisition occurs through a continuous build up of number of instances in memory (Logan, 1988; Compton & Logan, 1991, Logan 2002). Our work would suggest that this theory does not fully explain the state changes we identified, or the changes in neural activation that cannot be solely accounted for by practice. In contrast, our conclusions are consistent with Rickard's CMPL with respect to the first two stages, although we found a third stage consistent with Fitts & Posner's (1967) theory and its later ACT-R instantiation. We consider our results in terms of CMPL and ACT-R.

Both perspectives agree that calculation and retrieval are better modeled by separate learning curves. Concurrent reports indicate State 1 problems are solved with a mix of calculation and partial strategies, which in turn, are supported by an effect of problem difficulty that is only apparent in this state. Taken together, the data suggest that when participants are in State 1, they are performing computations to answer the problems. This is consistent with both the first state identified by CMPL and the Cognitive Phase of skill acquisition described by Fitts and Posner and later modeled in ACT-R. Results from the ROI analysis support this view as well. The increased likelihood of being in State 1 rather than State 2 is associated with greater activation in both the HIPS and LIPFC. Previous studies have

shown that HIPS activation in adults is associated with the calculation of arithmetic and solving more complex mathematical expressions (e.g. Delazer et al., 2003; Delazer et al., 2005; Grabner et al., 2009). The LIPFC has consistently been activated with increased retrieval demands (e.g. Wagner, Maril, Bjork, & Schacter, 2001; Dobbins, Foley, Schacter, & Wagner, 2002). In contrast to State 1, concurrent reports and a lack of a problem difficulty effect indicate a shift to retrieval in State 2. Given this interpretation, it might seem strange that we observe the increased likelihood of State 1 in comparison with State 2 as being associated with greater activation in the LIPFC. The multiple computations engaged in solving Pyramid problems (i.e. 5\$3 one might retrieve both 5+4 and 9+3), however, require many retrievals while a single fact retrieval is required for State 2 (i.e. recognizing 5\$3 as equaling 12).

Along with a similar identification the cognitive processes within State 1, CMPL and ACT-R also share similar explanations of the learning mechanisms supporting within-state speedup. The mixed effects model we used to analyze the ROIs allows us to disassociate the discrete shifts in cognitive processing associated with the change in state from the gradual shifts in cognitive processes associated with practicing a problem. Our ROI analysis indicates that while the changes between states accounts for the decreases in both HIPS and LIPFC activation, the impact of practice within-state also accounts for some of the decrease in LIPFC activation as well. This finding is in alignment with what we would expect given the mechanism of associative strengthening, a mechanism supported by both theories. ACT-R additionally supports the presence of knowledge compilation during the first Learning Phase (State 1) with the emergence of calculation short cuts. In this regard,

the prevalence of partial reports during State 1 suggests that these short cuts are present and causing speedups. CMPL, on the other hand, considers strategy choice as an issue that becomes important close to the boundaries between learning states. The simplicity of our concurrent reports, however, gives us little insight as to whether this occurs in our task.

It is notable that the predictions both theories make about the cognitive processes that occur in State 1 do not apply to novel problems. CMPL predicts problem-specific speedups due to practice, but not general speedup (Rickard, 1997). ACT-R theory would predict some general speedup due to overlapping practice of sub-procedures; however, some evidence suggests that this is less likely to occur when using a small highly practiced set of items (e.g. Anderson, Fincham & Douglass, 1997; Rabinowitz & Goldberg, 1995). Even though the learning exponent we estimated was shallow, it would predict a 2 second speed up in State 1 times over the course of the experiment, which is much more than the very minimal speed for novel problems.

Our HMM identifies a State 3, but unlike the ACT-R theory of skill acquisition (Anderson, 1982), the CMPL does not distinguish between the memory processes involved in retrieval and automaticity (Rickard, 1997). From our latency and concurrent reports, it is difficult to tell if the distinction between State 2 and 3 marks a meaningful change in the cognitive process being used; however, our fMRI results suggest a decrease in cognitive engagement associated with the switch to State 3. With enough practice the associative phase of skill acquisition is theorized to transition into the autonomous phase. This shift can be interpreted as the

transition from a stimulus- retrieval- response to a stimulus-response process (Fitts and Posner, 1967; Anderson, 1982). To test the hypothesis that State 2 reflects the associative phase and State 3 the autonomous phase, we performed a ROI analysis using an area involved in effortful retrieval *LIPFC,* and areas associated with visual recognition, *left fusiform gyrus,* and *left motor cortex.* In line with our predictions, we find that even when accounting for the effects of practice, the likelihood of being in State 2 as opposed to State 3 is associated with greater activity in the LIPFC. On the other hand, increased activation in the left motor cortex is associated with a higher likelihood of being in State 3 over State 2. This suggests that the cognitive processes in State 2 involve the retrieval of answer, whereas, by the time participants reach State 3, this retrieval has either been highly strengthened or even phased out. In effect, the problem solving is reduced to a stimulus-response task. Future work could benefit from developing methods for verifying whether the third learning state identified by our models is in fact the automization phase.

*Assumptions of our method, limitations, and future work*

In our HMM after transition to a new learning state, there is no option for backward transition to a prior state. We make our assumption to simplify our model. We do not have to include further estimated transition parameters to represent these possibilities. In order to keep track of different combinations of practice in different states, backward regressions would make the number of HMM states we need for n learning states equal to $36^n$ rather than 36n. While our model does not allow for backward transitions, there is evidence that they occasionally occur. Our concurrent reports contain instances of participants reporting

calculation after having reported retrieval. While not every report may be veridical, some occurred after our HMM assigned items to State 2 or 3 and are associated with long latencies. For instance, out of the 1080 trials with concurrent reports, there are 14 cases where probability of being in State 1 is less than 10%  (mean .8%), yet participants reported calculation, and had latencies longer than 4 seconds (mean 6.05) seconds. In the most extreme case, one participant on the last trial was assigned to State 3 with a probability of 1, and yet, reported calculation with a mean latency of 10.65 sec. While such forgetting might have been relatively rare in our experiment, accounting for forgetting within the model will be important in tasks that either span multiple days (e.g. Fincham, Anderson & Douglass, 1997), or contain a large amount of interference between trained problems (e.g. Anderson, 1981). Future work incorporating forgetting into the HMM will, however, provide valuable insight into understanding the role of forgetting within skill acquisition.

Another assumption we made is that while there is variability between individuals, there is also some universality in the range of time it takes to perform cognitive functions. This leads us to assume that the same function could describe learning for both items and individuals. Counter to this assumption, our investigation of the impact of problem difficulty shows that all items are not the same in State 1. There is previous work fitting functions to individual items (e.g. Heathcote, Brown & Mewhort, 2000; Myung, Kim & Pitt, 2000). It would not be trivial to extend that to our approach with multiple states because, as discussed earlier, we use the behavior across items to estimate the likelihood of various transition patterns for a particular item. Recent work by Anglim and Wynton (2015)

has suggested Hierarchical Bayesian modeling as a means of fitting item level data while also incorporating metrics penalizing for model flexibility, which allows comparison between models. Future work modeling states of skill acquisition could benefit from this technique.

*Conclusion*

The current study diverges from previous work fitting theory-driven models to data by using an unsupervised modeling technique that allowed the data to aid in the construction of the best model. We found that even accounting for between state changes, our model was improved by the inclusion of the power function to explain within-state learning speedup. Using behavioral and brain evidence, we illustrated that the three-state model identified by this paper captured both quantitative and qualitative changes in the cognitive processes required to perform the task. Distinction between State 1 and 2 aligned with predictions made by both CMPL and ACT-R, however the distinction between State 2 and State 3, which suggests a shift from retrieval to automization, is unique to ACT-R. While this method offers strengths well suited to a simple training task, future improvements could render greater resolution and sensitivity for detecting changes associated with the acquisition of cognitive skills.

**Acknowledgements**

References

Anderson, J. R. (1981). Interference: The relationship between response latency and response accuracy. *Journal of Experimental Psychology: Human Learning and Memory*, *7*(5), 326.

Anderson, J. R. (1982). Acquisition of cognitive skill. *Psychological Review*, 89(4), 369-406.

Anderson, J. R., Fincham, J. M., & Douglass, S. (1997). The role of examples and rules in the acquisition of a cognitive skill. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, *23*(4), 932.

Anderson, J. R., Bothell, D., Byrne, M. D., Douglass, S., Lebiere, C., & Qin, Y. (2004). An integrated theory of the mind. *Psychological Review,* 111(4), 1036-1060.

Anderson, J. R., (2007). *How can the human mind occur in the physical universe?* New York: Oxford University Press.

Anderson, J. R., Betts, S., Ferris, J. L., & Fincham, J. M. (2011). Cognitive and metacognitive activity in mathematical problem solving: prefrontal and parietal patterns. *Cognitive Affective Behavioral Neuroscience*, 11(1), 52–67.

Anderson, J. R., & Fincham, J. M. (2014, November). Extending problem-solving procedures through reflection. *Cognitive Psychology*, 74, 1–34.

Anglim, J., & Wynton, S. K. (2015). Hierarchical Bayesian Models of Subtask Learning. *Journal of Experimental Psychology: Learning, Memory, and Cognition, 41(4), 957-974.*

Badre, D., Wagner, A. D. (2005). Frontal lobe mechanisms that resolve proactive interference. C*erebral Cortex*, 15(12), 2003-2012.

Barr, D. J., Levy, R., Scheepers, C., & Tily, H. J. (2013, April). Random effects structure for confirmatory hypothesis testing: Keep it maximal. *Journal of Memory and Language,* 68(3), 255–278.

Bates, D., Maechler, M., & Bolker, B. (2011). Lme4: Linear mixed-effects models using S4 classes*.* R package version 0.999375-39.

Bajic, D., & Rickard, T. C. (2009, January). The temporal dynamics of strategy execution in cognitive skill learning. *Journal of Experimental Psychology: Learning, Memory, and Cognition,* 35(1), 113–121.

Bajic, D., & Rickard, T. C. (2011, October). Toward a generalized theory of the shift to retrieval in cognitive skill learning. *Memory & Cognition*, 39(7), 1147–1161.

Birn, R. M., Smith, M. A., Jones, T. B., & Bandettini, P. A. (2008). The respiration response function: the temporal dynamics of fMRI signal fluctuations related to changes in respiration. *NeuroImage*, 40(2), 644-654.

Chen, G., Saad, Z. S., Britton, J. C., Pine, D. S., & Cox, R. W. (2013). Linear mixed-effects modeling approach to FMRI group analysis. *NeuroImage*, 73, 176-190.

Clauset, A., Shalizi, C., & Newman, M. E. J. (2009). Power-law distributions in empirical data. *SIAM Review,* 51(4), 661-703*.*

Codding, R. S., Shiyko, M., Russo, M., Birch, S., Fanning, E., & Jaspen, D. (2007, December). Comparing mathematics interventions: Does initial level of fluency predict intervention effectiveness? *Journal of School Psychology*, 45(6), 603–617.

Compton, B. J., & Logan, G. D. (1991, March). The transition from algorithm to retrieval in memory-based theories of automaticity. *Memory & Cognition*, 19(2), 151–158.

Cox, R. W. (1996). AFNI: software for analysis and visualization of functional magnetic resonance neuroimages. *Computers and Biomedical Research*, 29(3), 162-173.

Dehaene, S., Piazza, M., Pinel, P., & Cohen, L. (2003, May). Three parietal circuits for number processing. *Cognitive Neuropsychology*, 20(3), 487-506.

Delaney, P. F., Reder, L. M., Staszewski, J. J., & Ritter, F. E. (1998). The strategy specific nature of improvement: The power law applies by strategy within task. *Psychological Science*, 9, 1-7.

Delazer, M., Domahs, F., Bartha, L., Brenneis, C., Lochy, A., Trieb, T., & Benke, T. (2003). Learning complex arithmetic—an fMRI study. *Cognitive Brain Research*, 18(1), 76–88.

Delazer, M., Ischebeck, A., Domahs, F., Zamarian, L., Koppelstaetter, F., Siedentopf, C. M., Kaufmann, L., Benke, T., & Felber, S. (2005). Learning by strategies and learning by drill--evidence from an fMRI study. *NeuroImage*, 25(3), 838–849.

Dobbins, I. G., Foley, H., Schacter, D. L., & Wagner, A. D. (2002). Executive control during episodic retrieval: multiple prefrontal processes subserve source memory. *Neuron*, 35(5), 989-996.

Fitts, P. M., & Posner, M. I. (1967). Human performance. Brooks/Cole.

Friston, K. J., Fletcher, P., Josephs, O., Holmes, A., Rugg, M. D., & Turner, R. (1998). Event-related fMRI: characterizing differential responses. *NeuroImage*, 7(1), 30-40.

Grabner, R. H., Ischebeck, A., Reishofer, G., Koschutnig, K., Delazer, M., Ebner, F., & Neuper, C. (2009). Fact learning in complex arithmetic and figural-spatial tasks: the role of the angular gyrus and its relation to mathematical competence. *Human Brain Mapping*, 30(9), 2936–2952.

Haider, H., & Frensch, P. A. (2002). Why aggregated learning follows the power law of practice when individual learning does not: Comment on Rickard (1997, 1999), Delaney et al. (1998), and Palmieri (1999). *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 28(2), 392–406.

Heathcote, A., Brown, S., & Mewhort, D. J. (2000, June). The power law repealed: the case for an exponential law of practice. *Psychonomic Bulletin & Review*, 7(2), 185–207.

Hecht, S. A. (2006). Group differences in adult simple arithmetic: Good retrievers, not-so-good retrievers, and perfectionists. *Memory & cognition*,*34*(1), 207-216.

Jones, R. H. (2011). Bayesian information criterion for longitudinal clustered data. *Statistics in Medicine,* 30 (25)*, 3050-3056.

Judd, C. M., Westfall, J., & Kenny, D. A. (2012, July). Treating stimuli as a random factor in social psychology: a new and comprehensive solution to a pervasive but largely ignored problem. *Journal of Personality and Social Psychology*, 103(1), 54–69.

Kass, R. E., & Raftery, A. E. (1995). Bayes factors. *Journal of the American Statistical Association*, *90*(430), 773-795.

Kenward, M. G., & Roger, J. H. (1997). Small sample inference for fixed effects from restricted maximum likelihood. *Biometrics*, 53(3), 983–997.

Kilpatrick, J., Swafford, J., & Findell, B. (2001). *Adding it up: Helping students learn mathematics*. Washington, DC: National Academy Press.

Kliegl, R., Masson, M. E., & Richter, E. M. (2010). A linear mixed model analysis of masked repetition priming. *Visual Cognition*, 18(5), 655-681.

Logan, G. D. (1988). Towards an instance theory of automatization. *Psychological Review*, 95(4), 492–527.

Logan, G. D. (2002). An instance theory of attention and memory. *Psychological Review*, 109(2), 376–400.

Myung, I. J., Kim, C., & Pitt, M. A. (2000, July). Toward an explanation of the power law artifact: insights from response surface analysis. *Memory & Cognition*, 28(5), 832–840.

McCandliss, B. D., Cohen, L., & Dehaene, S. (2003). The visual word form area: Expertise for reading in the fusiform gyrus. *Trends in Cognitive Sciences*, 7(7), 293-299.

Newell, A. & Rosenbloom, P. S. (1981). Mechanisms of skill acquisition and the law of practice. In J.R. Anderson (Ed.), *Cognitive Skills and Their Acquisition* (pp.1-55). Hillsdale, NJ: Erlbaum.

Nosofsky, R. M., & Palmeri, T. J. (1997). An exemplar-based random walk model of speeded classification. *Psychological Review*, 104(2), 266–300.

Palmeri, T. J. (1999). Theories of Automaticity and the Power Law of Practice. *Journal of Experimental Psychology*: *Learning, Memory, and Cognition*, 25(2), 543–551.

Rabiner, R. E. (1989). A tutorial on Hidden Markov Models and selected applications in speech recognition. *Proceedings of the IEEE,* 77(2), 257-286.

Rabinowitz, M., & Goldberg, N. (1995). Evaluating the structure-process hypothesis. In F. Weinert and W. Schneider (Eds.), *Memory performance and competencies: Issues in growth and development*, 225-242. Hillsdale, NJ:Erlbaum.

Rickard, T. C. (1997). Bending the power law: A CMPL theory of strategy shifts and the automization of cognitive skills. *Journal of Experimental Psychology: General,* 126(3), 288-311.

Rickard, T. C. (2004). Strategy execution in cognitive skill learning: an item-level test of candidate models. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 30(1), 65–82.

Shrager, J., & Siegler, R. (1998). SCADS : A Model of Children ' s Strategy Choices and Strategy Discoveries. *Psychological Science*, 9(5), 405–410.

Siegler, R. S. (1988). Individual differences in strategy choices: Good students, not-so-good students, and perfectionists. *Child Development*, 833-851.

Siegler, R. S., & Shipley, C. (1995). Variation, selection, and cognitive change. In Simon, T.& Halford, G. (Eds.) *Developing cognitive competence: New approaches to process modeling* (pp 31-76).

Singer-Dudek, J., & Greer, R. D. (2005). A long-term analysis of the relationship between fluency and the training and maintenance of complex math skills. *Psychological Record,* 55(3), 361-376.

Skinner, C. H., Fletcher, P. A., & Henington, C. (1996). Increasing learning rates by increasing student response rates: A summary of research. *School Psychology Quarterly*, *11*(4), 313.

Spiess, A. N., & Neumeyer, N. (2010). An evaluation of R2 as an inadequate measure for nonlinear models in pharmacological and biochemical research: a Monte Carlo approach. *BMC pharmacology*, *10*(1), 6.

Taatgen, N. A., & Lee, F. J. (2003, January). Production compilation: a simple mechanism to model complex skill acquisition. *Human Factors*, 45(1), 61–76.

Tenison, C., Fincham, J. M., & Anderson, J. R. (2014, February). Detecting math problem solving strategies: an investigation into the use of retrospective self-reports, latency and fMRI data. *Neuropsychologia*, 54, 41–52.

Touron, D. R., & Hertzog, C. (2004). Strategy shift affordance and strategy choice in young and older adults. *Memory & Cognition*, *32*(2), 298-310.

Wagner, A. D., Maril, A., Bjork, R. A., & Schacter, D. L. (2001). Prefrontal contributions to executive control: fMRI evidence for functional distinctions within lateral prefrontal cortex. *NeuroImage*, 14(6), 1337-1347.

Appendix A

*Simulation Analysis*

We ran a simulation analysis to explore the sensitivity of our HMM modeling technique for discovering the true number of states underlying a dataset, and for identifying the shape of the learning within the states. We simulated data for a three- parameter power function, a three-parameter exponential function and a flat line. We did not explore the APEX function because of its over-complexity and low performance in our initial model fitting (Table 2). For each function we generated samples reflecting one-,two- and three- state models of skill acquisition. The parameters used to generate the simulated data were the best-fitting parameters for the actual data. The transition probabilities were used to determine when participants switched from one state to the next. The simulated latency data were generated by sampling data distributed around previously fit functions according to a gamma with shape of 3 and scale of 1/3 predicted response latency. We generated 'problem solving latencies' for 60 items which were practiced 36 times for one- through three- state models which used each of the three functions of interest. There were 100 such Monte Carlo simulations for each of the 9 models (3 numbers of states by 3 learning functions). We then fit the data using the same HMM fitting technique described in the Results and initialized by the same parameters as we used when modeling our experimental data.

We used BIC to evaluate the performance of the models for each of the 100 items in a simulated set. The actual data were best fit by the three- state exponential and power functions with almost no difference between the two. Table A1 shows the results from the fits to the 900 simulations (three- states times three-learning

functions times 100 datasets). The simulation results indicate that when a dataset is best fit by one of these two models it was always generated by such a model, although there is relatively poor ability to distinguish between the two. Interestingly, 83 of the simulated three-state exponential data and 46 of the simulated three-state power datasets are best fit by 3 state flat no-learning functions. This reflects the rather shallow learning parameters that we estimated. In any case, we can conclude that our result of best fits by 3 state-exponential and power functions are unlikely to have been generated by fewer states or by a flat learning function.

Table A
Counts of best fitting models for simulated datasets

| | | | Generating Function & Generating Number States | | | | | | | | |
| | | | Exponential | | | Power | | | No Learning | | |
| | | | 1 State | 2 States | 3 States | 1 State | 2 States | 3 States | 1 State | 2 States | 3 States |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Testing Function & Testing number of States | Exponential | 1 State | 98 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | | 2 States | 0 | 75 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | | 3 States | 0 | 0 | 15 | 0 | 0 | 20 | 0 | 0 | 0 |
| | Power | 1 State | 0 | 0 | 0 | 99 | 0 | 0 | 0 | 0 | 0 |
| | | 2 States | 0 | 0 | 0 | 0 | 78 | 1 | 0 | 0 | 0 |
| | | 3 States | 0 | 0 | 2 | 0 | 0 | 33 | 0 | 0 | 0 |
| | No Learning | 1 State | 0 | 0 | 0 | 0 | 0 | 0 | 100 | 0 | 0 |
| | | 2 States | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 100 | 0 |
| | | 3 States | 2 | 25 | 83 | 1 | 22 | 46 | 0 | 0 | 100 |

Appendix B

The HMM calculates the probability of the latencies for an item given a set of

parameters. Through expectation maximization iteration it determines the

parameters that maximize the likelihood of the data on all the trials. These

calculations must consider all the possible trajectories of the item through the states

of the HMM. It can do so with reasonable efficiency given the dynamic programming

algorithm that is being used.

The probability of a set of $N$ observed latencies $t_1 \ldots t_N$ for an item can be

specified with a set of recursive equations. Let $Pr(j,i)$ denote the probability of the

latencies from trials $j$ to the last trial given that learning state $i$ is entered on trial $j$.

Since the item starts in learning state 1 on trial 1 the probability of the N latencies

for that item is $Pr(1,1)$. The following is the recursive formula for $Pr(j,i)$:

$$
\begin{aligned}
\Pr(j,i) = \sum_{k=1}^{N-j} & \left[ (1-\pi_{i,i+1})^{k-1} \pi_{i,i+1} \left( \prod_{m=1}^{k} g(t_{j+m-1}, 3, \frac{T_{im}}{3}) \right) \Pr(j+k, i+1) \right] \\
& + (1-\pi_{i,i+1})^{N-j} \left( \prod_{m=1}^{N-j+1} g(t_{j+m-1}, 3, \frac{T_{im}}{3}) \right)
\end{aligned}
$$

This equation adds the probabilities of the observed latencies given various number

of trials in state $i$. The expression within the summation sign deals with all the cases

where there is a transition to the next state and the final term deals with the case

where the item stays in state $i$ until the last trial $N$. Within the summation sign, *(1-*

$\pi_{i,i+1})^{k-1}\pi_{i,i+1}$ is the probability of spend $k$ trials in learning state $i$; $g(t_{j+m-1},3,T_{im}/3)$ is the gamma-3 probability of the observed latency $t_{j+m-1}$ on trial $j+m-1$ given the expected latency $T_{im}$ for the $m^{\text{th}}$ trial in state $i$; $\Pr(j+k,i+1)$ is the probability of the latencies from trials $j+k$ to $N$ given that the item transition to the next state $i+1$ on trial $j+k$.

Note in the case of the last state the probability, $\pi_{i,i+1}$ of transitioning to a further state, is 0 and the whole expression just reduces to the last term outside of the summation sign for spending the remaining trials in that state. These equations make reference to expected times, $T_{im}$, for different numbers of trials in a state. These expected times are calculated according to the particular model (no learning, exponential, power, APEX) being considered.

Appendix C

To provide an impression of how the different 1, 2 and 3 State HMMs fit the

item level data, we present several plots imposing the predicted model fits on the

item level latency scores. We wanted to provide a balanced impression of the trends

present within our sample of 60 items (20 participants x 3 problems each). We used

k-means clustering to divide the items into three groups that spent similar amounts

of time within each of the three states. Figure C shows the latency over the 36

practice opportunities of 6 items. We chose two items that were closest to the

centroid of each cluster. Cluster 1(Figure C.A and C.B), featured items (n= 29) that

had the most observations per item in State 3 (M 28.8, SD=4). Items (n=15) in

Cluster 2 (Figure C.C and C.D) were observed nearly equally between State 2

(M=15.8 SD=4.4) and State 3 (M=15.5, SD=4.4). Cluster 3 (Figure C.E and C.F)

favored items (n=16) were observed most in State 2 (M=25.5, SD=8.9). Figure C also

shows the predicted latencies for 1, 2 and three-state power functions assuming

that the item was in its most probable state on a trial. It needs to be emphasized that

in some cases there is considerable uncertainty about what state an item was in on a

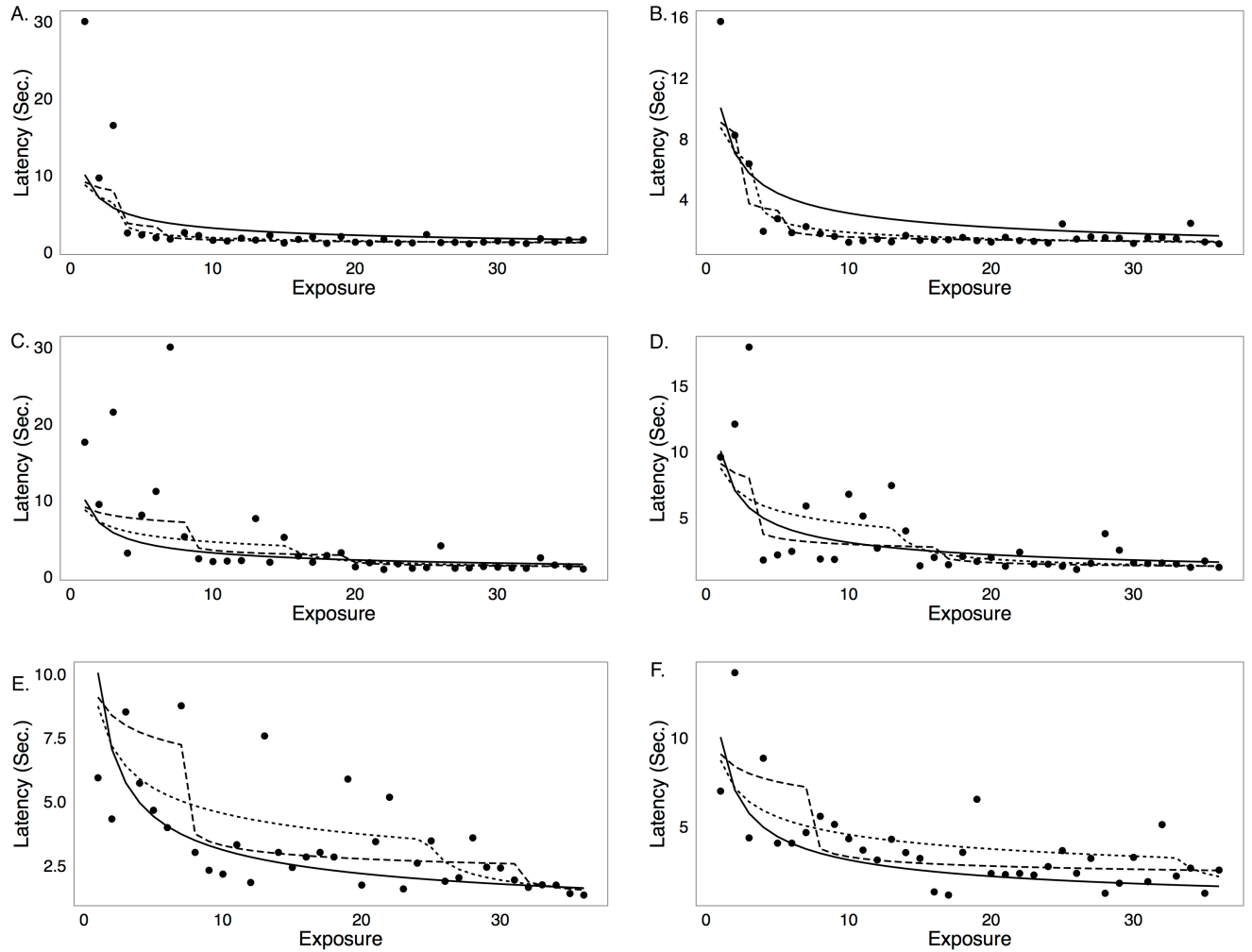trial and there are other likely state trajectories besides the ones displayed.

*Figure C.* Latency plots for two practiced items (columns) from each of the three clusters (rows). For each item, the latency scores for each practice opportunity are plotted along side predicted latency generated by a one-state (solid line), a two-state (dotted line) and a three-state model (dashed line).

# Appendix D

Table D
List of practiced and novel problems used in the experiment

|  |  | Height 3 | Height 4 | Height 5 |
|---|---|---|---|---|
|  | 4 | 4$3 = 9 | x | x |
|  | 5 | 5$3 = 12 | 5$4 = 14 | x |
|  | 6 | 6$3 = 15 | 6$4 = 18 | 6$5 = 20 |
| Base | 7 | 7$3 = 18 | 7$4 = 22 | 7$5 = 25 |
|  | 8 | 8$3 = 21 | 8$4 = 26 | 8$5 = 30 |
|  | 9 | 9$3 = 24 | 9$4 = 30 | 9$5 = 35 |
|  | 10 | 10$3 = 27 | 10$4 = 34 | 10$5 = 40 |
|  | 11 | 11$3 = 30 | 11$4 = 28 | 11$5 = 45 |