

A Simulation of Memory for Computer Programs

FERNAND GOBET AND IAIN OLIVER

ESRC Centre for Research in Development, Instruction and Training

School of Psychology, University of Nottingham

University Park, Nottingham NG7 2RD

United Kingdom

email: frg@psychology.nottingham.ac.uk

ESRC Centre for Research in Development, Instruction and Training

March 2002

Technical Report No. 74

Abstract

Both for practical and theoretical reasons, it is important to understand the type of knowledge acquired by computer programmers, and how it is used to solve problems. A substantial amount of research has been carried out to characterise programmers' knowledge and memory, a subset of which is reviewed in this paper. While several theoretical explanations have been advanced to explain the knowledge and memory of programming experts, no computational model has been developed so far to simulate the empirical data. We describe how a simplified version of CHREST, a computational model originally developed to account for perception, learning, and memory in chess, can simulate the main effects found in research into computer-program memory. While it captures the skill differences found even in cases where substantial randomisation had been done, it fails to capture within-skill effects of variations at the semantic level.

Email: frg@psychology.nottingham.ac.uk
Phone 0115 9515402, Fax 0115 9515324

1 INTRODUCTION

Computer programming involves a variety of skills: the ability to understand programs written by others, to design, write and debug one's own programs, and to use problem-solving strategies to turn a set of constraints and desiderata into a correct and running program. Several general theories have been proposed to account for the presence of these abilities. Some authors (e.g., Adelson, 1981; McKeithen, Reitman, Rueter & Hirtle, 1981; Ye & Salvendy, 1994) have proposed that semantic knowledge plays an essential role. Others have proposed that analysing the characteristics of the (programming) environment offers the best starting point (e.g., Vicente & Wang, 1998). Finally, it has also been proposed that expertise in programming, like in other domains, stems from the acquisition of a large number of perceptual chunks, which are the building blocks on which later semantic and procedural knowledge is constructed (e.g., Chase & Simon, 1973; Gobet & Simon, 1996).

Chase and Simon's theory was mostly grounded in data from research into chess expertise. Studying a Master, a good amateur and a novice in a variety of tasks, they found good evidence for the psychological reality of chunks. In particular, they found that there was a massive skill effect for the recall of positions taken from Master games, but that this effect disappeared with random positions. (It has recently been shown that chess Masters keep a small, but reliable, superiority even with random positions (Gobet & Simon, 1996), a result that will be of interest in our review of papers devoted to memory for computer programs.)

Probably due to the fact that it was more specified than alternative explanations—a subset of the theory was implemented in a computer program by Simon & Gilmarin (1973)—Chase and Simon's chunking theory has spawned a large number of experimental studies. Several of these studies have been carried out in the domain of computer programming, and the importance of chunking in programming is generally accepted (e.g., Adelson, 1981; Bartfield, 1986; McKeithen et al., 1981). In addition, Schmidt (1986) has found that recall of computer programs correlates with their comprehension.¹ Given that high comprehension is a distinguishing feature of expertise in programming, this correlation suggests that chunks, as measured by the recall task, may play a causal role. However, no computational model has been developed so far to simulate the empirical data about memory for computer programs. The goal of the present paper is to fill in this gap, using as a modelling environment the CHREST architecture, which emphasises the role of perceptual chunks.

¹A similar result has also been found in chess research (e.g., De Groot & Gobet, 1996). However, an important difference between the two domains should be noted: due to their layout and their lack of semantic transparency, computer programs cannot be grasped 'at a glance' during a brief presentation the way chess positions can.

²This is far from being a foolproof method. From research into other domains of expertise, it is known that experience correlates only imperfectly with expertise level (Ericsson, 1991; Richman et al., 1996). For example, in spite of decades of practice and formal instruction, many a chessplayer still plays weakly.

For all the experiments, the materials were examples of genuine computer code. Each experiment only used one programming language to draw its examples from, even if some of the participants knew more than one language. But unfortunately, there are no two experiments using the same language so as to allow direct comparison of results. Indeed, differences in the languages, experimental designs, and scoring methods make detailed comparisons awkward.

Experimental Material

Unlike similar endeavours into the domain of chess, there is no standard rating scale that measures the level of expertise of computer programmers. For the four experiments to be reviewed, the level of expertise was determined by the participants' experience of programming.² In general, the participants in the novice group had some experience of the programming language used in the experiment. The expert group usually consisted of programmers that had completed, or lectured on, courses in the language. Some programmers were rated as experts because they had experience in more languages than the target one in the experiment.

Assessing Programming Ability

Several studies have been carried out to investigate memory for computer code by individuals of different levels of expertise. While some of the studies were also interested in cognitive processing differences, this review will focus upon the studies where De Groot's (1965) recall task has been used—that is, a brief presentation of material taken from the domain of expertise, and a subsequent test of memory. In selecting the studies, we have also used the criterion that the experiment should compare the memory recall ability of expert and novice programmers, and that some measure of performance (e.g., percentage of lines correctly recalled) was provided. These criteria resulted in the selection of four studies (Adelson, 1981; Bartfield, 1986; Bateson, Alexander & Murphy, 1987; and Guerin & Matthews, 1990). Important features in these experiments include the participants' skill level and the type of stimulus given to them.

2 MEMORY FOR PROGRAMS: A BRIEF REVIEW

The paper is organised as follows. First, we briefly review research on memory for computer programs. Second, we describe a computer simulation using the CHREST architecture. Third, we compare the results of the simulations with those obtained with humans. Finally, we reflect upon the impact of our results upon research into expertise.

Some authors (e.g., Guerin & Matthews, 1990) even go so far as to criticise other experimenters for their choice of target language.

Summaries of Experiments

Adelson (1981). This experiment addressed the question of how experts represent and use programming concepts. It tried to show that experts use a hierarchy to organise their information and base its structure upon functional aspects of programming. This is opposed to novices who, according to Adelson, organise information based on the program syntax.

The experiment used Polymorphic Programming Language (PPL). PPL is a variant of PL/I, which is a combination of FORTRAN, ALGOL and COBOL (see Schmidt, 1986). The novices were five undergraduates who had completed a course in PPL, and the experts were five lecturers in that language. Sixteen lines of PPL code taken from three separate, complete programs were used as stimuli. Each line of code was presented separately on a screen. Lines were presented in a random order and each line was visible for 20 seconds. After all lines had been presented, the participants had 8 minutes to recall the code. This procedure was repeated for nine trials.

Experts recalled more than the novices (see Figure 1). Thus, the results were inconsistent with the conclusions drawn from Chase and Simon's chess studies, where no skill difference was found with random material. Adelson suggests that the discrepancy between the Chase and Simon chess data and her experiment comes from the fact that the code consists in lines taken from three complete programs and not in lines randomly selected from 16 different programs. Moreover, these results are in line with the recent finding in chess expertise of a skill effect even with briefly-presented random positions (Gobet & Simon, 1996).

Adelson also looked at the size of chunks used in recall, defining a chunk as a sequence of items recalled in succession with less than a 10-second pause between them. Experts' chunk size was greater than novices' (on average 3.5, and 2.4 items, respectively). Based on these and additional results obtained with multidimensional scaling and cluster analysis, Adelson concluded that experts organise information using functional principles, while novices categorise on a more syntactic (surface) basis.

Barfield (1986). Barfield was interested in distinguishing novice from expert problem-solving behaviour and knowledge acquisition, and concentrated on chunking as the main process that separates experts from novices. He suggested that programmers take in the complicated stimuli as meaningful chunks before they are processed.

Four levels of expertise were used. Naïve participants (n=42) had not completed any programming courses. Novices (n=80) had completed just one course in BASIC. Intermediates

(n=73) had completed a minimum of one BASIC course plus two or three courses in other languages. Experts (n=26) were graduates in computer science as well as having at least one course in BASIC.

The material consisted of one 25-line program written in BASIC. The experimenters identified likely modules within the program, but no visible boundaries were marked (i.e. no spaces between lines). The experiment had three conditions: the stimulus could be presented either (a) in executable order, (b) with the order of lines randomised, or (c) with the order of modules randomised (in that case, the lines within a module preserved their order). The participants were allowed three minutes to study the stimulus and four minutes to recall it.

The results are summarised in Figure 1. Naïve and novice participants obtained the same level of performance regardless of stimulus type, indicating that little, if any, of the chunk knowledge possessed by experts is present with novices. According to the results, intermediates can chunk together lines of code as long as they are in executable order. As expected, randomising the lines did affect the performance of the experts, although they still did better than Novices and Naïve participants. The randomising of modules did not affect the expert performance and this was taken as support for Barteld's chunking explanation. However, Guerin and Matthews (1990) argue that Barteld is measuring recall and not comprehension, so even though the semantic structure of the program is tampered with, it will not show in the results because BASIC programmers are not as sensitive to the semantic complexity as programmers using other languages. They also criticise Barteld for not randomising lines within modules to complete his experimental design.

Bateson, Alexander and Murphy (1987). This paper aimed to expose expert-novice differences in syntactic and semantic memory, and tactical and strategic skill. The gist of Bateson's experiment is to demonstrate the importance of semantics in gauging programmer differences over tasks that use measures of memory and chunk size in syntactic recall. Only the first of Bateson's battery of tasks, the syntactic memory task, is of interest to this review.

Two groups were used for this task: novices (n=20), who had completed no more than three programming courses, and experts (n=30), who had completed more than three courses. All participants had completed 12 weeks in an introductory FORTRAN class. The material used for the experiment was four short programs of equal length written in FORTRAN. Two of the four programs had lines randomised.

Participants were given only one normal and one random program, and were given three minutes to study a program and then allowed four minutes of free recall to write down what they remembered. The means for the proportion of total program recall are shown in Figure 1. As with Adelson (1981), there is a skill effect even when the order of the lines is randomised.

material is robust.

³Doing a meta-analysis of these and other studies, Feddon (2000) has found that this skill effect with scrambled

The net is grown by two EPAM-like learning mechanisms, *familiarisation* and *discrimination*. When a new object is presented to the model, it is sorted through the discrimination net. When the net or constructed (see Gobet & Simon, 2000, for detail).

largest chunk met at any point in time (the *hypothesis*), is kept in STM until a larger chunk is STM, which consists of four chunks, is mostly a queue (first-in first out). However, the CHERST consists of the following components: recognition LTM, semantic LTM, and STM. of language (see Gobet et al., 2001, for a review).

memory, use of multiple representations in physics, letter perception, spelling and acquisition used to account for domains such as verbal behaviour, chess memory, expert digit-span handling information in short-term memory (STM). Together, EPAM and CHERST have been long-term memory (LTM) through the construction of a discrimination net and mechanisms for Simon 1995). At the core of EPAM and CHERST lie mechanisms for encoding chunks into Memorizer) cognitive architecture (Feigenbaum & Simon, 1984; Richman, Staszewski & 1993; Gobet & Simon, 2000) is an expansion of the EPAM (Elementary Perceiver And CHERST (Chunk Hierarchy and Retrieval Structures; De Groot & Gobet, 1996; Gobet, emphasises perceptual chunking, can account for the empirical data.

we wish to explore to what extent a computational implementation of the theory, which explanation of the skill effect found in memory tasks for computer programs. In this section, As mentioned above, the chunking theory has often been proposed, in its informal form, as an

3 THE CHERST ARCHITECTURE

experts have larger chunks than novices. correlation between recall and comprehension. Finally, some evidence has been uncovered that also present when the code has been scrambled in various ways.³ Third, there is a robust experts always do better than non-experts with executable code. Second, the expertise effect is **Summary.** Several phenomena clearly stand out from these experiments. First, as expected,

randomised. better performance than novices could not be used when the lines of the program were superior semantic knowledge and comprehension skills of experts that allowed them to obtain of the program and of how this purpose was achieved. Guerin and Matthews conclude that the comprehension, which was measured by having participants write a summary of the purpose Lines and Modules (see Figure 1). Finally, they found that recall correlated highly with

language. During training, CHREST is given computer programs in FORTRAN—a naturalistic material—as input so that the vocabulary of the language as well as some sequences of vocabulary items can be learnt. The lines of code had a mean length of 7 words. Elements (this includes numbers, punctuation and other special characters) are recognised as distinct, individual items by the model. This type of input affords the model the ability to build a discrimination net that encodes both the primitive items and legal strings from the computer

Training Phase

For training and testing the model, a large collection of data was gathered from a variety of sources. Using the internet and some reference books, a corpus of about one hundred different FORTRAN programs was built.

4 SIMULATIONS

Cognitive operation	Duration
creating an LTM chunk	8 s
familiarising an LTM chunk	2 s
placing a symbol into STM	50 ms
comparing two symbols	50 ms
carrying out a test in the net	10 ms

Table 1: Main time parameters used in CHREST

data. Therefore, we did not use templates in the simulations. This paper, we are primarily interested in how far perceptual chunks can account for the human the theory) allowing information to be rapidly encoded in slots (Gobet & Simon, 2000). In CHREST incorporates mechanisms for incrementally creating schemas (known as templates in familiarising it occur in parallel to the other operations.

stimulus in the test phase. Note that creating an LTM chunk, adding a new link to a chunk, or processing can be done both during the training phase and during the presentation of the 2000) and are important in that they impose stringent constraints on how much information previous work (Feigenbaum & Simon, 1984; De Groot & Gobet, 1996; Gobet & Simon, Table 1 shows the key time parameters used with CHREST. These parameters are taken from a node is reached, the object is compared with the *image* of the node, which is its internal representation. If the image under-represents the object, new features are added to the image (familiarisation). If the information in the image and the object differ on some feature or some sub-element, a new node is created (discrimination).

The same basic model is used to simulate different levels of ability; that is, only the amount of input is varied, and no other mechanisms or parameters are altered. This study will focus upon the difference between *novices* and *experts*, who are simulated by passing CHESTST either one program or a corpus of eighty-eight programs during training.

As each item is passed to CHESTST, the model constructs its discrimination net. The net initially starts with an empty root node. Primitive items are usually the first to be added to the net. Then, after a period of learning, the images at the nodes will come to represent sequences of items.

Test Phase

Once the appropriate training had been done (study of one program for novices and eighty-eight for experts), the novice and expert level models were tested using twelve new FORTRAN programs that were not included in the training set. Test programs were selected that did not have too many "print" statements in them and so that they were all of roughly the same number of lines and words. The simulations were run like an experiment with human participants, with a presentation time of five seconds. Various levels of randomisation were applied to the test programs before they were passed to the models. In addition to the conditions used in the studies reviewed above, we also thought it interesting to use a condition in which all elements of a program were randomised. There were therefore five conditions in total:

1. **Normal.** The sequence of the program is unaltered.
2. **Random Modules.** Segments of the program are randomised, but the line order within a segment is retained.
3. **Random Lines.** The lines of the program are randomised. Information within a line is unaltered.
4. **Modules and Lines.** Both modules and lines within modules are randomised.
5. **All Random.** All elements within a program are randomised, yielding a total randomisation.

To achieve the Random Modules, lines of code that acted as a meaningful unit of instruction were grouped. For example, lines belonging to declaration statements would be retained together as a module. For the All Random condition, the maximum and minimum line lengths of the original program were first noted; then, all the elements in the program were randomised and lines of random length were constructed, with the condition that the values fell between the original lengths.

Each model received two "practice" problems which were always the same. The first practice program was a Normal type and the second was an All Random type. The test programs were

then presented. In addition to the practice problems, each model received two examples of each condition, thus making a total of twelve test programs. To control for random variation due to the order of programs in the learning set, CHREST was run with forty simulations per skill level. The random order in which the programs appeared, the random order in which conditions appeared and the randomisation of programs were all reset for each simulation.

For each program, CHREST read the program line by line, storing recognised chunks into STM, and, when applicable, using the following learning mechanisms. First, as described before, CHREST can add a chunk as a test to another chunk. It takes 8 seconds to carry out this discrimination operation. Typically, a new test is added to the hypothesis. Second, for chunks that have been in STM for at least 8 seconds, a new branch is added to access them by a novel path; this essentially means that episodic cues that permit access to this node are added to the discrimination net (Gobet & Simon, 2000). Such nodes can be recalled during the reconstruction phase even if they are no longer in STM.

During the recall phase, CHREST could output the information held in STM and in the nodes that had been created or for which new access links had been created. Recall was scored in the following way. A list of items that CHREST had recalled, in the order they were retrieved, was collected from the model. This recall list was matched alongside the original stimulus that was presented to CHREST. The first line in the stimulus was compared to the recall list and if the first items matched, then the lines were compared to find out how many of the items were recalled correctly before a mismatch occurred; both lines were then discarded and the next stimulus line compared to the recall. If the stimulus line did not match the recall line, then the line was discarded and the next one matched against the recall list, until all the stimulus lines were used. Not only does this method show how many items were recalled correctly, but it shows how many errors of commission the model made.

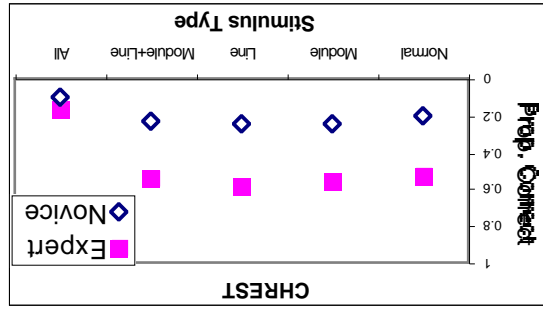


Figure 2. CHREST's memory for computer code as a function of level of expertise and type of randomisation.

Results

Figure 2 illustrates the results of the simulations. We can see that for all cases, CHREST predicts a skill effect. These predictions are borne out by the data, with the exception of Guerin and Matthews' (1990) Module+Line condition, where no effect was found with the human participants. None of the studies reviewed incorporated the All condition, where the order of all elements of the code is randomised. Although the difference between the Novice and Expert models is small for the All condition, it is statistically reliable ($t(78) = 3.71, p < .001$). To test this counter-intuitive prediction of the model, we collected data with C programmers ($n=9$) and novices ($n=9$); the results have supported the prediction. Given that we found a skill effect with full randomisation, it is unclear why no such effect was found with Guerin and Matthews' Module+Line condition, which destroys less structure than our method. While the simulations show a differential effect for the type of randomisation, this effect is limited to the All conditions vs. the other conditions. The model fails to capture the differential recall shown by humans from the Normal condition to the Module+Line conditions. It is likely that humans pick up semantic information from the modules, as proposed for example by Adelson (1981), which are beyond the essentially perceptual knowledge that this simplified version of CHREST can store.

5 CONCLUSION

The experimental studies reviewed in this paper clearly show that randomisation of aspects of computer code affects recall, while still preserving a skill effect in most cases. In order to account for these data, we have described a simplified version of the CHREST simulation model, where emphasis was given to perceptual chunking. We have shown that the model accounts for the skill differences in the data, which supports Simon and Gobet's (2000) contention that theories based on chunking mechanisms can account for skill effects in memory for computer programs. However, the model did not show similar differences between the randomisation conditions, as humans did. Further work will establish whether the presence of templates (Gobet & Simon, 2000) can help the model capture this result, which is often claimed to tap into differences in high-level, semantic knowledge (e.g., Adelson, 1981).

REFERENCES

Adelson, B. (1981). Problem solving and the development of abstract categories in programming languages. *Memory and Cognition*, 9, 422-433.

Barfield, W. (1986). Expert-novice differences for software: Implications for problem-solving and knowledge acquisition. *Behaviour and Information Technology*, 5, 15-29.

Bateson, A. G., Alexander, R. A., & Murphy, M. D. (1987). Cognitive processing differences between novice and expert computer programmers. *International Journal of Man-Machine Studies*, 26, 649-660.

Chase, W. G., & Simon, H. A. (1973). Perception in chess. *Cognitive Psychology*, 4, 55-81.

de Groot, A. D. (1965). *Thought and choice in chess*. The Hague: Mouton. (First edition in Dutch, 1946).

de Groot, A. D., & Gobet, F. (1996). *Perception and memory in chess*. Assen: Van Gorcum.

- Ericsson, K. A., & Smith, J. (1991). Prospects and limits of the empirical study of expertise: An introduction. In K. A. Ericsson & J. Smith (Eds.), *Studies of expertise: Prospects and limits*. Cambridge: Cambridge University Press.
- Feddon, J. (2000). *Chunking in programming: A meta-analysis*. Unpublished manuscript. Dept. of Psychology, University of Tallahassee (Florida).
- Feigenbaum, E. A., & Simon, H. A. (1984). EPAM-like models of recognition and learning. *Cognitive Science*, 8, 305-336.
- Gobet, F. (1993). A computer model of chess memory. *Proceedings of 15th Annual Meeting of the Cognitive Science Society*, (pp. 463-468). Hillsdale, NJ: Erlbaum.
- Gobet, F., Lane, P.C.R., Croker, S., Cheng, P.C-H., Jones, G., Oliver, I., & Pine, J. (2001). Chunking mechanisms in human learning. *Trends in Cognitive Sciences*, 5:236-243.
- Gobet, F., & Simon, H. A. (1996). Templates in chess memory: A mechanism for recalling several boards. *Cognitive Psychology*, 31, 1-40.
- Gobet, F., & Simon, H. A. (2000). Five seconds or sixty? Presentation time in expert memory. *Cognitive Science*, 24, 651-682.
- Guerin, B., & Matthews, A. (1990). The effects of semantic complexity on expert and novice computer program recall and comprehension. *The Journal of General Psychology*, 117, 379-389.
- McKeithen, K. B., Reitman, J. S., Rueter, H. H., & Hirtle, S. C. (1981). Knowledge organisation and skill differences in computer programs. *Cognitive Psychology*, 13, 307-325.
- Richman, H. B., Gobet, F., Staszewski, J. J., & Simon, H. A. (1996). Perceptual and memory processes in the acquisition of expert performance: The EPAM model. In K. A. Ericsson (Ed.), *The road to excellence*. Mahwah: Erlbaum.
- Richman, H. B., Staszewski, J., & Simon, H. A. (1995). Simulation of expert memory with EPAM IV. *Psychological Review*, 102, 305-330.
- Schmidt, A. L. (1986). Effects of experience and comprehension on reading time and memory for computer programs. *International Journal of Man-Machine Studies*, 25, 399-409.
- Simon, H. A. & Gobet, F. (2000). Expertise effects in memory recall: A reply to Vicente and Wang. *Psychological Review*, 107, 593-600.
- Simon, H. A., & Gilmarin, K. (1973). A simulation of memory for chess positions. *Cognitive Psychology*, 5, 29-46.
- Vicente, K. J., & Wang, J. H. (1998). An ecological theory of expertise effects in memory recall. *Psychological Review*, 105, 33-57.
- Ye, N., & Salvendy, G. (1994). Quantitative and qualitative differences between experts and novices in chunking computer software knowledge. *International Journal of Human-Computer Interaction*, 6, 105-118.