

The Development of Explicit Rule-Learning

Martin Redington (m.redington@ucl.ac.uk)

Department of Psychology, University College London,
26, Bedford Way, London, WC1E 6BT, UK

Elliot Ronald (elliot.ronald@corpus-christi.oxford.ac.uk)

Department of Experimental Psychology, University of Oxford,
South Parks Road Oxford, OX1 3UD, UK

Abstract

Implicit and explicit learning were originally distinguished in terms of accessibility to verbal report. We identify evidence for the proposal that the implicit/explicit contrast corresponds to a divide between connectionist and symbolic representations. We show that explicit learning shows marked improvement between 4 and 8 years of age. This finding contrasts against very early implicit learning abilities, and concurs with other evidence on the progressive development of symbolic reasoning abilities.

The identification and study of human learning mechanisms is a central concern of psychology and cognitive science. An important contemporary debate in this area concerns the distinction between implicit and explicit learning. Generally, implicit and explicit learning mechanisms have been distinguished in terms of the accessibility of the knowledge acquired to conscious awareness, as assessed by verbal report (e.g., Reber, 1967).

At first this division appears to be relatively convincing: In implicit learning, by definition, participants' verbally reported knowledge is insufficient to account for their performance on some task. For example, participants who memorise strings generated by a finite state grammar are subsequently able to classify strings as obeying or violating the grammar to a significant degree, despite being unable to report the rules of the grammar verbally (e.g., Reber, 1967; Reber & Allen, 1978).

By implication, explicit learning is defined as those cases where participants are able to verbally report sufficient knowledge to account for their performance. For example, if a person can play a legal game of chess, and can also report the rules of chess, then one might conclude that their ability to play a legal chess game was based on this explicit knowledge.

However, the claim that dissociations between verbal report and performance mark the boundary between two distinct learning mechanisms, differing in the accessibility of their knowledge to conscious awareness, has proved problematic. The strongest critics, Shanks and St. John (1994), argue that the insensitivity of verbal report, and the problems of relating participants' reports to the knowledge representations underlying their performance render apparent dissociations between performance and verbal report suspect. Shank and St. John do however endorse the view

that learning mechanisms can be distinguished in terms of the representational form of the knowledge acquired, contrasting exemplar or instance-based learning mechanisms, for which connectionism provides a natural model, against processes of hypothesis testing and rule-discovery, best described by symbolic mechanisms.

In connectionist learning mechanisms, knowledge and cognition are embodied in patterns of activation of many simple units, and in the flow of activations between those units. Learning is the modification of the strengths of the connections between units. In symbolic mechanisms, knowledge is embodied by statements or rules composed of arbitrary symbols, and interpreted according to a consistent syntax and semantics. Learning is the addition of new statements or rules.

Dienes and Perner (1996) suggest that viewing implicit and explicit learning in terms of a divide between connectionist and symbolic mechanisms explains the differing availability of implicit and explicit knowledge to verbal report. The form of knowledge in a connectionist network—the strengths of interunit connections—does not lend itself to verbal communication. In contrast, symbolic knowledge can be easily communicated, and utilised by the receiver.

In this paper we present evidence in support of an implicit/explicit divide based on representational form. We first describe a recent study which dramatically contrasts explicit and implicit learning. We then present evidence on the development of explicit learning, showing marked developmental changes between four and eight years of age. This contrasts against developmental evidence on implicit learning, which appears to function in mature form from the first year of life. We discuss similar findings on the development of symbolic reasoning abilities from other paradigms.

Two dissociable human learning systems

Shanks, Johnstone and Staggs (1997, Experiment 4) report a study where, as in implicit learning studies, they presented participants with a set of rule-governed training strings, and subsequently tested their ability to distinguish between test strings which obeyed or violated the rules underlying the training strings.

However, the materials used by Shanks et al. (1997)

were very different to those of implicit learning studies.¹ Generally in artificial grammar learning studies, materials are drawn from complex grammars, with many (e.g., ten) rules, which specify relationships between adjacent letters. The distinction between grammatical and nongrammatical strings is usually correlated with simple local distributional properties of the training materials, such as the frequency of letter pairs and triples.

Shanks et al. (1997) drew their training and test strings from a crypto-grammar: Grammatical strings possessed the structure 1234.1234, with each number being replaced by the two halves of a pair of letters according to the following three rules: D ↔ F, G ↔ L, K ↔ X. Thus a typical string might be DFGK.FDLX. Nongrammatical strings violated one or more of these rules. Additionally, the training and test items were painstakingly constructed so that only conformance to the rules distinguished grammatical and nongrammatical strings: Local distributional properties provided no useful information.

Shanks et al. (1997) utilised two different training conditions. In the *match* condition, participants were shown a single grammatical training item, and then had to match that example to one of a display of five training items. This is akin to the memorisation training usually used in studies of implicit learning, with participants uninformed that the training stimuli were rule-governed. In the *edit* training condition, participants were informed of the rule-governed nature of the materials. They were shown items which violated the rules, and were required to indicate which elements (letters) were correct and which were not. After each item they were shown the correct string and given feedback as to the actual errors present in the string. This training was intended to facilitate rule-discovery processes.

At test, the match group showed no ability to appropriately classify the test items as obeying or violating the rules. Participants in the edit group fell into two distinct subgroups. One subgroup scored at or around chance on the grammaticality judgment test, while, the other subgroup scored at or near 100% of classifications correct. The manipulation of materials and training conditions appears to flip participants between implicit and explicit learning modes, with the latter resulting in a distinctive bimodal pattern of performance.²

The Shanks et al. (1997) results also provide support for the view that the differences between implicit and explicit knowledge are best explained in terms of the contrast between connectionist and symbolic representations.

¹The materials and training conditions used by Shanks et al. (1997) were based on those of a similar study by Mathews, Buss, Stanley, Blanchard-Fields, Cho and Druhan (1989, Experiment 4). However, the contrast between implicit and explicit learning is much clearer in the Shanks et al. study.

²It is the interaction of materials and training conditions that is important: In an earlier study Shanks et al. (1997, Experiment 3) used the same training conditions with typical artificial grammar learning materials (with many complex rules, governing local dependencies). Both match and edit groups showed typical implicit learning effects (i.e., performance was unimodal and imperfect).

The failure of the match group, who were presented with a typical implicit learning paradigm, to make accurate grammaticality judgments concurs with accounts of implicit learning which stress the learning of local distributional properties (e.g., Perruchet & Pacteau, 1990; Redington & Chater, 1996). In the Shanks et al. (1997) materials, by design, local distributional properties of the materials give no cue to grammaticality. However, in most artificial grammar learning studies, such distributional properties are strong predictors of grammaticality, and sensitivity to such properties is sufficient to account for human performance. Connectionist models, such as the simple recurrent network (Elman, 1990), provide a natural framework for learning of this kind, and are able to capture much of the data on artificial grammar learning (Redington, 1996).

As for the the edit group, three features of their performance contrast clearly against implicit learning, and suggest a process of symbolic, rule-based learning:

1. The step function of participants' performance: Participants either discover the correct rules, or they do not. With "typical" artificial grammar learning materials (e.g., Reber, 1967, or the commonly used set from Reber & Allen, 1978), participants' performance is unimodal, and imperfect: Participants' classification scores exceed chance and untrained controls, but never approach 100% (60–70% correct is a typical score).
2. The ability to capture relationships between "arbitrary" (non-adjacent) elements is consonant with a rule-based representation. Evidence from both the Shanks et al. (1997, Experiment 4) study and St. John & Shanks (1997) suggests that implicit learning is limited to local dependencies (between adjacent or near-adjacent letters).
3. A hitherto unmentioned manipulation in the Shanks et al. (1997) crypto-grammar study was that test items were either similar specific training items (two letters different to), or dissimilar (at least four letters different from) to all of the training items. The edit group were equally likely to classify both kinds of items as grammatical. The absence of effects of surface similarity is often proposed as an indicator of rule-learning, and contrasts against studies such as Vokey and Brooks (1992), where under implicit conditions, participants showed clear effects for the similarity of training and test items (which can generally be explained in terms of similarity in terms of distributional properties).

The Shanks et al. (1997) effects appear to be robust: We replicated the crypto-grammar study, using only the edit training condition. Using the exact same procedure and stimuli, six of our participants ($n = 12$) clearly showed evidence of rule-learning (near-ceiling performance), while the remainder scored at or around chance. Verbal reports and a post-task questionnaire provide convergent evidence that this study captures the implicit/explicit distinction: Participants who showed near-ceiling classification were

able to report the rules of the crypto-grammar without error, whereas those who scored near chance were unable to report the rules.

Alternative hypotheses

Although the Shanks et al. (1997) effects do point towards the operation of two distinct learning mechanisms, and two different forms of representation, the possibility that a single learning mechanism (and representational form) underlies performance on both explicit and implicit tasks remains.

To sketch one possible alternative hypothesis, match and edit training might encourage the consideration of different hypotheses, or different orderings of hypotheses, by a single, symbolic, learning mechanism. Edit training might encourage the initial consideration of nonlocal relationships, while memorisation training might limit hypotheses to local dependencies. With a small number of rules, participants might be able to discover and report them all (as for the learners in the edit condition). When the number of rules is large (as in the relatively complex artificial grammars), participants may well fail to discover every rule, and the sheer number of rules might preclude accurate reporting of every rule that has been learnt. These additional assumptions about the effect of training conditions and the relationship between the complexity of the knowledge base and accessibility to verbal report allow the evidence to be reconciled with a single symbolic learning mechanism.

In general, the problem of distinguishing between different kinds of representational form is very difficult (e.g., see Barsalou, 1990). Additional assumptions will always permit apparent dissociations to be reconciled with a single learning mechanism or representational form. However, it may be possible (and possibly necessary) to support the case for distinct learning mechanisms by appealing to multiple lines of converging evidence. For example, if implicit and explicit learning mechanisms show different profiles of development, this would reinforce the case that there are two distinct mechanisms. Below we present some evidence from a developmental study of the explicit learning effects found by Shanks et al. (1997).

Developmental evidence

As far as the development of implicit learning is concerned, preliminary evidence suggests that implicit learning mechanisms are functioning to a considerable degree within the first year of life. For instance, Saffran, Newport, and Aslin (1996) found that 8-month-old infants exposed to a stream of phonemes of an artificial language subsequently exhibited sensitivity to the sequential structure of the sequence in a preferential listening paradigm. Gomez and Gerken (1997) have observed artificial grammar learning effects in 11- to 13-month old infants, again using a preferential listening paradigm. In both of these studies, performance was unimodal and far from ceiling. Although cross-sectional studies remain to be performed, the indication is that im-

PLICIT learning mechanism functions in essentially the adult form from very early on.

The question we investigated here was the developmental profile of the explicit rule-learning effects found by Shanks et al. (1997). We assessed the performance of four, six, and eight-year-olds on a variation of the edit training condition. We predicted that if explicit learning reflected the the action of a separate learning mechanism, then we would observe a substantial increase in the proportion of learners (participants showing near-perfect classification performance) with increasing age. As well as the classification task, we also administered the Test for Reception of Grammar (TROG test, Bishop, 1981), in order to provide some measure of the cognitive development of each child.

The participants ($n = 36$) were all students at a Hertfordshire primary school. There were 12 participants for each age group, with mean ages of 61, 82, and 102 months.

The materials were based on those used by Shanks et al. (1997, Experiment 4). We simplified the material in order to reduce the effect of differences in cognitive ability or memory capacity due to age. The materials were expressed in terms of animals and fruit in order to provide an engaging task for four- to eight-year-olds.

Grammatical strings all followed the pattern 123.123, with numbers being replaced by elements of the following rules: elephant → orange, lion → apple, rabbit → banana. Note that unlike the Shanks et al. materials, the rules here are unidirectional: Grammatical strings always follow the pattern "animal animal animal fruit fruit fruit" as in the sample test item shown in Figure 1. The divider between the animals and fruits in Figure 1 was present in both training and test stimuli, in order to make the division more salient for participants.

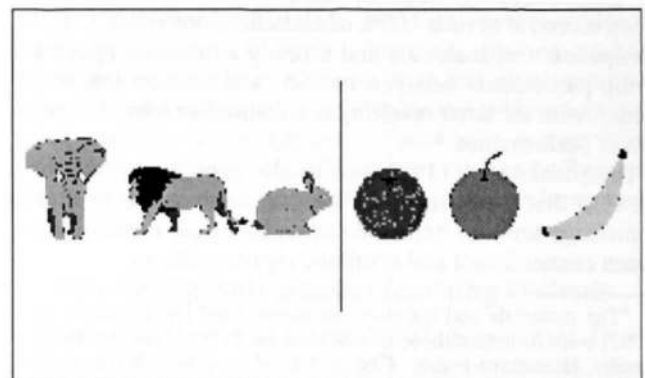


Figure 1: A sample test item (27.5% of actual size).

We constructed two sets of materials (shown in Table 1). Each test set consisted of 12 grammatical and 12 non-grammatical strings. The grammatical and nongrammatical strings had the same combinations of animals, but the combinations of fruit in the nongrammatical strings violated one, two, or three of the above rules. The distributional properties (letter pairs and triples) of the grammatical and nongrammatical items were roughly balanced, and all letter pairs or triples in the nongrammatical strings were also present in the grammatical strings, and vice versa. The non-grammatical items from each test set served as the training items for participants tested on the opposite set. Training and test items were presented to participants on individual sheets of paper.

Table 1: The two sets of test items. A, B, and C refer to elephant, lion and rabbit, and 1, 2, and 3 refer to orange, apple and banana respectively. Underlining indicates a violation of the rules. The training items for each test set were the nongrammatical items from the opposite set.

Set 1		Set 2	
AAB.112	<u>AAB.113</u>	AAC.113	<u>AAC.112</u>
ACA.131	<u>ACA.133</u>	ABA.121	<u>ABA.122</u>
BAA.211	<u>BAA.213</u>	CAA.311	<u>CAA.312</u>
BBC.223	<u>BBC.113</u>	BBA.221	<u>BBA.331</u>
BAB.212	<u>BAB.332</u>	BCB.232	<u>BCB.112</u>
CBB.322	<u>CBB.132</u>	ABB.122	<u>ABB.321</u>
CCA.331	<u>CCA.223</u>	CCB.332	<u>CCB.113</u>
CBC.323	<u>CBC.231</u>	CAC.313	<u>CAC.231</u>
ACC.133	<u>ACC.321</u>	BCC.233	<u>BCC.312</u>
ABC.123	<u>ABC.323</u>	ACB.132	<u>ACB.232</u>
BCA.231	<u>BCA.321</u>	BAC.213	<u>BAC.231</u>
CAB.312	<u>CAB.123</u>	CBA.321	<u>CBA.132</u>

The procedure of the study was as follows: Participants were tested individually in a quiet room, performing the TROG test first, and then performing the edit training and classification test. Prior to edit training, the experimenter informed the child about the task, using the following words:

I am now going to show you some pictures. Each animal likes only one kind of fruit. I will show you three animals and three fruit. In each of these pictures some of the animals don't get the fruit they like. Can you tell me which animals get the fruits they like and which don't?

If you think the animal gets the fruit they like then tick the box under the animal and the fruit. If you think they don't then cross the box under the animal and the fruit.

The order of the animals and the fruit is important. They must be in the right place to get the fruit they

like. At first you will be guessing. I will put the correct answers in the grey boxes to help you learn the rules.

The two different sets of materials shown in Table 1 were counterbalanced within each age group. The training items were presented in random order. Each training item consisted of a nongrammatical item, with a white and grey box under each animal/fruit, for the child's response and feedback from the experimenter. After the child had indicated which animals and fruits they thought were correct, the experimenter would write the correct sequence of ticks and crosses. This is very similar to the edit task used by Shanks et al. (1997), with the exception that the corrected sequences were not presented with the feedback. Participants also received verbal feedback intended to encourage rule-discovery. If the child correctly identified all of errors, they were told "Well done. You got all of them right." In items containing two errors, the experimenter would draw the child's attention to the correct animal-fruit pair, by asking "What do you think the X likes?" In items containing only one error, the experimenter would draw the child's attention to the two "satisfied" animals and their food, by stating "The X gets what they like and the Y gets what they like. The A is eaten, and so is the B." The set of 12 non-grammatical training items was presented twice, for a total of 24 training trials. After training, participants received the following instructions before proceeding to the test phase.

Now you will see some more pictures of the same animals and fruit. I want you to tell me if each page is right. It is right if each animal gets the fruit they like, and wrong if any of the animals don't get the fruit they like.

The 24 test items were then presented in random order. The experimenter recorded the children's responses, but gave no feedback as to whether they were correct or not.

Results

We first analysed the training data. For each training item, the child's response was correct if they correctly identified all rule violations, and incorrect otherwise. The 24 training trials were divided into four blocks of six items each. A $2 \times 3 \times 4$ (Training Set \times Age Group \times Block) mixed model ANOVA revealed effects for Training Set, $F(1, 30) = 11.27, p = 0.002$, Age Group, $F(2, 30) = 13.04, p = 0.0001$, and Block, $F(3, 90) = 5.12, p = 0.0026$. The Training Set \times Age interaction was also reliable, $F(2, 30) = 5.72, p = 0.0079$. The effect of Training Set, and the reliable interaction suggest that Training Set 2 was more likely to encourage rule-learning than Training Set 1, moreso with increasing age, despite the identical nature of the two sets of materials. In fact, it appears that despite random allocation of participants to materials within each age group, participants trained and tested on Set 2 were significantly more advanced, as measured by their TROG scores, than those tested on Set 1,

associative mechanism, and learn an EDS faster than an RS. Adult humans learn by forming "mediated concepts" and learn a RS faster than an EDS.

As for the rule-discovery task reported here, child performance on the discrimination-shift task shows a clear developmental progression, with the proportion of children showing adult-like performance in the majority only after around 6 years of age. Prior to this, most children show discrimination-shift performance characteristic of animals, learning a EDS faster than an RS.

Raijmakers and Molenaar (1996) report that feedforward connectionist networks perform like animals or young children on the discrimination-shift task. In contrast, symbolic representations provide a natural metaphor for mediating concepts. Future studies might assess children's performance in both the rule-discovery and discrimination-shift tasks, in order to see if the apparent correspondence between development in these two tasks is present within particular individuals.

Another interesting line of potential research concerns the performance of amnesic and elderly patients on the explicit rule-discovery task. Implicit learning appears relatively unimpaired in amnesic patients (Knowlton, Ramus & Squire, 1992), despite their known deficits in declarative memory. This tallies neatly with the properties of connectionist and symbolic representations. Connectionist representations degrade gracefully in the face of damage, because their knowledge is distributed over many interunit connections. Symbolic mechanisms are brittle, because their knowledge is concentrated in discrete statements, interpreted by a single processing mechanism. Damage to either rules or processor can have drastic consequences. The finding of impaired explicit rule-learning in amnesic patients would further reinforce the case for representationally distinct implicit and explicit learning mechanisms.

Acknowledgements

Thanks to David Shanks and Theresa Johnstone, for their advice and generous provision of their materials and experimental software. We are grateful to all of the children and teachers for their cooperation during this study. Thanks also to Hilary Leever for advice on appropriate materials for use with children, and to Giles Shilson, Annette Karmiloff-Smith, and three anonymous reviewers for comments on an earlier version of this paper. This research was supported in part by the U.K. Economic and Social Research Council (ESRC) Research Grant number R000236214.

References

- Barsalou, L. W. (1990). On the indistinguishability of exemplar memory and abstraction in category representation. In T. K. Srull & R. S. Wyer, Jr. (Eds.), *Advances in Social Cognition*, Vol. III. Hillsdale, NJ: LEA.
- Bishop, D. V. M. (1983). *Test for Reception of Grammar*. Medical Research Council, U.K.
- Dienes & Perner (1996). Implicit knowledge in people and connectionist networks. In G. Underwood (Ed.), *Implicit Cognition*, (pp. 227–256). Oxford, England: Oxford University Press.
- Elman, J. L. (1990). Finding structure in time. *Cognitive Science*, 14, 179–211.
- Gomez, R. L. & Gerken, L. A. (1997). Artificial grammar learning in one-year-olds: Evidence for generalization to new structure. In E. Hughes, M. Hughes, A. Greenhill (Eds.), *Boston University Conference on Language Development 21*, p. 194. Somerville, MA: Cascadilla Press.
- Kendler, T. S. (1995). *Levels of Cognitive Development*. Mahwah, NJ: Erlbaum.
- Knowlton, B. J., Ramus, S. J. & Squire, L. R. (1992). Intact artificial grammar learning in amnesia: Dissociation of classification learning and explicit memory for specific instances. *Psychological Science*, 3, 172–179.
- Mathews, R. C., Buss, R. R., Stanley, W. B., Blanchard-Fields, F., Cho, J. R. & Druhan, B. (1989). Role of implicit and explicit processes in learning from examples: A synergistic effect. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 15, 1083–1100.
- Perruchet, P., & Pacteau, C. (1990). Synthetic grammar learning: Implicit rule abstraction or explicit fragmentary knowledge? *Journal of Experimental Psychology: General*, 119, 264–275.
- Raijmakers, M. E. J. & Molenaar, P. C. M. (1995). How to decide whether a neural representation is a cognitive concept? *Behavioral and Brain Sciences*, 18, 641–642.
- Raijmakers, M. E. J. & Molenaar, P. C. M. (1995). An experimental test of rule-like network performance. In G. Cottrell (Ed.), *Proceedings of the Eighteenth Annual Conference of the Cognitive Science Society* (p. 827). Mahwah, NJ: Erlbaum.
- Reber, A. S. (1967). Implicit learning of artificial grammars. *Journal of Verbal Learning and Verbal Behaviour*, 5, 855–863.
- Reber, A. S. & Allen, R. (1978). Analogy and abstraction strategies in synthetic grammar learning: A functional interpretation. *Cognition*, 6, 189–221.
- Redington, M. (1996). *What is learnt in artificial grammar learning*. Unpublished doctoral dissertation, Department of Experimental Psychology, University of Oxford.
- Redington, M. & Chater, N. (1996). Transfer in artificial grammar learning: A Reevaluation. *Journal of Experimental Psychology: General*, 125, 123–138.
- St. John, M. F. & Shanks, D. R. (1997). *Implicit Learning from an Information Processing Standpoint*. In D. Berry (Ed.), *How implicit is implicit learning?*, (pp. 162–194). Oxford, England: Oxford University Press.
- Saffran, J. R., Aslin, R. N. & Newport, E. L. (1996). Statistical cues in language acquisition: Word segmentation by infants. In G. W. Cottrell (Ed.), *Proceedings of the Eighteenth Annual Conference of the Cognitive Science Society*, (pp. 376–380). Mahwah, NJ: Lawrence Erlbaum Associates.
- Shanks, D. R., & St. John, M. F. (1994). Characteristics of dissociable human learning systems. *Behavioral and Brain Sciences*, 17, 367–447.
- Shanks, D. R., Johnstone, T. & Staggs, L. (1997). Abstraction processes in artificial grammar learning. *Quarterly Journal of Experimental Psychology*, 50A, 216–252.
- Vokey, J. R., & Brooks, L. R. (1992). Salience of item knowledge in learning artificial grammar. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 20, 328–344.