# UC Merced
## Proceedings of the Annual Meeting of the Cognitive Science Society

**Title**
Dual Processes Mediate Discrimination and Generalization in Humans

**Permalink**
https://escholarship.org/uc/item/67r6k222

**Journal**
Proceedings of the Annual Meeting of the Cognitive Science Society, 28(28)

**ISSN**
1069-7977

**Authors**
Livesey, E.J.
Mansi, C.S.S.A.
McLaren, I.P.L.

**Publication Date**
2006

Peer reviewed

# Dual Processes Mediate Discrimination and Generalization in Humans

**E. J. Livesey (el253@cam.ac.uk)**
**C. S. S. A. Mansi (cm403@cam.ac.uk)**
**I. P. L. McLaren (iplm2@cam.ac.uk)**
Department of Experimental Psychology, University of Cambridge
Downing Street, Cambridge CB2 3EB. UK.

## Abstract

The similarities between learned behavior in animals and humans have often led researchers to conclude that associative processes form the basis of many aspects of human learning. Such an argument has been applied in the past to discrimination learning and the generalization of discriminative responses to new stimuli. However, several experiments that have used two-choice categorization to study post-discrimination generalization in humans have found markedly different results to those found with animals. We argue that this difference occurs mainly because these categorization procedures reflect rule-governed behavior rather than responses based on associative learning. In the current experiment, participants learned to discriminate two very similar training stimuli, differing slightly in hue, and were then tested across a broader range of colored stimuli. Those participants that were able to accurately describe the difference between the training stimuli and also reported using an appropriate rule, produced the same pattern of results as previous categorization experiments of this nature. Those who were not able to identify the difference and did not report using the correct rule produced results that better conform to the predictions of an associative model.

## Dual processes in learning

Within human learning, the distinction has often been made between what could be termed associative learning and rule-governed learning. Dual process models of learning have theoretical and explanatory appeal and although the particulars vary substantially from one model to the next, they generally incorporate one key dichotomy that can be summarized as follows. On the one hand, relatively automatic learning processes result in the formation of associations between representations based on the surface features of the stimuli to which the organism is exposed. The formation, and changes in the strengths of these connections occurs via mechanisms that operate regardless of the higher cognitive processes that the organism has available to it, and can be described by well-specified learning rules (e.g. Blough, 1975; Rescorla and Wagner, 1972). On the other hand, humans also engage in intentional acts of reasoning, deduction and inference, devising verbally-mediated rules to govern their behavior and applying such rules when deemed appropriate. Higher order cognition of this sort is limited by the attentional and cognitive capacities of the individual and influenced by motivational and contextual factors. It might be argued, for instance, that an undergraduate student who is aware that they are participating in a psychology experiment and that their responses are being recorded, will attempt wherever possible to find a rational solution or rule to govern their actions. This, of course, makes it particularly difficult to study associative learning if one cannot be sure that a subjects' responses are not rule-based. It is only in situations where an appropriate rule cannot be learned or confidently applied that we might expect to see a pattern of results that reflects more primitive learning processes.

**Stimulus discrimination and the generalization gradient**
Within animal conditioning, discrimination learning has been studied for almost a century and the form of the post-discrimination generalization gradient is well documented. In such experiments, animals are typically trained to discriminate two similar stimuli by reinforcing responses to one (S+) but not the other (S-). When tested on a wider range of stimuli that differ along the same dimension as S+ and S-, responding typically peaks either for a stimulus very similar to S+ but slightly further removed from S- (i.e. a peak-shifted gradient), or for S+ itself but with the mean number of responses either side of S+ shifted away from S- (i.e. a mean-shifted gradient). Importantly, the peak of the response gradient occurs at, or close to, S+ and then declines as one moves to more extreme values along the dimension. This general form of the post-discrimination gradient has been reliably found in several species and using many different stimulus dimensions (for a recent review, see Ghirlanda and Enquist, 2003). Such results are accurately predicted by simple models of associative learning which represent stimuli as graded activation over a series of input units based on the dimensional qualities of the stimulus (e.g. Blough, 1975; Ghirlanda and Enquist, 1998).

In contrast, results from human experiments seem to be highly contingent upon the type of task employed and quite often produce a strikingly different pattern of results. The specific example that is relevant to the current study is two-alternative forced-choice categorization, where subjects are trained to make one response to one training stimulus, and a different response to a second training stimulus, and are then required to respond to each of a wider range of test values. An associative model designed to simulate dimensional discrimination in animal conditioning can easily be adapted for this type of task. As depicted in Figure 1, such models use the graded activation of a series of

sensory feature units, each "tuned" to a given value along a particular stimulus dimension. Connections to two category units are then modified via a simple error correcting learning algorithm such as the delta rule. Such a model predicts highest accuracy either for the training stimuli (S) or for similar values that are slightly further removed from S. The exact location of the peak depends on the amount of overlap between the activations of the sensory units for each training stimulus. But in any case, a decline towards chance responding, or 50% accuracy, would be expected towards the extremes of the test range, provided the test range is sufficiently broad.
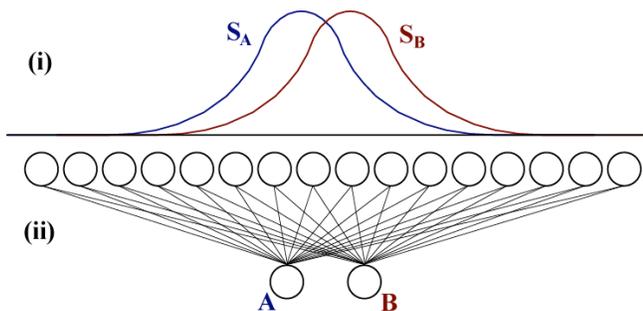


Figure 1. An associative learning model, depicting (i) hypothetical patterns of activation for two training stimuli ($S_A$ and $S_B$) associated with categories A and B, and (ii) a delta-rule network with a series of sensory feature units linked to two category units.

However, most human categorization experiments using simple stimulus dimensions have produced monotonic gradients where some differential responding is evident for the training stimuli but accuracy continues to increase the further one moves along the dimension. Response accuracy peaks at the most extreme values of the test range rather than stimulus values close to S. Such gradients have been obtained on numerous stimulus dimensions including stimulus location (La Berge, 1961), auditory frequency (Cross & Lane, 1962), lifted weight (Capehart & Pease, 1968), light intensity or brightness (Hebert, 1970), and line orientation (Thomas, Lusky & Morrison, 1992).

The difference between gradients produced in human categorization and animal conditioning does not appear to be merely a consequence of having two responses rather than one, as similar experiments with animals have produced peak-shifted, rather than monotonic, gradients (e.g. Blough, 1973). This should, of course, come as no surprise – while there is considerable debate over the extent to which animals can learn abstract relationships between stimuli, there is no doubt that normal humans understand the dimensions along which the simple discriminative stimuli in these experiments differ. The monotonic gradient found in human categorization conforms *exactly* to what would be expected if subjects were learning a relational rule (i.e. press

Left if brighter, heavier, higher pitched, etc.; press Right if duller, lighter, lower pitched, etc.) and applying such a rule consistently to the novel test stimuli. The rule applies best to the stimuli that most obviously possess the relative characteristics on which that rule is based, which are usually the extremes of the test range.

The abstraction of such a relational rule presumably requires the individual to accurately perceive the difference between the initial training stimuli, to have the capacity to describe the abstract relationship between them, and to be able to derive and apply a response strategy to novel stimuli. Where this is indeed possible, it immediately becomes unclear what, if anything, has been learned via more primitive learning processes – as any associatively based generalization and discrimination may be obscured by the fact that participants' responses reflect the application of a cognitive strategy. However, it could still be claimed that in situations where rule-based learning is restricted, one should observe a pattern of behavior based on elementary associative processes. Evidence in support of this view comes from two-choice discrimination experiments where attempts have been made to prevent participants from learning and applying an appropriate response rule. Some studies have used complicated artificial dimensions where the relationship between successive stimuli is not easily verbalized (e.g. Wills and Mackintosh, 1998). Typically, these experiments have obtained peak shifted gradients, where peak accuracy occurs at intermediate stimuli relatively close to S. Others have used simple stimulus dimensions, and have attempted to limit the opportunity for the subject to identify or verbally characterize the relationship between the discriminative stimuli and consequently prevent the learning of an appropriate rule. Jones and McLaren (1999), for instance, reported two experiments where participants were trained on two intermediate levels of brightness of a green stimulus, and then tested across six stimuli covering a much wider range of intensities. Under normal conditions, participants produced monotonic gradients as would be expected from a rule-based analysis, and indeed reported using the correct rule in a post-experiment questionnaire. However, when the number of training trials was halved, or when the contingency between the training stimuli and the correct response was reduced, participants were generally unable to identify the appropriate rule and produced peak shifted gradients that conformed more to the predictions of an associative learning model.

The current study aimed to provide further evidence for dual processes in discrimination learning, using a metathetic stimulus dimension based on hue. The dimension was designed to co-vary with the wavelength of light, a stimulus dimension extensively used in animal experiments on post-discrimination generalization experiments, and which has, with human subjects, produced results similar to animal conditioning studies, albeit under very different task

requirements (Doll and Thomas, 1967). Participants first learned to categorize two very similar shades of green, one slightly more yellow, the other slightly more blue. Participants in the "Easy" group were given an initial discrimination with two colors that, although very similar, differed sufficiently that the difference between them could be identified and readily verbalized. Participants in the "Hard" group were given an initial discrimination based on more similar colors such that it was possible to learn to discriminate between them, but the difference was difficult to characterize and thus an appropriate response rule was difficult to establish. They were then given a range of 15 test stimuli, shown in Figure 2, including the training stimuli (S), and were required to respond to each individually in the absence of any feedback. Finally, all participants were given a structured post-experiment questionnaire in which they were asked to report the relationship between the training stimuli, and the strategies they used to make their judgments.
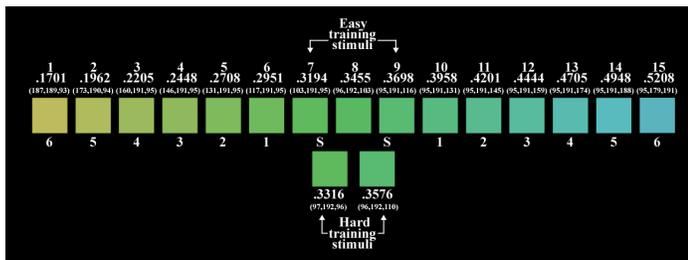


Figure 2. Stimuli along the hue based test dimension. Numbers correspond to ordinal stimulus value, Hue parameter, the Red, Green, Blue color components, and ordinal position in terms of distance to nearest training stimulus (S). Note that the subjective properties of the stimuli differ according to viewing conditions and are presented here for illustrative purposes only.

In this situation, associative and rule-based analyses predict markedly different generalization gradients across the full dimension, as depicted in Figure 3. If participants learn, for instance, that Category A is more yellow and Category B is more blue, and use this relational rule to govern their responses, then one would expect a monotonic gradient with highest accuracy at the extremes of the test range (where the stimuli are most yellow or most blue). Collapsing the dimension around the trained stimuli (S), one would expect accuracy to increase towards ceiling as the distance from the S increases (as shown in Figure 3ii). On the other hand, an associative model such as that illustrated in Figure 1 would predict highest accuracy for stimuli near S, with accuracy falling to chance for stimulus values at the extremes of the test range. It was predicted that the Easy group would be more likely to both identify the relationship between the training stimuli correctly and report using a hue-based strategy or rule to make their judgments. It was also expected that those participants who reported the correct relational rule would show post-discrimination judgments

that matched the predictions of the rule-based analysis, while those that did not report the rule would produce judgments that better matched the predictions of the associative model. In this respect, it was predicted that the relationship between accuracy and distance from S along the dimension would differ according to whether or not rule learning had occurred.
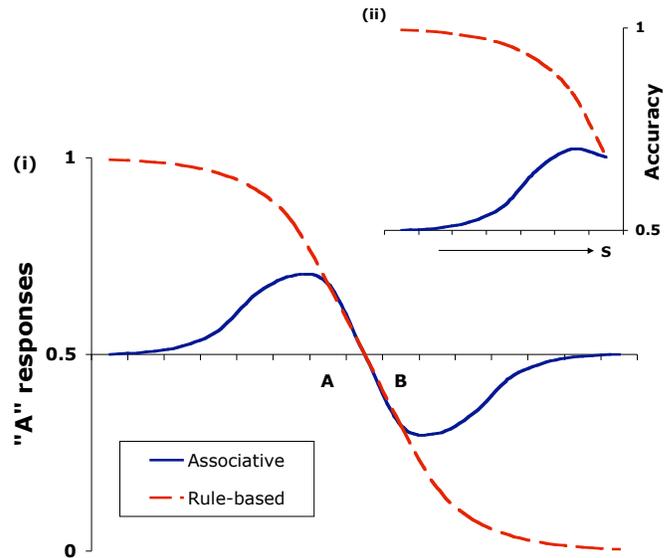


Figure 3. Predictions derived from an associative model and a rule-based analysis: (i) across the full dimension, plotting proportion of category A responses as a function of stimulus value, (ii) collapsed to simplify analysis, plotting accuracy as a function of ordinal distance form nearest S.

## Method

**Subjects and Apparatus**  43 undergraduate students from the University of Cambridge served as participants in the experiment. All reported having normal color vision. They were tested individually in a dimly lit room. The experiment was run on a Macintosh G4 Power PC with 17-inch monitor, using REALbasic software. Participants were assigned to group according to time of arrival.

**Stimuli**  The training and test stimuli were squares of uniform color, measuring 6cm by 6cm and appearing individually in the centre of the computer screen. A set of 15 such stimuli (see Figure 2 above) were created by modifying the Hue parameter in a set of color coordinates based on Hue, Saturation, and Brightness. The Hue values were equally spaced and ranged from .1701 (a greenish yellow) to .5208 (a greenish blue) in approximately equal steps (adjusted slightly to fit the nearest exact Red, Green, and Blue coordinates for displaying on screen). The Hue parameter is designed to keep the Saturation (or "vividness") and Value (or brightness) of the colors

approximately equal, and these separate values were set at 0.5 and 0.75 respectively. Slight adjustments were made to the exact Red, Green, Blue color coordinates of some stimuli following an initial experiment designed to match the subjective luminance of the colors as closely as possible. The colors were therefore considered to be sufficiently isoluminant when presented on screen. Of the 15 stimulus values, the 7th and 9th were used as training stimuli for the Easy condition. For the Hard condition, these training stimuli were replaced with new values between the 7th and 8th stimuli and between the 8th and 9th stimuli, as detailed in Figure 2. For the Hard condition, these adjusted values were used in both the training phase and in replacement of stimulus values 7 and 9 during the test phase.

A concurrent task, intermixed with the hue-based trials, was used to prevent direct comparison of the hue stimuli over successive trials. This task involved the categorization of complex patterns of abstract shapes. Details of the construction and design of these filler stimuli were reported previously by Livesey, Broadhurst and McLaren (2005).

**Procedure** Participants were first given written and verbal instructions explaining that they were required to complete two concurrent, but unrelated categorization tasks, one involving plain colored squares and the other abstract patterns. They were told that in the first phase they would need to learn, through a process of trial and error, which of two categories each stimulus belonged to and to respond accordingly with two keys, "x" and "." (full stop), on the keyboard. They were also informed that a test phase would follow in which they would be given new stimuli and asked to categorize them according to what they had already learned, but would not be given any more feedback.

The training phase involved presentation of 48 hue-based stimuli (24 of each training stimulus) intermixed with 48 filler trials. Trial order was randomized in blocks containing 3 presentations of each of the hue-based and filler training stimuli (12 in total), with the restriction that hue and filler trials alternated. The correct category and response for each of the two training stimuli was counterbalanced so that for half the participants in each group the more yellow stimulus belonged to category A and for the other half, the more blue stimulus belonged to category A. Participants were given up to four seconds to respond with either key, and on doing so received feedback saying either "correct" or "WRONG" (accompanied by a computer beep). If the trial timed out, then "no response" was displayed and the next trial followed immediately.

In the test phase, participants were given five presentations of each of the 15 test values, and an equal number of filler test trials (detailed in Livesey, Broadhurst and McLaren, 2005). Presentation of stimuli was randomized in blocks containing one each of the hue and filler test stimuli (30 in total), again with the restriction that hue trials and filler trials alternated. Participants were not given any feedback in this phase.

At the completion of the experiment, participants were given a post-experiment questionnaire containing two questions regarding the possibility of rules and abstract relations for the hue stimuli:

"Did you notice a difference between the colored rectangle stimuli that belonged to different categories in the training phase? If so, describe the difference between them."

"What strategies did you use to make your judgments about the colored rectangle stimuli in the training phase and the test phase?"

**Analysis** In order to ensure that all participants did in fact learn something during the discrimination phase, a criterion of 55% accuracy or higher across the second half of training was used. Participants who did not meet this criterion were replaced and their data discarded, such that final analyses were conducted on the first 16 participants to perform above the criterion in each group.

In order to classify the responses obtained from the written questionnaire, the questions were independently analyzed by the first and middle authors. The responses were classified in terms of whether the participant had described the difference between the training stimuli in terms of hue (i.e. with reference to their "yellowness", "blueness" or similar relevant adjectives to describe variation in color) and whether the participant had reported using a rule or response strategy based on hue. Participants were classified as "rule learned" if they were able to accurately verbalize the difference between the discriminative stimuli in terms of at least one of the color components or if they reported making their judgments during the test phase according to color.

To simplify the analysis of the post-discrimination gradient, the full dimension was collapsed around the training stimuli as described above, with values grouped according to their ordinal position along the dimension, and expressed in terms of their distance from the training stimuli (S). Thus positions 6 and 10 are one step removed from S, 5 and 11 two steps removed, etc. Accuracy was measured as the proportion of total responses for that stimulus position that were appropriate for the nearest category. Thus "A" responses were correct for values 1-7, while "B" values were correct for values 9-15. Stimulus value 8, which fell between the training stimuli, was not used in the analysis. Analyses were conducted over these points, as well as on accuracy during the training phase, with both group and classification as "rule learned" or "rule not learned" used as factors. It was predicted that increasing distance from S would be associated with a rise in accuracy for those that were using an appropriate rule-based strategy but would instead be associated with a decline in accuracy for those

that did not identify such a rule. Thus linear trend analyses were used, with predictions of opposite trends in each classification and an interaction between them.

## Results

Of the 16 participants in each group, 11 of the Easy group participants and 5 of the Hard group participants reported either identifying or using the correct rule appropriately and were classified as "rule learned". Easy group participants were significantly more likely to report the correct relationship or rule than Hard group participants ($\chi^2 = 4.5$, p < .05). Of those that did not report noticing the correct relationship between the training stimuli, the most frequently reported difference was one of perceived brightness. Of those participants who did not report a hue-based rule or strategy, several other strategies were reported, including simply guessing, with the most frequent again being based on brightness. Importantly, many of these subjects reported being unsure of these strategies and even abandoning them or using them inconsistently during the test phase. The pattern of results did not seem to be systematically affected by any of these incorrect strategies. For instance, the pattern of results for those who reported a strategy based on stimulus brightness was approximately the same as the gradient for those who reported noticing no difference between the training stimuli or reported guessing during the test phase.

Figure 5 shows mean accuracy over the course of the training phase for the Hard and Easy groups, split according to whether the hue-based rule was identified in the post-experiment questionnaire. It is evident that participants in both groups were able to acquire the discrimination and reached a relatively high level of accuracy regardless of whether they were able to accurately verbalize the difference between the training stimuli or report a hue-based response rule. Analysis of variance on the overall training accuracy, with group (Hard vs Easy) and rule ("rule learned" vs "rule not learned") as between subjects factors revealed that the Easy group performed significantly better overall than participants in the Hard group ($F_{1,28} = 8.992$, p = .006). However there was no significant effect of rule ($F_{1,28} = .642$, p = .43), or an interaction between group and rule ($F_{1,28} < .001$, p = .993). Examining each group separately, there was no significant effect of rule in either the Hard group ($F_{1,14} = .287$, p = .601) or Easy group ($F_{1,14} = .366$, p = .555). These results suggest that the Easy discrimination was indeed easier than the Hard discrimination, but clearly both groups were able to discriminate between the training stimuli by the end of training, and whether or not the stimulus difference was correctly identified appears to have had little effect on overall accuracy.
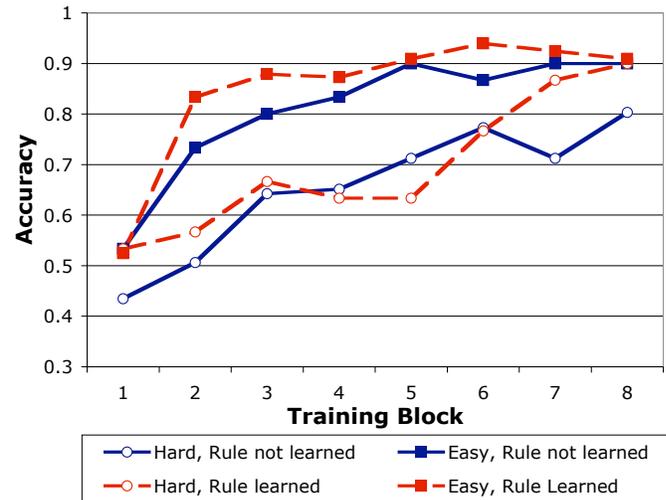


Figure 5. Mean accuracy over blocks of six trials during the training phase.

Figure 6 shows the post-discrimination gradients for groups Hard and Easy, sub-divided according to whether the hue-based rule was identified in the post-experiment questionnaire. The generalization gradients for the "rule learned" participants increase to ceiling accuracy as the distance from S increases. The gradients for the "rule not learned" participants peak in accuracy at or near S and gradually decline as the distance from S increases. ANOVA with stimulus (i.e. ordinal distance from S) as a within subjects factor, and group and rule as between subjects factors revealed a significant main effect of stimulus ($F_{6,168} = 2.624$, p = .019), and rule ($F_{1,28} = 113.256$, p <.001), but no significant effect of group ($F_{1,28} = 1.191$, p = .284). The interaction between stimulus and rule was significant ($F_{6,168} = 14.195$, p < .001), but no interaction with group approached significance (max. $F_{6,168} = 1.241$, p = ns.).
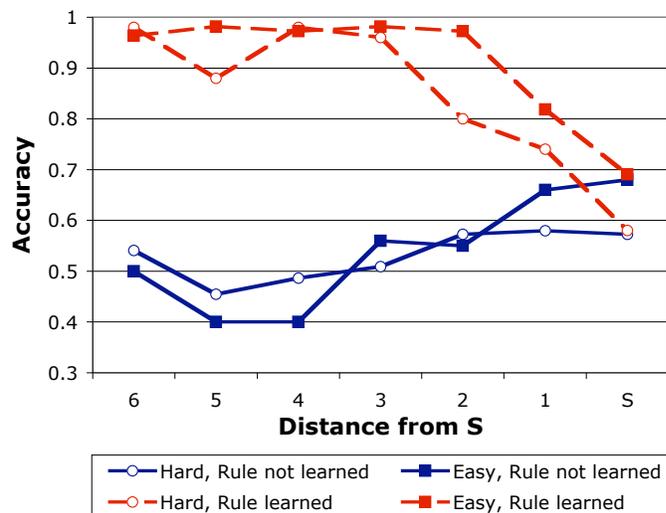


Figure 6. Post-discrimination gradients across the collapsed dimension, with accuracy expressed as a function of ordinal distance from the nearest training stimulus (S).

These results suggest that significant variation in accuracy occurred according to the distance of the test stimuli from S,

and that this pattern differed according to whether the appropriate hue-based rule had been identified. As predicted, a linear trend analysis showed a significant interaction between stimulus and rule ($F_{1,28}$ = 29.924, $p$ < .001). Again there was no significant interaction between stimulus and group or three-way interaction (larger $F_{1,28}$ = 2.528, $p$ = .123). Analyzing the 'rule learned' participants separately, there was no interaction between stimulus and group ($F_{6,84}$ = 1.78, $p$ = .113) but a significant main effect of group ($F_{1,14}$ = 6.646, $p$ = .022) suggested that participants from the Easy group who learned the rule displayed higher overall accuracy than participants from the Hard group who learned the rule. Averaged over all of the 'rule learned' participants, test accuracy was significantly above chance at every point along the collapsed dimension (smallest $t_{15}$ = 5.168, $p$ < .001). Planned linear and quadratic trend analyses both revealed significant main effects (smaller $F_{1,14}$ = 57.05) and no significant interactions with group (larger $F_{1,14}$ = 1.919). These suggest that for the 'rule learned' participants, accuracy generally increases as the distance from S increases and the resulting gradient is negatively accelerated. These results thus closely match the predictions of our rule-based analysis.

Analyzing 'rule not learned' participants separately, there was still no main effect of group ($F_{1,14}$ = .007, $p$ = .936) or interaction between stimulus and group ($F_{6,84}$ = .689, $p$ = .659). Test accuracy was only significantly greater than chance at S ($t_{15}$ = 2.959, $p$ = .01) and positions one step removed from S ($t_{15}$ = 2.58, $p$ = .021). A linear trend analysis revealed a significant main effect ($F_{1,14}$ = 5.154, $p$ = .04) and no significant interactions with group ($F_{1,14}$ = 1.158, $p$ = .300). This suggests that for the 'rule not learned' participants, accuracy generally decreases as the distance from S increases. Although there is no peak shift evident over these test points, this pattern of results still fits well with the predictions of an associative analysis, which would predict highest accuracy at or close to S and a gradual decline in accuracy as the distance from S increases.

## Discussion and Conclusion

There is a clear difference in the generalization gradients produced by those that reported using the correct rule and those that did not. Not surprisingly, the generalization gradients of those participants who reported the rule match the predictions of a verbally mediated cognitive strategy. Learning still occurred in those that were unable to identify the characteristic by which the training stimuli differed, as evidenced by the fact that the participants who did not report the correct rule or the correct stimulus relationship still acquired the discrimination and were still significantly better than chance for the training stimuli and immediately neighboring stimuli during the test phase. However, for the "rule not learned" participants, accuracy declined as the distance from the training stimuli increased, and this pattern of test results fits with an associative analysis. Thus, we

conclude that two learning systems can be involved in discrimination learning in humans, one based on simple associative principles the other reflecting higher order cognition and rule-governed behavior.

## References

Blough, D. S. (1973). Two-way generalization peak shift after two-key training in the pigeon. *Animal Learning & Behavior, 1*, 171-174.

Blough, D. S. (1975). Steady state data and a quantitative model of operant generalization and discrimination. *Journal of Experimental Psychology: Animal Behavior Processes, 1*, 3-21.

Capehart, J., & Pease, V. (1968). An application of adaptation-level theory to transposition responses in a conditional discrimination. *Psychonomic Science, 10*, 147-148.

Cross, D. V., & Lane, H. L. (1962). On the discriminative control of concurrent responses: The relations among response frequency, latency, and topography in auditory generalization. *Journal of Experimental Analysis of Behavior, 5*, 487-496.

Doll, T. J., & Thomas, D. R. (1967). Effects of discrimination training on stimulus generalization for human subjects. *Journal of Experimental Psychology, 75*, 508-512.

Ghirlanda, S., & Enquist, M. (1998). Artificial neural networks as models of stimulus control. *Animal Behaviour, 56*, 1383-1389.

Ghirlanda, S., & Enquist, M. (2003). A century of generalization. *Animal Behaviour, 66*, 15-36.

Hebert, J. A. (1970). Context effects in the generalization of a successive discrimination in human subjects. *Canadian Journal of Psychology, 24*, 271-275.

Jones, F., & McLaren, I. P. L. (1999). Rules and associations. In *Proceedings of the Twenty-First Annual Conference of the Cognitive Science Society*. Mahwah, NJ: Erlbaum.

LaBerge, D. (1961). Generalization gradients in a discrimination situation. *Journal of Experimental Psychology, 62*, 88-94.

Livesey, E. J., Broadhurst, P. J. C., & McLaren, I. P. L. (2005). Discrimination and Generalization in Pattern Categorization. In *Proceedings of the Twenty-Seventh Annual Conference of the Cognitive Science Society*. Mahwah, NJ: Erlbaum.

Thomas, D. R., Lusky, M., & Morrison, S. (1992). A Comparison of Generalization Functions and Frame of Reference Effects in Different Training Paradigms. *Perception & Psychophysics, 51*, 529-540.

Wills, S., & Mackintosh, N. J. (1998). Peak shift on an artificial dimension. *Quarterly Journal of Experimental Psychology Section B- Comparative and Physiological Psychology, 51*, 1-32.