

# Stereotype Formation: Biased by Association

Mike E. Le Pelley  
Cardiff University

Stian J. Reimers  
University College London

Guglielmo Calvini and Russell Spears  
Cardiff University

Tom Beesley  
University College London

Robin A. Murphy  
University of Oxford

We propose that biases in attitude and stereotype formation might arise as a result of learned differences in the extent to which social groups have previously been predictive of behavioral or physical properties. Experiments 1 and 2 demonstrate that differences in the experienced predictiveness of groups with respect to evaluatively neutral information influence the extent to which participants later form attitudes and stereotypes about those groups. In contrast, Experiment 3 shows no influence of predictiveness when using a procedure designed to emphasize the use of higher level reasoning processes, a finding consistent with the idea that the root of the predictiveness bias is not in reasoning. Experiments 4 and 5 demonstrate that the predictiveness bias in formation of group beliefs does not depend on participants making global evaluations of groups. These results are discussed in relation to the associative mechanisms proposed by Mackintosh (1975) to explain similar phenomena in animal conditioning and associative learning.

*Keywords:* stereotypes, attitudes, associative learning, attention, bias

Stereotypes are beliefs that traits or behavioral dispositions are shared by members of a social group. Much research has been conducted in an effort to understand the consequences of, and factors controlling, the expression of existing stereotypes (see Fiske, 1998; Hamilton & Sherman, 1994, for reviews). The mechanisms by which stereotypes initially form on the basis of experience of group-related information, however, have come under less empirical scrutiny (see Hamilton & Sherman, 1994; Hilton & von Hippel, 1996, for reviews). The issue of how stereotypes form is an important one not least because debates have raged about whether stereotypes, which are often linked to prejudice and discrimination, are necessarily bad, or even biased, representations (e.g., Eagly & Steffen, 1984; Fiske, 1998; Hamilton, 1981; Oakes, Haslam, & Turner, 1994). This article describes a novel approach to the study of stereotype formation that may shed some light on these issues (to which we return in due course). Through a series of studies of stereotype formation, we use an approach based on associative learning theory in order to address two of the key

issues that have confronted research on stereotyping: (a) the question of why certain social categories seem to become salient and support stereotypes while others do not and (b) the question of whether (and, if so, why) such stereotypes are systematically biased rather than reflecting the reality of the groups depicted.

## Category Selection in Stereotype Formation

The first of these issues is famously illustrated by an example from Jean-Paul Sartre's (1948) *Anti-Semite and Jew*. He relates the story of an anti-Semitic woman whose prejudice toward Jews stems from prior negative experiences with Jewish furriers. But why (as Sartre noted) has she learned to hate Jews rather than furriers? In more general terms, a target individual can belong to several groups simultaneously (in terms of gender, age, height, hair color, etc.). Stereotype formation can therefore be seen as a categorization problem, in which people learn to associate certain category features (but not others) with behavioral dispositions. What, then, determines the extent to which a given category feature engages the stereotype formation process? Experimental studies of stereotype formation have tended to skirt this issue by providing participants with information relating to target individuals who are each described as belonging to only a single group (e.g., Hamilton & Gifford, 1976; McGarty, Haslam, Turner, & Oakes, 1993; Pryor, 1986; Stroessner, Hamilton, & Mackie, 1992), and hence these studies are unable to address this question. In contrast, certain studies of stereotype *activation* or *expression* have used multiply categorizable targets, for example, Asian women, who might be categorized on the basis of race or gender (Gilbert & Hixon, 1991; Macrae, Bodenhausen, & Milne, 1995; Shih,

---

Mike E. Le Pelley, Guglielmo Calvini, and Russell Spears, School of Psychology, Cardiff University, Cardiff, Wales; Stian J. Reimers and Tom Beesley, Department of Psychology, University College London, London, England; Robin A. Murphy, Department of Psychology, University of Oxford, Oxford, England.

This work was supported by Grant RES062230182 from the United Kingdom Economic and Social Research Council.

Correspondence concerning this article should be addressed to Mike E. Le Pelley, School of Psychology, Cardiff University, Tower Building, Park Place, Cardiff, United Kingdom CF10 3AT. E-mail: lepelleyME@cf.ac.uk

Ambady, Richeson, Fujita, & Gray, 2002; Smith, Fazio, & Cejka, 1996; Zárate & Smith, 1990). Such studies have investigated the factors that determine which of the multiple sets of stereotypes supported by a given target will be activated at a given time, and hence they differ from studies of stereotype formation, which examine how those sets of attitudes form in the first instance. The issue of stereotype formation versus expression is taken up again in the General Discussion.

Associative models first designed to address phenomena of animal conditioning might shed light on this issue. It is well established that experience of the predictive validity of a conditioned stimulus (CS; e.g., a tone) with regard to an unconditioned stimulus (US; e.g., food) can affect the rate at which an animal learns about that CS on subsequent trials (see Le Pelley, 2004, for a review). For instance animals show slower conditioning to a tone CS that has previously been established as nonpredictive of a food US (Mackintosh, 1973). Several formal models of associative learning have been developed to account for such *learned predictiveness* effects (e.g., Kruschke, 2001, 2003; Le Pelley, 2004; Mackintosh, 1975; Pearce & Hall, 1980). In the current article we focus on Mackintosh's (1975) model, which is similar to more recent models developed by Kruschke (2001, 2003). Mackintosh's model includes a stimulus-specific associability factor (sometimes referred to as an attentional factor) that influences the rate of associative learning about a CS. Specifically, a CS maintains a high associability to the extent that it is a better predictor of the US with which it is paired than are other presented CSs. The result is that stimuli that have in the past been relatively accurate predictors of outcomes will be learned about more rapidly than stimuli that have been inaccurate predictors.

Although the Mackintosh (1975) model has its origins in animal conditioning research, recent studies of human learning have also found evidence consistent with this theory (Bonardi, Graham, Hall, & Mitchell, 2005; Griffiths & Le Pelley, 2009; Le Pelley & McLaren, 2003; Livesey, Harris, & Harris, 2009). Moreover, certain phenomena of stereotyping have proved amenable to an associative analysis (Murphy, Schmeer, Mondragon, Vallee-Tourangeau, & Hilton, 2009; Smith & DeCoster, 1998; Van Rooy, Van Overwalle, Vanhooymissen, Labiouse, & French, 2003). On this approach, stereotype formation is modeled as the formation of an association between a mental representation of a group and a representation of a trait or attribute. For example, formation of the stereotype "Members of Group X are lazy" would be modeled as learning of an association between a representation of Group X and a representation of laziness. Once this association is learned, encountering a new member of Group X will tend to activate the idea of laziness; that is, the stereotype will be activated.<sup>1</sup>

Taking this associative approach raises the possibility that the biases anticipated by the Mackintosh (1975) model might also be observed in stereotype formation. That is, biases in the specific features around which stereotypes are formed might arise from learned differences in the associabilities of those features. Specifically, we might be more likely to develop stereotypes regarding a feature (e.g., gender) that has in the past been found to be predictive of behavioral or physical properties, than one (e.g., eye color) that has been less predictive. Although some researchers have proposed that stereotyping is more likely for categories that are more socially meaningful or valued (e.g., Tajfel, 1982), there has been little if any attention to the psychological mechanisms un-

derlying this proposal that would redress the speculative and potentially circular nature of such claims. The current research provides an empirical test of the hypothesis that stereotype formation might be selectively influenced by our previous experience, although in this case not through prior experience of social meaning or value: Rather, in the current experiments participants' previous experience is with information that is independent of the social meaning or value of the later stereotype-relevant information.

Before proceeding, we ought to clarify an issue of terminology. Most of the current experiments measure the influence of prior predictiveness on the rate of development of evaluations (i.e., liking or disliking) of different groups; hence it could be argued that these experiments are more accurately described as testing formation of attitudes or prejudice rather than stereotypes per se—although stereotypes can be evaluative, this need not be the case, as they can also convey purely descriptive meaning, and thus differentiate between groups on descriptive dimensions, hence the analytic distinction with prejudice, or *stereotypic prejudice*. The associative account, however, takes a general-purpose approach: The rules governing formation of associations to valenced information are assumed to be the same as those involved in learning about nonvalenced attributes. In recognition of the idea that this approach places no special importance on the evaluative dimension (an idea that, to anticipate, is supported by the results of Experiments 4 and 5), for simplicity we refer to *stereotypes* throughout this article.

### The Experimental Paradigm

In all experiments, participants were provided with information concerning the behavior of individuals each belonging simultaneously to two groups. This stereotype-formation stage was preceded by an independent training phase designed to foster differences in the learned predictiveness of those groups.

Table 1 shows the design of Experiment 1A, which forms the basis of all of the current experiments. On each trial of Stages 1 and 2, participants read a description of an individual who belonged to two different gangs. Symbols G1–G12 refer to the 12 different gangs to which individuals could belong. During Stage 1, following this description, participants were asked to decide which of two pictures showed the person described. These pictures differed only in the color of the person's clothing. Symbols gr (green) and ye (yellow) in Table 1 refer to the color worn by the "correct" figure for each combination of gangs. For any individual described as a member of Gang G1 or G4 the correct figure was always in green; for members of Gang G2 or G3 the correct figure was always in yellow. Membership of Gangs G5–G8, in contrast, provided no basis for discrimination. Half of the individuals belonging to Gang G5 wore green, while the other half wore yellow;

<sup>1</sup> An alternative approach would involve learning associations from representations of *individuals*, who are each members of Group X, to the representation of laziness. Once these associations are learned, encountering a new member of Group X will tend to activate the representations of these previous individuals (by virtue of this new member's similarity to these previous individuals in terms of membership in Group X), which will activate the idea of laziness. For all present purposes, these two alternatives are functionally equivalent.

Table 1  
Design of Experiment 1A

Stage 1	Stage 2	Test
G1,G5 → gr	G1,G7 → positive	Likeability for G1–G8
G1,G6 → gr	G2,G8 → negative	
G2,G5 → ye	G3,G5 → positive	
G2,G6 → ye	G4,G6 → negative	
G3,G7 → ye	<i>G9,G10 → positive</i>	
G3,G8 → ye	<i>G11,G12 → negative</i>	
G4,G7 → gr		
G4,G8 → gr		

*Note.* Symbols G1–G8 represent different gangs to which target individuals were described as belonging. gr and ye represent different colors of clothing worn by these target individuals (green and yellow, respectively). Positive and negative refer to the affective valence of behavior statements that were attributed to target individuals. Filler trials are shown in italics. On test, participants rated the likeability of each group on a scale from 0 (*strongly dislike*) to 10 (*strongly like*).

the same applied for Gangs G6–G8. As such, Gangs G1–G4 were predictive of clothing color (and hence, according to the Mackintosh [1975] model, would maintain high associability), whereas Gangs G5–G8 were nonpredictive (such that their associability should decline). Therefore we might expect Gangs G1–G4 to maintain a ready ability to engage in new learning, while the corresponding ability of Gangs G5–G8 would decline.

In Stage 2 participants encountered individuals who belonged to pairs of gangs that had not been presented in combination before. Following the description of each individual, participants read a statement describing a behavior performed by that individual. Considering the first four Stage 2 trial types in Table 1, the individuals defined by combinations G1,G7 and G3,G5 performed exclusively positive behaviors; individuals defined by combinations G2,G8 and G4,G6 performed exclusively negative behaviors. Thus gangs paired with green in Stage 1 were equally likely to be paired with positive or negative behaviors in Stage 2; the same applied for gangs paired with yellow in Stage 1. This renders behavior valence statistically independent of clothing color, so neither of the Stage 1 outcomes was predictive of the valence of the statements used on the different trial types of Stage 2. For example, knowing that members of a particular gang wear green tells a perceiver nothing about the valence of behaviors performed by members of that gang.

The behavior statements of Stage 2 were intended to lead participants to form evaluative stereotypes concerning the different gangs (measured in a subsequent test phase using likeability ratings). The question of interest was whether the stereotypes formed with regard to the two different gangs presented on each trial would be of equal strengths, or whether a stronger stereotype would form for one of these gangs than the other. Mackintosh's (1975) theory predicts that predictive gangs should begin Stage 2 with a higher associability than nonpredictive gangs. This would promote learning of evaluative stereotypes regarding predictive gangs relative to nonpredictive gangs during Stage 2. Thus we expected that, following Stage 2, participants would have strong positive stereotypes regarding Gangs G1 and G3, strong negative stereotypes regarding G2 and G4, weak positive stereotypes regarding G5 and G7, and weak negative stereotypes regarding G6 and G8.

The remaining two “filler” trial types in Stage 2 (G9,G10 and G11,G12) involved novel gangs not encountered in Stage 1. It may seem tempting to use these novel gangs as a baseline against which to assess any difference in stereotype strength between predictive and nonpredictive gangs. Any such comparison is ambiguous, however, because Gangs G9–G12 differ from G1–G8 not only in their predictiveness but also in their novelty, and learning can be influenced by novelty independently of predictiveness (see Lubow & Gewirtz, 1995). Hence any difference in stereotype formation between G9–G12 and G1–G4, for example, could reflect a difference in the experienced predictiveness of these cues (with G1–G4 experienced as predictive and hence undergoing a change in associability that would not apply to G9–G12) but could equally reflect the fact that G1–G4 were experienced many times during Stage 1 whereas G9–G12 were not. In contrast, the comparison between predictive and nonpredictive gangs does not admit this confound in terms of novelty, as all of these gangs were experienced an equal number of times during the experiment. Consequently these novel gangs are not discussed further in this article.

### Bias in Stereotype Formation

The second focus of this article is the question of whether (and, if so, why) stereotypes are systematically biased rather than reflecting the reality of the groups depicted. There has been considerable debate in the stereotyping literature between those who have taken the notion of bias in stereotyping as a given (going back to Lippmann, 1922) and others who have argued that stereotyping reflects reality at some level (e.g., Eagly & Steffen, 1986; Oakes et al., 1994; Sherif, 1967). This latter view is particularly associated with self-categorization theorists. For example, Oakes et al. (1994) have argued that “stereotyping is psychologically rational, valid and reasonable, that it provides veridical social perception (i.e., it reflects reality accurately)” (Oakes et al., 1994, p. 187; see also Sherif, 1967, p. 27).

Perhaps the key phenomenon used to support the claim of bias in stereotype formation is the so-called illusory correlation effect (Hamilton & Gifford, 1976), wherein participants develop differing evaluations of two groups as a result of differences in the relative frequency of the two groups, even though they are described by evaluatively equivalent information. Subsequent researchers, however, have noted that appropriate models can explain the illusory correlation effect without any appeal to bias or illusory effects in learning or memory (Fiedler, 1991, 1996; Klauer & Meiser, 2000; McGarty & de la Haye, 1997; Smith, 1991). That is, these researchers have suggested that the stereotype formation mechanism underlying the illusory correlation effect is itself not biased, but that biases can arise from this mechanism when the “environmental input” to the mechanism is itself biased, in terms of unequal frequencies or skewed distributions (see Fiedler, 1996).

With regard to the current experiments, the Mackintosh (1975) model predicts that unequal stereotypes will form to predictive and nonpredictive gangs, despite them being paired with identical evaluative information; for example, Gang G1 is paired with the same behaviors as Gang G7 during Stage 2. Hence this model anticipates that stereotype formation will be biased by differences in the previously experienced predictiveness of groups, even though this predictiveness is established with respect to a

property (clothing color) that is statistically independent of, and hence unrelated to, behavior valence. This pattern of results is quite different from that anticipated on the basis of any *unbiased* model of learning or reasoning (e.g., Allan, 1980; Fiedler, 1991, 1996; Klauer & Meiser, 2000; McGarty & de la Haye, 1997; Schaller, 1994; Smith, 1991). Given that the objective contingency between groups and behavior valence are identical for predictive and nonpredictive gangs, and that prior predictiveness is established with respect to a property that is statistically independent of behavior valence, such theories must predict that stereotypes will be formed equally with regard to both predictive and nonpredictive gangs. As such, we believe that demonstration of this predictiveness effect would represent an unequivocally illusory effect in stereotype formation, that is, a clear example of a bias.

## Experiment 1A

### Method

**Participants, apparatus, and stimuli.** Twenty-seven Cardiff University students (18 women, 9 men) took part in exchange for £5. Participants were tested individually using a standard PC. The 12 gang names were six-letter nonsense syllables all ending in *-s*, for example, Dreebs, Stooks. These names were randomly assigned to Gangs G1–G12 for each participant. The two choice options on each Stage 1 trial were computer-generated pictures of a male figure differing only in color of clothing, one in green and the other in yellow. The sentences describing 18 moderately positive and 18 moderately negative behaviors that were used in Stage 2 were taken from Murphy et al. (2009). Participants in a pilot study had been asked to rate a large corpus of sentences on a scale from 1 (*very positive*) to 7 (*very negative*). The mean judgments for positive and negative sets were  $M = 2.10$ ,  $SE = 0.60$ , and  $M = 5.75$ ,  $SE = 0.85$ , respectively. Examples include “He gave good advice to a friend in trouble” (*positive*) and “He trespassed on private property” (*negative*). The order in which the various positive and negative statements were presented was randomized for each participant.

**Procedure.** On-screen instructions to participants are included in the Appendix. Notably these instructions made no reference to stereotypes: The experiment was described as a study of how people “retain and process visual information.”

On each Stage 1 trial, participants read a person description of the form “[Pair of initials] is a member of the [Gang X], and a member of the [Gang Y].” Initials were generated randomly on each trial, with no two trials using the same pair. Participants were asked to select which of the two pictured figures showed the person described by clicking on that figure. Immediate feedback was provided—the word “Correct” or “Wrong” appeared, and a blue border framed the correct picture. Stage 1 comprised 16 blocks, with each of the eight trial types shown in Table 1 appearing once per block in random order. Presentation order of the two gangs for each trial type was counterbalanced across blocks; for example, for trial type G1,G5 → gr, there were eight presentations with G1 before G5 in the person description and eight presentations with G5 before G1 (the order of these presentations was randomized).

Each Stage 2 trial displayed a person description of the form “[Pair of initials] is a member of the [Gang X], and a member of the [Gang Y],” followed by a statement describing a behavior performed by that person, for example, “He trespassed on private property.” After 8 s a button appeared to allow participants to move to the next trial. This enforced 8-s period was used to ensure that participants attended to the behavior statements that were not accompanied by any form of feedback. Stage 2 comprised six blocks, with each of the six trial types in Table 1 appearing once per block in random order. Presentation order of gangs was counterbalanced as for Stage 1.

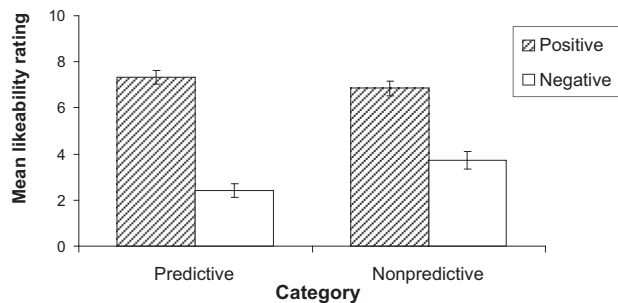
In the final, test stage, participants provided likeability ratings for each of Gangs G1–G8 individually, in random order. On each trial the message “How much do you like members of the [gang]?” appeared above 11 radio buttons, labeled from 0 to 10. A panel at the leftmost end of the scale (below the 0 button) read *strongly dislike* and a panel at the rightmost end (below the 10 button) read *strongly like*. Participants entered their rating by clicking the appropriate button.

### Results and Discussion

Participants were clearly able to learn the Stage 1 gang–color contingencies: Mean percent correct rose steadily across the 16 blocks of Stage 1, reaching 88.9% in the final block. Figure 1A shows mean likeability ratings for the gangs on test. The strength of stereotype formation is indicated by the extent to which gangs paired with positive behaviors elicit higher likeability ratings than gangs paired with negative behaviors. These data were analyzed using analysis of variance (ANOVA) with factors of predictiveness and valence. Significance in all analyses was assessed against a Type I error rate of  $\alpha = .05$ . Crucially, the interaction was significant,  $F(1, 26) = 10.38$ ,  $MSE = 2.10$ , indicating that the extent to which likeability ratings discriminated between positive and negative gang cues depended on the predictive history of those cues. That is, consistent with the central prediction of the Mackintosh (1975) theory, predictive gangs formed significantly stronger stereotypes than did nonpredictive gangs, despite the fact that both classes of cues were paired with equivalent evaluative information. This analysis also revealed a significant main effect of valence,  $F(1, 26) = 73.34$ ,  $MSE = 5.97$ , and no main effect of predictiveness,  $F(1, 26) = 2.17$ ,  $MSE = 2.43$ ,  $p = .13$ .

In the remainder of this article we explore the necessary conditions for this predictiveness bias in stereotype formation. One possibility relates to the nature of the property with respect to which predictiveness is developed. Experiment 1A used a “behavioral” property in Stage 1, clothing color. That is, people make a behavioral choice over the color of clothes that they wear. It seems plausible that participants might view predictiveness established with respect to one aspect of behavior (clothing color) to be a marker of predictiveness with respect to a second aspect of behavior (valence, as studied in Stage 2), and hence the influence of predictiveness might be particularly likely to transfer between the two. In order to test whether this commonality of behavioral outcomes is necessary for the occurrence of predictiveness bias, Experiment 1B used Stage 1 outcomes that were nonbehavioral, namely differences in height.

### A. Experiment 1A



### B. Experiment 1B

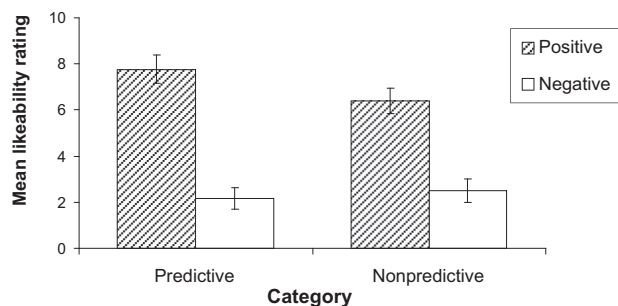


Figure 1. Mean likeability ratings ( $\pm$ SEM) for gangs cues in Experiments 1A (Panel A) and 1B (Panel B), broken down by prior predictiveness (predictive and nonpredictive; this refers to predictiveness during Stage 1) and Stage 2 valence (positive and negative). Data are averaged over gangs from the same prior predictiveness condition that were paired with the same valence of behaviors in Stage 2 (e.g., Gangs G1 and G3 were both predictive gangs paired with positive behaviors in Stage 2).

## Experiment 1B

### Method

**Participants, apparatus, stimuli, and procedure.** Seventeen Cardiff University students (16 women, 1 man) took part in exchange for £5. The picture of the figure in green clothes was replaced with a grayscale picture of a tall person (10.7 cm high onscreen), and the picture of the figure in yellow clothes was replaced with the same picture scaled to show a short person (9.3 cm high onscreen). All other details were as in Experiment 1A.

### Results and Discussion

Results were very similar to those of Experiment 1A. Mean percent correct reached 91.9% in the final block of Stage 1. Figure 1B shows mean likeability ratings for the gangs on test. An ANOVA revealed a significant Predictiveness  $\times$  Valence interaction,  $F(1, 16) = 8.55$ ,  $MSE = 1.45$ , indicating stronger stereotypes for predictive gangs than for nonpredictive gangs. There was a significant main effect of valence,  $F(1, 16) = 26.47$ ,  $MSE = 14.58$ , and the main effect of predictiveness approached significance,  $F(1, 16) = 4.18$ ,  $MSE = 1.02$ ,  $p = .058$ , with a trend toward higher likeability ratings for predictive gangs.

To assess whether the commonality of behavioral outcomes in the pretraining and stereotype formation stages exerted a significant influence on the magnitude of the predictiveness bias, a combined analysis of Experiments 1A and 1B was conducted with experiment as a between-subjects factor (and predictiveness and valence as within-subject factors). The three-way interaction was nonsignificant ( $F < 1$ ), indicating that the influence of prior predictiveness on stereotype formation was equivalent in Experiments 1A and 1B. Hence the biasing influence of prior predictiveness did not depend on the commonality of behavioral outcomes in the pretraining and stereotype formation stages.

Experiments 1A and 1B demonstrate that the experienced predictiveness of groups can influence the formation of stereotypes regarding those groups and that this influence can result in a nonnormative pattern of stereotypes. That is, the bias generated by the influence of prior predictiveness meant that the pattern of stereotypes formed in Stage 2 deviated systematically from the veridical contingencies between cues and behaviors (which were identical for both classes of cues). Consistent with the predictions of the Mackintosh (1975) model, previously predictive groups formed significantly stronger stereotypes than did previously nonpredictive groups. Computational simulations with a version of the Mackintosh model have confirmed that it is indeed able to provide an accurate description of these data. Details of these simulations are available from M. E. Le Pelley on request.

Before proceeding we emphasize the difference between the predictiveness bias observed in Experiment 1 and two established social phenomena—*discounting* and *accentuation*.

**Predictiveness bias versus discounting.** Discounting (Kelley, 1972) refers to a decrease in the evaluation of the strength of Cause X when learned about in the presence of an alternative established Cause Y. Phrased in terms of the current experiment, suppose that during Stage 1 training participants were taught that individuals who belong to Gang G1 perform positive behaviors. Then, in Stage 2, participants encounter individuals belonging to Gangs G1 and G7, who also perform positive behaviors. This might cause perceivers to discount the influence of G7 during Stage 2 training, as there exists an alternative cause (membership in G1) that already explains the observed outcome (positive behavior). Discounting (which is equivalent to the “blocking” effect observed in animal conditioning) can be predicted by unbiased statistical models that use conditional computations of probability and by reasoning-based inference theories (see De Houwer & Beckers, 2002)—exactly the types of theory that we have argued are unable to match the predictions of the Mackintosh (1975) model in Experiment 1.

The fundamental difference between a study of discounting and Experiment 1 is that in our experiment, at the outset of Stage 2 neither of the two cues presented on each trial (e.g., Gangs G1 and G7) is more predictive of behavior valence than the other. The differential predictiveness of the cues during Stage 1 is established with regard to properties that are independent of the valence of Stage 2 behaviors. Consequently, during Stage 2 there is no reason for Gang G1 to lead to discounting of Gang G7 any more than G7 leads to discounting of G1. Hence an account based on discounting cannot explain the systematic difference in the strengths of stereotypes formed by these gangs.

A similar argument distinguishes predictiveness bias from the *expectancy-based illusory correlation effect* (Chapman & Chap-

man, 1967; Hamilton & Rose, 1980). Hamilton and Rose (1980) found that participants' recollections of stimulus information were influenced by their existing knowledge of stereotypic relationships. For example, in a recall test conducted after reading several statements describing individuals, participants overestimated the frequency of statements regarding wealthy doctors, presumably because the concepts "wealthy" and "doctor" benefit from a pre-existing stereotypic connection. In the current experiments, however, the predictive validity of the cues was established with respect to information (clothing color or height) that was independent of the stereotypic content of the Stage 2 information. Hence at the outset of Stage 2 there could not be any difference between the cues in terms of the expected behavior valence with which they would be paired.

**Predictiveness versus accentuation.** Tajfel (1957), in his accentuation theory, has suggested that the discrimination of stimulus objects in a task-relevant dimension, *Z*, is strengthened or accentuated if another, irrelevant dimension, *Y*, also discriminates between the same stimuli (Tajfel, 1957; Tajfel & Wilkes, 1963; see also Eiser & Stroebe, 1972; Fiedler, 1996). That is, accentuation of perceived differences will occur when there is (and relies on there being) a correlation between Dimensions *Z* and *Y*. Accentuation theory implies that the psychological similarity of stimuli belonging to the same category (within-group homogeneity) will increase, while the psychological differences between stimuli belonging to different categories (between-group heterogeneity) will increase. In terms of the current experiment, Stage 1 experience that Gang G1 consistently predicts green clothes while Gang G2 consistently predicts yellow clothes will tend to accentuate the difference between G1 and G2. This might aid participants' discrimination of these gangs during Stage 2, when G1 is paired with positive behaviors and G2 is paired with negative behaviors, which could consequently enhance a difference in likeability of these groups.

However, accentuation theory cannot account for the predictiveness bias observed in Experiment 1, because the dimensions on which participants categorize groups in Stages 1 and 2 are statistically independent (i.e., uncorrelated). On the one hand, accentuation of differences between G1 and G2 (based on clothing color) will aid stereotype formation in Stage 2, because these groups are paired with a different valence of behavior. Accentuation of differences between G3 and G4 will also help stereotype formation, as these groups are also paired with a different valence of behavior. However, there will be equal accentuation of differences between G1 (which predicts green clothes) and G3 (which predicts yellow clothes), and yet both of these groups are paired with positive behaviors in Stage 2; hence this accentuation will tend to hinder formation of distinct evaluations of these groups. Likewise there will be accentuation of differences between G2 and G4, but both of these groups are paired with negative behaviors in Stage 2, so once again this accentuation will hinder stereotype formation. Furthermore, because G1 and G4 predict the same clothing color during Stage 1, accentuation theory anticipates that these groups will come to be seen as more similar; this will again hinder stereotype formation during Stage 2, when these two groups are paired with different outcomes. Similarly, G2 and G3 also predict the same clothing color but different behavior valences. Consequently, accentuation of differences between predictive groups on the basis of the clothing color that they predict (and minimization of differences when predictive groups predict the same color) will produce no net benefit in stereotype formation during Stage 2 relative to

nonpredictive groups; if anything, the number of sources of hindrance to stereotyping of predictive groups arising from accentuation processes outweighs the number of sources that will help. As a result, this account cannot explain the advantage in stereotype formation for predictive gangs over nonpredictive gangs observed in Experiment 1.

Instead, the advantage in stereotype formation for predictive gangs seems to stem not from the specific values on a particular dimension that they predict during Stage 1 but rather from the fact that they predict any value at all; two gangs can predict the same value, or different values, but crucially both are predictive in each case. That is, it seems that predictiveness in general increases the extent to which a stimulus is subsequently learned about, as suggested by associability-based models of associative learning.

In fact, while accentuation theory is unable to provide a valid account of the present findings, it is possible that the class of predictiveness mechanism that is suggested by our data may contribute to the accentuation effect. That is, learning that two stimuli, *A* and *B*, predict different outcomes will, according to Mackintosh's (1975) model, increase the processing resources that are devoted to these two stimuli. It does not seem unreasonable to suggest that this may increase the subsequent discriminability of these stimuli in another task. Indeed, Sherman et al. (2009) have recently advocated an interpretation of the accentuation effect (and, incidentally, the illusory correlation effect) in terms of Kruschke's (2001, 2003) attentional theory, which is formally very similar to Mackintosh's model.

## Experiment 2

Experiment 1 demonstrates that learned predictiveness can bias the extent to which stereotypes are formed regarding different exemplars of the same categorization dimension (gangs). Experiment 2 extends this idea by looking at a situation in which target individuals can be categorized on the basis of two different dimensions (the gang they belong to and the suburb they live in) to see if predictiveness exerts a similar bias when the predictive and nonpredictive cues belong to different dimensions. This leads on to the possibility that there might exist preexperimental biases that influence learning about these different cue dimensions and that such preexperimental biases might interact with the effect of experimentally defined predictiveness.

In addition, the stereotype learning phase of Experiment 1 was unrealistic, in that groups were consistently paired with a single valence of behavior. Behavior by real-world groups is more likely to constitute a mixture of positive and negative. In Experiment 2 each group was paired with such a mixture, with one valence in the majority (e.g., 70% positive, 30% negative). This requires participants to make an overall evaluation of the group rather than basing their judgments on any one statement.

Table 2 shows the design of Experiment 2. Target individuals were described as belonging to a particular gang (G1–G4) and coming from a particular suburb (S1–S4). For participants in condition GANG-P (gang predictive), gang cues were predictive of clothing color during Stage 1, while suburbs were nonpredictive. For participants in condition SUBURB-P (suburb predictive), suburbs were predictive, while gangs were nonpredictive.

Stage 2 was similar to Experiment 1, but all cue combinations were paired with a mixture of positive and negative behaviors, either 70% positive or 70% negative. Given the predictiveness bias

Table 2  
Design of Experiment 2

Stage 1		Stage 2	Test
Condition GANG-P	Condition SUBURB-P	Both conditions	Both conditions
S1,G1 → gr	S1,G1 → gr	S1,G3 → 7 pos, 3 neg	S1?
S2,G1 → gr	S2,G1 → ye	S2,G4 → 3 pos, 7 neg	S2?
S1,G2 → ye	S1,G2 → gr	S3,G1 → 7 pos, 3 neg	S3?
S2,G2 → ye	S2,G2 → ye	S4,G2 → 3 pos, 7 neg	S4?
S3,G3 → ye	S3,G3 → ye		G1?
S4,G3 → ye	S4,G3 → gr		G2?
S3,G4 → gr	S3,G4 → ye		G3?
S4,G4 → gr	S4,G4 → gr		G4?

*Note.* Symbols S1–S4 represent different suburbs in which target individuals were described as living; symbols G1–G4 represent different gangs to which target individuals were described as belonging. gr and ye represent different colors of clothing worn by these target individuals (green and yellow, respectively). Pos (positive) and neg (negative) refer to the affective valence of behavior statements that were attributed to target individuals. Participants in conditions GANG-P (gang predictive) and SUBURB-P (suburb predictive) differ only in the training they received during Stage 1: For condition GANG-P, gang cues were predictive, and suburb cues nonpredictive, to the Stage 1 discrimination, whereas for condition SUBURB-P this was reversed. On test, participants rated the likeability of each group on a scale from 0 (*strongly dislike*) to 10 (*strongly like*).

observed in Experiment 1, we expected cues that were predictive in Stage 1 to form stronger evaluative stereotypes than those that were nonpredictive. Thus for participants in condition GANG-P we expected more extreme likeability ratings for gang cues than for suburbs; in condition SUBURB-P we expected the opposite. That is, the pattern of likeability ratings should differ systematically in the two conditions as a result of differences in their Stage 1 training.

**Method**

**Participants, apparatus, stimuli, and procedure.** Fifty-two Cardiff University students (30 women, 22 men) took part in exchange for £5. Participants were randomly and evenly assigned to conditions GANG-P and SUBURB-P. The four gangs were Buzzards, Eagles, Falcons, and Kestrels; the four suburbs were Hammerton, Kinford, Oakeshott, and Redville. All person descriptions were of the form “[Pair of initials] lives in [suburb] and is a member of the [gang].”

As indicated by Table 2, some of the Stage 2 trial types were paired with a majority of positive behaviors (seven positive, three negative) and others with a majority of negative behaviors (three positive, seven negative). The first two behavior statements encountered for each trial type described behaviors of the majority valence for that trial type. For example, the first two statements for S1,G3 individuals were positive, with the latter eight statements containing five positive and three negative statements in random order. Other details were as for Experiment 1.

**Results and Discussion**

Figure 2 shows the percentage of correct responses during Stage 1; learning appears more rapid in condition GANG-P. An ANOVA

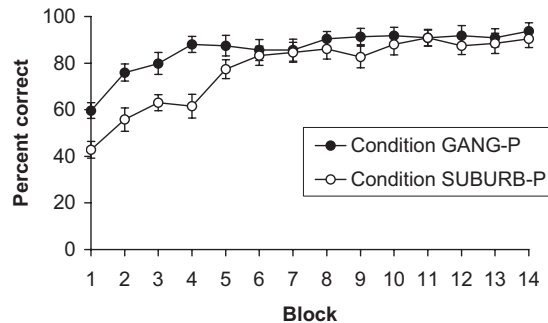


Figure 2. Percent correct responses ( $\pm$ SEM) for the two conditions of Experiment 2 over the 14 blocks of Stage 1. Results are averaged across all eight trial types in each block. GANG-P = gang predictive; SUBURB-P = suburb predictive.

with factors of block and condition revealed significant main effects of block,  $F(13, 650) = 32.16$ ,  $MSE = 225.21$ , and condition,  $F(1, 50) = 4.40$ ,  $MSE = 3,073.90$ , confirming an advantage for condition GANG-P. The interaction was also significant,  $F(13, 650) = 3.93$ ,  $MSE = 225.21$ . Nevertheless, by the end of Stage 1 performance was similar in both conditions—a  $t$  test using data from the final block revealed that the two conditions did not differ reliably ( $t < 1$ ). It is possible that a performance ceiling is masking differences in learning between the two conditions at the end of Stage 1 training. However, a between-condition difference in the extent of Stage 1 learning could not produce the selective effects on likeability ratings that are observed in this experiment—poor learning of Stage 1 relationships will weaken any effects observed on test; it cannot manufacture them.

Figure 3 shows mean likeability ratings for the cues on test, which were analyzed using an ANOVA with factors of category (gang vs. suburb), valence (mainly positive vs. mainly negative), and condition (GANG-P vs. SUBURB-P). Crucially, the three-way interaction was significant,  $F(1, 50) = 4.94$ ,  $MSE = 3.39$ , indicating that the relative extremity of likeability ratings for gang and suburb cues differed reliably in the two conditions. In other words,

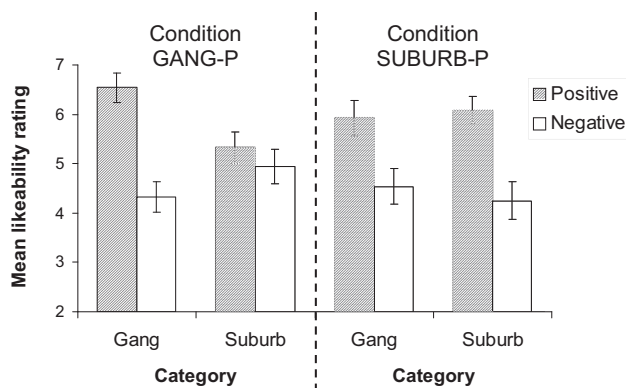


Figure 3. Mean likeability ratings ( $\pm$ SEM) for the test cues of Experiment 2 for conditions GANG-P (gang predictive) and SUBURB-P (suburb predictive). Results are averaged across cues from the same category (gangs or suburbs) that were paired with the same behavior valence in Stage 2.

the different predictive histories of the two categories of cues experienced by the different conditions during Stage 1 exerted selective influences on the ability of those cues to form stereotypes. Figure 3 indicates that, collapsed across both conditions, stereotypes appeared stronger for gangs than for suburbs (mean difference between positive and negative cues was 1.80 for gangs and 1.10 for suburbs). This tendency, assessed by the Category  $\times$  Valence interaction, failed to reach significance,  $F(1, 50) = 1.84$ ,  $MSE = 3.39$ ,  $p = .18$ . The main effect of valence,  $F(1, 50) = 35.73$ ,  $MSE = 3.07$ , was significant; all other effects were non-significant ( $F_s < 1$ ).

The significant three-way interaction legitimizes the calculation of separate two-way effects (Category  $\times$  Valence simple interaction effects) for each condition, with error terms adjusted accordingly. Looking first at the data for condition GANG-P, gang cues paired with positive and negative behaviors elicited more extreme likeability ratings than did suburb cues. Consistent with this suggestion, a two-way ANOVA revealed a significant Category  $\times$  Valence interaction,  $F(1, 25) = 6.88$ ,  $MSE = 3.16$ . The pattern of results from condition SUBURB-P is quite different. Rather than gang stereotypes being considerably stronger than suburb stereotypes, condition SUBURB-P showed marginally more extreme likeability ratings for suburbs than for gangs. For this condition, the Category  $\times$  Valence interaction was nonsignificant ( $F < 1$ ).

Overall, these results show a predictiveness effect similar to that of Experiment 1. Both conditions demonstrated the formation of stronger stereotypes to previously predictive than to previously nonpredictive cues, although this interaction was nonsignificant in condition SUBURB-P. Nevertheless, the significant three-way interaction reveals that the predictiveness of the cue categories experienced during Stage 1 had a reliable influence on stereotype formation regarding those cues during Stage 2, resulting once again in a nonnormative pattern of stereotypes.

It seems likely (and intuitively plausible) that there is a preexperimental bias toward learning about gangs more than about suburbs, as indicated by (a) the significant advantage for condition GANG-P during Stage 1 learning and (b) the trend, although nonsignificant, for gangs to support stronger stereotypes than suburbs when collapsed across both conditions (to anticipate, Experiment 4 provides further evidence consistent with this idea). Thus the results of Experiment 2 seem to represent the combination of two factors: a preexisting bias toward learning about gangs and the influence of Stage 1 predictiveness on associability. In condition GANG-P the preexisting bias would sum with the influence of associability to ensure considerably stronger stereotypes to gangs than suburbs. In condition SUBURB-P, the associability changes (promoting learning about suburbs) would counteract the preexisting bias (promoting learning about gangs), with the result that both categories form stereotypes of similar strength.

The question then arises as to where this preexisting bias toward gangs comes from. One possibility would be to assume that the associability of gang cues was higher at the outset of the experiment than that of suburb cues. Hence the assumption would be that, prior to the experiment, participants had experienced gang membership to be more predictive than area of residence with respect to behavioral and physical properties. Computational simulations using a version of the Mackintosh (1975) model confirm that such an approach can provide an accurate account of the results of Experiment 2. Alternatively, the preexisting bias toward

gangs might have its root in additional, nonassociative processes that also contribute to stereotyping; this issue is taken up again the General Discussion, where various candidates are outlined.

### Experiment 3

We have argued that the predictiveness bias observed in Experiments 1 and 2 conflicts with the predictions of unbiased models of learning and reasoning (e.g., Allan, 1980; Fiedler, 1991, 1996; Klauer & Meiser, 2000; McGarty & de la Haye, 1997; Schaller, 1994; Smith, 1991). Instead these results seem to demand that the stereotype formation mechanism is susceptible to bias, in terms of the potential influence of previously established predictiveness on stereotyping.

This leads on to the issue of how this bias is best characterized. Up to this point we have characterized predictiveness bias in terms of automatic, associative processes, via the concept of associability. It is also possible, however, that predictiveness bias could reflect a biased “higher level” process of statistical reasoning (cf. Allan, 1980; Fiedler, 1991; Schaller, 1994) or rational/Bayesian inference (cf. De Houwer, Vandorpe, & Beckers, 2005; Gopnik & Schulz, 2007). But if this is the case, what would be the source of this higher level bias?

One possibility relates to differences in the entitativity of the groups. *Entitativity* refers to the degree to which a group is perceived as having “the nature of an entity, of having real existence” (Campbell, 1958, p. 17; see also Crawford, Sherman, & Hamilton, 2002). A highly entitative group is one with a high coherence and internal consistency, that is, a group in which all members are alike. Perhaps prior experience that all members of Group X wear the same color, while members of Group Y wear different colors, endows Group X with higher entitativity than Group Y. Consequently participants might reason that all members of Group X are likely to behave in a similar way, while members of Group Y do not, resulting in formation of stronger evaluative stereotypes for Group X than Group Y.

One aspect of the current data is, at least on the surface, problematic for an entitativity account. Such an account sits well with the case in which Stage 1 predictiveness is established with respect to a behavioral property such as clothing color—one might suppose that, if all members of Group X behave in a similar way with respect to the clothing they choose to wear, they may also behave in a similar way with respect to valence. It seems less natural, however, to suppose that similarity in terms of height (over which group members have no choice) will lead participants to view a group as an entity in terms of behavior (over which members do have a choice). A natural interpretation, then, suggests that the influence of differences in entitativity on formation of evaluative stereotypes would be reduced in this latter case, and yet the bias was equally strong in Experiments 1A and 1B. This issue cannot be solved with an entitativity-based account wherein participants believe that people of the same height will tend to behave in a similar way, as height was statistically independent of behavior in Experiment 1B—overall, individuals belonging to “short” groups were equally likely to perform positive and negative behaviors in Stage 2. However, the results are open to an account in which participants are prepared to accept the entitativity of a specific group established with respect to a nonbehavioral property (height) as indicative of the entitativity of that same group with



respect to behavior. The plausibility of such an account remains open to debate.

Perhaps more problematic for the entitativity approach is the demonstration by Le Pelley and McLaren (2003) of a similar predictiveness bias in a food allergy learning task, in which common foods act as cues and types of allergic reaction (e.g., nausea) act as outcomes. These foods, and their similarity to and relationships with one another, are well-known to participants before the experiment. Hence it seems unlikely that learning that, say, oranges are predictive of nausea while lemons are not would lead participants to see oranges as a more entitative category than lemons.

An alternative possibility is that participants' reasoning about stimuli is influenced by the attention paid to those stimuli, with attention in turn influenced by the cues' predictive ability in much the same way as suggested by Mackintosh's (1975) associative model. That is, paying more attention to predictive than nonpredictive gangs as a result of Stage 1 training could feasibly bias a reasoning-based or statistical process. Note that the crucial aspect of this approach is not that it involves *attention* rather than *associability*—indeed, Mackintosh uses the terms interchangeably (see the General Discussion). Rather, the central distinction between the approach discussed here and our associative account is that the former proposes that attention influences *statistical reasoning*, whereas the latter proposes that attention influences the formation of associations.

Such “higher level” accounts of Experiments 1 and 2 still rely on differences in the processing of cues during stereotype formation that depend on differences in their prior predictiveness. As such they do not undermine our general conclusions regarding the biasing influence of learned predictiveness on stereotype formation. In some sense, these accounts are redescription of the processes at work in the associative model of these effects, although without the formalized mechanism of the associative account. Nevertheless, Experiment 3 was designed as an empirical test of the extent to which predictiveness bias relies on higher level cognitive processes.

Higher level accounts of reasoning and inference assume that reasoning is based on knowledge of the frequencies (or probabilities) of co-occurrence of cues and outcomes (e.g., Allan, 1980; Schaller, 1994). In Experiments 1 and 2 this information was provided to participants on a piecemeal, intermixed, trial-by-trial basis. The use of reasoning in such tasks presupposes that participants can integrate the information from each separate trial in an appropriate way and retain that information in order to generate these frequencies internally. Given the complexity of the experimental designs, this will place great demands on memory. In contrast, in Experiment 3 participants were instead provided directly and explicitly with all of the frequency information on which reasoning presumably must be based, in verbal form, at the same time as they were required to provide judgment ratings. This manipulation will drastically reduce the cognitive load involved in reasoning; for example, the memory requirement will be reduced to zero. Consistent with this suggestion, several previous studies have established that providing summary information makes the use of reasoning or inference processes more likely (see Arkes & Harkness, 1983; Le Pelley, Oakeshott, & McLaren, 2005; Shanklee & Mims, 1982; Ward & Jenkins, 1965). Hence if the predictiveness bias ob-

served in Experiments 1 and 2 were a product of higher level inference (perhaps via the influence of entitativity or attention on reasoning), then providing summary information should if anything strengthen the effect. For example, the summary information format would make it easier to deduce the greater entitativity of Group X than Group Y; if entitativity differences are the source of predictiveness bias, this should promote differences in stereotype formation between them.

The alternative possibility, advanced above, is that predictiveness bias has its root in low-level associative processes. Our support for an associative account of this effect should not be taken as support for the view that human learning is *invariably* a result of associative processes. Demonstrations of rule use in human learning (e.g., Shanks & Darby, 1998) clearly indicate a role for inference and analogy-making. Likewise, it is clear that nonassociative processes can contribute to stereotype formation. For example, the stereotypes that we form can be influenced by our “folk theories” about the groups involved (Martin & Parker, 1995), which are not easily captured by an associative learning account. We agree with Shanks (2007) that the available evidence supports a dual-process view, wherein human learning can be a product of both associative processes (along with the biases inherent to such processes) and higher level cognitive processes (presumably based on some form of inferential reasoning and hence more normative), with each tending to dominate under certain circumstances. For example, it has been suggested that these latter effortful, controlled processes implementing reasoning will operate only to the extent that participants have the motivation and opportunity to engage in such reasoning (Sloman, 1996). Consistent with this suggestion is evidence indicating that increasing the cognitive load on participants during a learning task decreases their ability to reason, causing them to fall back on associative mechanisms that are less affected by load (Le Pelley et al., 2005). Similarly, the memory load involved in integrating information from the large numbers of different trial types used in Experiments 1 and 2 might render participants unable or unwilling to use effortful reasoning-based processes to deduce the normative relationships between cues and behaviors, forcing them to rely on more automatic associative mechanisms.

In a recent review concerning the distinction between associationism and higher level cognition, Shanks (2007) argued that “it is particularly compelling that a judgment is based on some non-cognitive [i.e., associative] process if participants, under less stressful conditions, behave differently” (p. 302). In terms of the current paradigm, the suggestion is that associative processes will produce predictiveness bias, while higher level reasoning processes will not. Following Shanks, the summary information provided in Experiment 3 would provide “less stressful conditions” than the trial-by-trial procedure of Experiment 2. That is, Experiment 3 was designed to emphasize higher level reasoning at the expense of associative processes.

If the predictiveness bias found in Experiment 2 were an associative, rather than a reasoning-based, phenomenon, we might therefore expect the results of Experiment 3 to differ from those of Experiment 2. In Experiment 3 the output of reasoning-based processes should outweigh any influence of associative mechanisms, which (assuming that these reasoning-based processes are more normative and hence will not themselves produce a predic-

tiveness bias) will reduce or eliminate the predictiveness bias that these associative mechanisms would otherwise generate. In contrast, if the root of the predictiveness bias observed in Experiments 1 and 2 were itself a product of (biased) reasoning, then if anything we would expect the summary information manipulation of Experiment 3 to strengthen this bias.

Table 3 shows the design of Experiment 3 for condition GANG-P, for which gangs were predictive of Stage 1 outcomes and suburbs were nonpredictive; for condition SUBURB-P, suburbs were predictive and gangs were nonpredictive. This design was as similar as possible to that of Experiment 2, the only difference being the specific combinations of suburbs and gangs experienced in Stage 2. Whilst in Experiment 2 the combinations experienced in Stage 2 were novel, those in Experiment 3 had previously been experienced in Stage 1. This allowed the design to be split in two such that the mental load on participants could be reduced as far as possible. Given that, to anticipate, Experiment 4 of the current article (and much unpublished work from our laboratory) demonstrates predictiveness bias using previously experienced compounds, we can be confident that this difference alone will have little impact on the results.

## Method

**Participants, apparatus, and stimuli.** Fifty-two Cardiff University students (33 women, 19 men) took part on a voluntary basis and were randomly and evenly assigned to conditions. Apparatus and stimuli were as for Experiment 2, and instructions to participants were similar.

**Procedure.** Participants dealt with the information from each half of the design shown in Table 3 separately—whether they dealt with the information from the upper half (relating to

Suburbs S1 and S2 and Gangs G1 and G2) first, or from the lower half (Suburbs S3 and S4 and Gangs G3 and G4) first, was determined randomly. For each half of the design, the participants first received information about the relationships between suburb–gang combinations and clothing color. Thus for participants in condition GANG-P dealing with the upper half of the design, the following information would appear at the top of the screen (with gang/suburb names replacing the placeholders S1, S2, G1, and G2):

You observe 50 people who live in S1 and are members of G1. All of them wear GREEN.

You observe 50 people who live in S2 and are members of G1. All of them wear GREEN.

You observe 50 people who live in S1 and are members of G2. All of them wear YELLOW.

You observe 50 people who live in S2 and are members of G2. All of them wear YELLOW.

The presentation order of these four statements was randomly determined for each participant. Below this, participants were asked to rate each of the two gangs and suburbs involved in the statements according to the clothing color of their members/inhabitants, one at a time, in random order. Ratings were on an 11-point scale, with the leftmost point labeled *much more likely to wear YELLOW*, the rightmost labeled *much more likely to wear GREEN*, and the midpoint labeled *equally likely to wear GREEN or YELLOW*.

Immediately after participants had rated each of the suburbs and gangs for clothing color, more information was provided. If participants were dealing with the top half of the design shown in Table 3, the following statements might appear:

You observe 10 people who live in Suburb S1 and who are members of Gang G1.

7 of them do NICE things (e.g., show their affection to a friend when they really need it).

3 of them do NASTY things (e.g., deliberately give wrong directions to someone who is lost).

You observe 10 people who live in Suburb S2 and who are members of Gang G2.

3 of them do NICE things (e.g., give up their seat to an elderly person on a crowded bus).

7 of them do NASTY things (e.g., do not support their friends when they are bullied).

The order of these statements was randomized, and the examples of behaviors used to illustrate each suburb–gang combination were randomly shuffled. The ratio of seven majority to three minority behaviors was as for Experiment 2. The statements describing the clothing color of the suburb–gang combinations remained visible at the top of the screen throughout. Below the behavior statements, participants were asked to provide likeability

Table 3  
Design of Experiment 3

Stage 1	Stage 2	Test
S1,G1 → gr	S1,G1 → 7 positive, 3 negative	S1?
S2,G1 → gr	S2,G2 → 3 positive, 7 negative	S2?
S1,G2 → ye		G1?
S2,G2 → ye		G2?
S3,G3 → ye	S3,G3 → 7 positive, 3 negative	S3?
S4,G3 → ye	S4,G4 → 3 positive, 7 negative	S4?
S3,G4 → gr		G3?
S4,G4 → gr		G4?

*Note.* Symbols S1–S4 represent different suburbs in which target individuals were described as living; symbols G1–G4 represent different gangs to which target individuals were described as belonging. gr and ye represent different colors of clothing worn by these target individuals (green and yellow, respectively). Positive and negative refer to the affective valence of behavior statements that were attributed to target individuals. The center line shows the division of the design into two parts as used in Experiment 3. Design is shown for condition GANG-P (gang predictive) only; condition SUBURB-P (suburb predictive) differed only in the pattern of outcomes experienced during Stage 1, where suburb cues were predictive of Stage 1 outcomes (S1 and S4 were consistently paired with gr; S2 and S3 were consistently paired with ye) and gang cues were nonpredictive. On test, participants rated the likeability of each group on a scale from 0 (*strongly dislike*) to 10 (*strongly like*).

ratings for each of the two suburbs and gangs mentioned in these statements, in random order, on a scale from 0 to 10. This procedure was then repeated for the cues from the other half of the experimental design.

**Results and Discussion**

The clothing color ratings allow us to verify that participants appreciated the differences in entitativity of the different categorization cues. For condition GANG-P, this would involve more extreme color ratings for gangs than for suburbs (G1 and G4 should receive high scores, G2 and G3 should receive low scores, and all suburb cues should receive intermediate scores). For condition SUBURB-P, ratings should be more extreme for suburbs than for gangs. Two participants, one from each condition, failed to discriminate clearly between predictive and nonpredictive cues; one had a discrimination difference of 0 in their ratings, the other a difference of 0.5 (the next lowest discrimination difference was 5). As a conservative measure the data from these two participants were excluded from further analyses.

Table 4 shows clothing color ratings. High values indicate that group members were perceived as more likely to wear green, low values that they were more likely to wear yellow, and a value of 5 indicates that both colors were equally likely. Results have been averaged for pairs of cues “1 and 4” (Gangs G1 and G4; Suburbs S1 and S4) and “2 and 3” (G2 and G3; S2 and S3). Both conditions showed clear discrimination between cues that were predictive of clothing color and no discrimination between nonpredictive cues. An ANOVA with factors of category (gang vs. suburb), cue (“1 and 4” vs. “2 and 3”), and condition found a significant three-way interaction,  $F(1, 48) = 1,822.57, MSE = 0.60$ , indicating that participants in the different conditions appreciated the differences in predictiveness of the cue categories. Aside from a significant main effect of cue,  $F(1, 48) = 1,802.85, MSE = 0.59$ , and a significant Category  $\times$  Condition interaction,  $F(1, 48) = 4.95, MSE = 0.21$ , all other effects were nonsignificant ( $F_s < 1$ ).

Figure 4 shows mean likeability ratings for the different cues. An ANOVA with factors of category, valence, and condition revealed a nonsignificant three-way interaction,  $F(1, 48) = 1.44, MSE = 0.54, p = .24$ . Thus the relative extremity of likeability

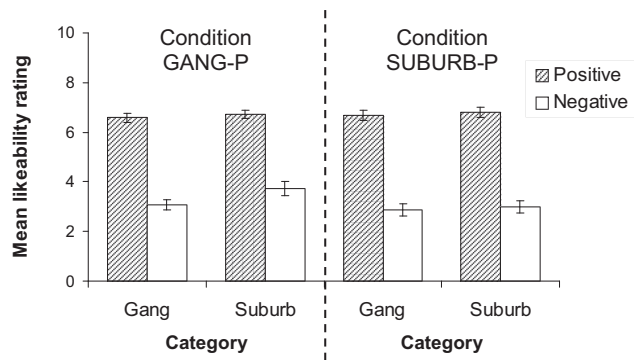


Figure 4. Mean likeability ratings ( $\pm SEM$ ) for the test cues of Experiment 3. Results are averaged across cues from the same category (gangs or suburbs) that were paired with the same behavior valence in Stage 2. GANG-P = gang predictive; SUBURB-P = suburb predictive.

ratings for gang and suburb cues did not differ in the two conditions. The main effect of valence was significant,  $F(1, 48) = 251.56, MSE = 2.48$ , as was that of category,  $F(1, 48) = 8.07, MSE = 0.44$ , with suburb cues receiving slightly higher likeability ratings than gang cues in general. All other effects were nonsignificant,  $F_{max}(1, 48) = 1.23, p = .27$ .

In Experiment 3, prior experience of the predictiveness of one category of cues with respect to clothing color did not influence the extent to which those cues supported stereotype formation. Hence under conditions designed to maximize participants’ use of higher level reasoning processes in interpreting the stereotype-relevant information, no predictiveness bias was observed. The most straightforward interpretation of this dissociation is that the predictiveness bias observed in Experiments 1 and 2 was not a product of reasoning-based mechanisms—if it were, then the effect should, if anything, be strengthened by the manipulation used in Experiment 3. Instead our findings are consistent with a view wherein different learning mechanisms were responsible for the results of Experiments 1 and 2 (which used a trial-by-trial procedure) and Experiment 3 (which used summary information). On this dual-process view, the high cognitive load imposed by the former would cause participants to rely on more automatic, associative processes to interpret the contingency information, along with the biases inherent to such processes. The low cognitive load of Experiment 3 would instead allow participants to use higher level cognitive processes (e.g., reasoning or statistical inference) that are not subject to predictiveness bias (e.g., Allan, 1980; Schaller, 1994; see also Smith, 1991), with these processes outweighing any influence of associative processes (which might otherwise tend to generate predictiveness bias). An alternative but related possibility is that participants in Experiment 3 were subject to the same automatic associative bias as those in earlier experiments but that the low cognitive load allowed them to bring controlled processes to bear, enabling them to recognize and correct for this bias. The current results do not allow us to choose between these alternatives.

A degree of caution is required in accepting the failure to find a predictiveness bias in Experiment 3 as evidence against a higher level account of this effect in Experiments 1 and 2, however. There are several differences between the procedures used in Experiment 2 and Experiment 3 beyond simply the extent to which participants

Table 4  
Mean Clothing Color Ratings (and SEM) for the Cues of Experiment 3

Condition and cue category	Pair of cues	
	1 and 4	2 and 3
GANG-P		
Gang cues	9.56 (0.19)	0.30 (0.13)
Suburb cues	5.06 (0.08)	5.16 (0.13)
SUBURB-P		
Gang cues	5.02 (0.07)	5.06 (0.06)
Suburb cues	9.58 (0.17)	0.28 (0.14)

Note. Pair of cues 1 and 4 refers to the average clothing color rating of gang cues G1 and G4 and of suburb cues S1 and S4. Pair of cues 2 and 3 refers to the average clothing color rating of gang cues G2 and G3 and of suburb cues S2 and S3. GANG-P = gang predictive; SUBURB-P = suburb predictive.

might be expected to use reasoning. For example, Figure 4 reveals a tendency for participants in Experiment 3 to give likeability ratings of 7 to groups paired with a majority of positive behaviors, and 3 to groups paired with a majority of negative behaviors, matching the numbers of positive behaviors provided in the descriptions of the groups—these proportions were not made explicit in Experiment 2. It is perhaps unsurprising that, if the participants of Experiment 3 are to give “valenced” likeability ratings, 7 and 3 are the numbers that they choose to use—the relative numbers of positive and negative behaviors provided by the experimenter are the only relevant information that they have on which to base these ratings. However, that is not to say that we would necessarily expect participants to give ratings of 7 and 3 to both predictive and nonpredictive groups. If one is told that individuals belonging to Group X and Group Y perform seven positive behaviors and three negative behaviors, one could quite feasibly assign all of the “responsibility” for valenced behavior to Group X (leading to a likeability rating of 7 for Group X) and none of that responsibility to Group Y (corresponding to a rating of 5 for Group Y). Thus even though explicit information on behavior is provided in Experiment 3, it remains open to the possibility of predictiveness bias, and hence the absence of such an effect is informative. As noted earlier, from the point of view of reasoning-based accounts the essential functional difference between a trial-by-trial procedure and a summary information procedure is simply that the latter “cuts out the middleman” by providing participants with the frequency information that they would otherwise need to calculate themselves.

The discussion above applies to an approach in which reasoning acts to determine judgments at the point of test, on the basis of unskewed (i.e., normative) frequency information calculated from experience of cue–outcome relationships. In other words, such an approach suggests that the *acquisition* of frequency information is normative (as assumed by Allan, 1980; Fiedler, 1991; Schaller, 1994; Smith, 1991) but that reasoned judgments based on that information might be subject to bias. However, one can imagine an alternative conceptualization of a reasoning-based account, which instead focuses on the role of reasoning in the acquisition of information, and this approach is more successful in accounting for our findings. Thus it is possible that, in Experiments 1 and 2, reasoning about the differential predictiveness of cues during Stage 1 (e.g., in terms of entitativity) leads to a bias in the integration of information during Stage 2, such that the cue–outcome frequencies on which participants base their final, reasoned, judgments are skewed toward previously predictive cues. Explicit provision of unskewed frequency information in Experiment 3 would then overcome any such effect, and hence no predictiveness bias would be expected. Although our current data cannot rule out such an account, it is of course essentially indistinguishable from the associative account outlined earlier, in that it allows predictiveness to bias learning about cues in the same way that associability influences learning in the associative model.

We should stress at this stage that our support for an associative account of predictiveness bias should not be taken as support for the idea that stereotype formation is entirely a product of associative learning processes. As noted earlier, it is clear that nonassociative factors also play a role in stereotyping. For example, instructions regarding the entitativity of groups have been shown to influence the extent to which those groups engage in stereotyping (Crawford et al., 2002), and this influence is not easily cap-

tured by an associative account. Thus it would seem that entitativity does influence stereotyping, but it does not seem to be the source of predictiveness bias. Our aim here is merely to show that some aspects of stereotype formation are explicitly anticipated by, and easily understood from the standpoint of, associative learning theory.

#### Experiment 4

We noted in the introduction that the associative account of the predictiveness bias in stereotype formation observed in Experiments 1 and 2 constitutes a general-purpose approach, wherein the associative rules governing formation of evaluations are the same as those governing formation of associations in any other categorization scenario. In other words, the account is equally applicable to the formation of attitudes or (stereotypic) prejudice (i.e., evaluation) as to formation of nonevaluative stereotypes. It is possible, however, that the evaluative content of attitudes distinguishes them, and the rules determining how they develop, from examples of categorization with nonevaluative stimuli. In the field of conditioning, where associative theory is particularly prevalent, it has been argued that the rules controlling learning with respect to evaluatively valenced outcomes might be different from those that apply to evaluatively neutral outcomes (e.g., Baeyens & De Houwer, 1995; De Houwer, Thomas, & Baeyens, 2001; but see also Davey, 1994; Field & Davey, 1997). This distinction between evaluative and nonevaluative learning has also been applied to biases in stereotype formation, in the context of the illusory correlation effect introduced earlier. Klauer and Meiser (2000) found a standard illusory correlation effect when assessing the attitudes that participants formed with respect to evaluative information but no illusory correlation effect with respect to nonevaluative information (the gender of the group members). They argued that these findings were inconsistent with general purpose accounts of illusory correlation and instead took them as support for the account by evaluative contrast proposed by McGarty and de la Haye (1997). Briefly, this states that the illusory correlation effect derives from participants’ predisposition to look for evaluative differences between groups, and hence (on Klauer and Meiser’s reading) illusory correlation will be observed only on evaluative dimensions. This finding led Berdsen, Spears, van der Pligt, and McGarty (2002) to suggest that “the evaluative dimension seems to be extremely important for the generation of illusory correlation” (p. 101).

We remain cautious of attributing undue weight to a null effect (the failure to find an illusory correlation effect with nonevaluative information), especially given that Klauer and Meiser (2000) did not verify that the salience of the nonevaluative (gender) information that they used was equivalent to that of the evaluative information. Nevertheless, given these previous claims of the importance of evaluation, we wished to test whether the predictiveness bias observed in Experiments 1 and 2 was specific to evaluative information. Although previous studies have demonstrated a predictiveness bias using nonevaluative stimuli (Le Pelley, Suret, & Beesley, 2009), these studies used a more “standard” categorization procedure with arbitrary stimuli. The aim of the current research is to demonstrate that similar learning phenomena can be observed in formation of group beliefs based on experience of individual members of those groups. Hence we need to demon-

strate that predictiveness bias with nonevaluative stimuli will occur when participants have experience of individuals engaging in nonevaluative behaviors and are then required to assess the likely behavior of new members of the group whom they have not previously encountered. This is not typically the case in standard human associative learning and categorization studies, where the trained and tested stimuli are usually the same. Moreover, we wished to demonstrate that a predictiveness bias using nonevaluative behaviors could be demonstrated under similar conditions to those observed with evaluative behaviors in Experiments 1 and 2, necessitating use of a similar procedure.

In order to provide an efficient method of data collection while accessing a more diverse participant population than the university students of Experiments 1–3, Experiment 4 was run over the Internet with no enforced restrictions on participation and no reward for participating. Our first aim (in Experiment 4A) was to replicate the predictiveness bias observed using evaluative information in Experiment 2. Experiment 4B then tested whether an equivalent effect could also be observed with nonevaluative information, while using the same general procedure.

Experiment 4 used an adapted version of the procedure of Experiment 2. To minimize the drop-out rate we used an accuracy threshold for continuation in Stage 1, and during this stage participants experienced all 16 pairings of the four gangs and four suburbs as shown in Table 5, rather than the subset of such pairings used in Experiment 2. To simplify Stage 2, each cue combination was paired with behaviors of a single valence (as in Experiment 1).

### Experiment 4A

#### Method

**Participants.** A total of 189 datasets were collected, with datasets recorded only for participants who completed the experiment. Thirty datasets derived from IP addresses for which data

had already been received and were excluded from analysis. Of the remainder, 103 participants entered their gender as female and 52 as male; four participants made no entry. Ages entered ranged from “under 18” to “61–70,” with the median being “25–30.” Participants were randomly assigned to conditions; of the datasets analyzed, 89 were for condition GANG-P and 70 were for condition SUBURB-P. Although a binomial test shows that this difference does not differ significantly from chance ( $p = .11$ ), the numeric difference in group size may indicate a higher drop-out rate in the harder SUBURB-P condition than the easier GANG-P condition, given the differences in rate of acquisition observed in Experiment 2 (see Figure 2). As noted earlier, such a difference cannot produce selective differences in the likeability ratings for different classes of cues between the two conditions.

**Procedure.** The experiment was run as an Adobe Flash movie contained within an HTML web page, which had links from several sites indexing online psychological experiments. Instructions preceding Stage 1 were similar to those of Experiment 2, with additional description of the criterion-based training in Stage 1, explaining that trials were grouped into blocks of 16 and that participants would need to get at least 13 out of 16 correct in a block to move on to the next stage. Each of the 16 trial types shown in Table 5 (or their equivalent for condition SUBURB-P) was shown once in each block of 16 trials, in random order. If participants had not reached criterion after seven blocks, they were moved on to Stage 2.

Stage 2 was as for Experiment 2, with the exception that the button used to progress to the next trial now appeared after 3 s. Participants experienced each of the four suburb–gang compounds shown in Table 5 10 times; the order of trials within this set of 40 was randomized for each participant. The likeability rating test phase was as for Experiment 2.

#### Results and Discussion

Results were discarded for participants who failed to reach criterion during Stage 1 (13 in each condition). For the remaining participants, learning of the cue–outcome relations in Stage 1 was rapid, reaching criterion after an average of 2.51 blocks in condition GANG-P and 2.95 blocks in condition SUBURB-P,  $t(131) = 1.66$ ,  $p = .10$ . Although nonsignificant, this trend toward faster learning of Stage 1 information in condition GANG-P agrees with that observed in Experiment 2.

Figure 5A shows mean likeability ratings. An ANOVA with factors of category, valence, and condition revealed a significant three-way interaction,  $F(1, 131) = 31.08$ ,  $MSE = 7.17$ , indicating that the different Stage 1 training received in the two conditions selectively influenced the resulting pattern of stereotype formation during Stage 2. The Category  $\times$  Valence interaction was also significant,  $F(1, 131) = 31.71$ ,  $MSE = 7.17$ , indicating that likeability ratings were more extreme for gangs than for suburbs. Aside from a significant main effect of valence,  $F(1, 131) = 308.46$ ,  $MSE = 8.05$ , all other effects were nonsignificant,  $F_{\max}(1, 131) = 2.21$ ,  $p = .14$ .

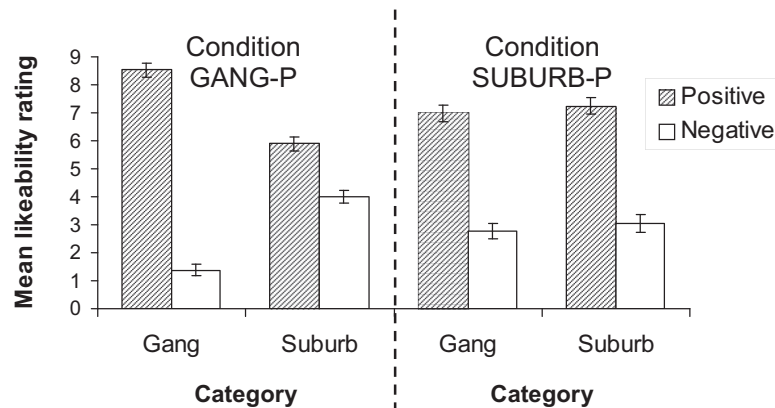
Simple interaction effects revealed that, in condition GANG-P, likeability ratings for gang cues were significantly more extreme than for suburb cues; Category  $\times$  Valence interaction,  $F(1, 75) = 87.27$ ,  $MSE = 6.02$ . In contrast, for condition SUBURB-P, dis-

Table 5  
Design of Experiment 4 (Condition GANG-P Only)

Stage 1		Stage 2	Test
S1,G1 $\rightarrow$ gr	S3,G1 $\rightarrow$ gr	S1,G3 $\rightarrow$ positive/sedan	S1?
S2,G1 $\rightarrow$ gr	S4,G1 $\rightarrow$ gr	S2,G4 $\rightarrow$ negative/hatchback	S2?
S1,G2 $\rightarrow$ ye	S3,G2 $\rightarrow$ ye	S3,G1 $\rightarrow$ positive/sedan	S3?
S2,G2 $\rightarrow$ ye	S4,G2 $\rightarrow$ ye	S4,G2 $\rightarrow$ negative/hatchback	S4?
S1,G3 $\rightarrow$ ye	S3,G3 $\rightarrow$ ye		G1?
S2,G3 $\rightarrow$ ye	S4,G3 $\rightarrow$ ye		G2?
S1,G4 $\rightarrow$ gr	S3,G4 $\rightarrow$ gr		G3?
S2,G4 $\rightarrow$ gr	S4,G4 $\rightarrow$ gr		G4?

*Note.* The first outcome listed for Stage 2 trial types (positive or negative) refers to the procedure of Experiment 4A; the second outcome (sedan or hatchback) refers to the procedure of Experiment 4B. Design is shown for condition GANG-P (gang predictive) only; condition SUBURB-P (suburb predictive) differed only in the pattern of outcomes experienced during Stage 1, where suburb cues were predictive of Stage 1 outcomes (S1 and S4 were consistently paired with gr [green]; S2 and S3 were consistently paired with ye [yellow]) and gang cues were nonpredictive. In Experiment 4A, on test participants rated the likeability of each group on a scale from 0 (*strongly dislike*) to 10 (*strongly like*). In Experiment 4B participants rated which type of car people from each group were likely to drive, on a scale from 0 (*always drive a hatchback*) to 10 (*always drive a sedan*).

## A. Experiment 4A



## B. Experiment 4B

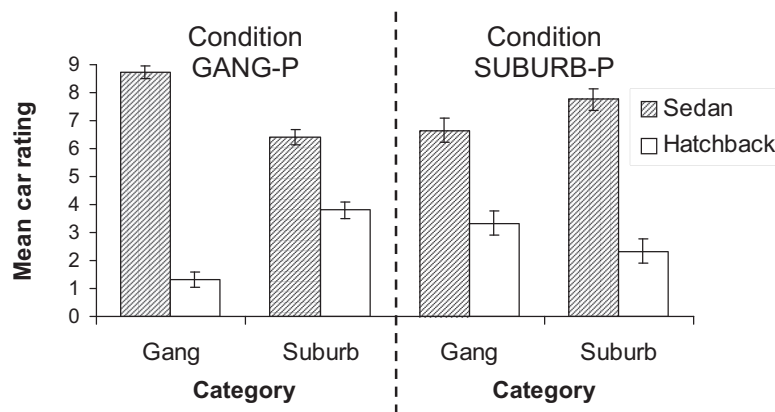


Figure 5. Panel A: Mean likeability ratings ( $\pm$ SEM) for Experiment 4A. Panel B: Mean car ratings ( $\pm$ SEM) for Experiment 4B. In both panels, results are averaged across cues from the same category (gangs or suburbs) that were paired with the same type of outcome in Stage 2. GANG-P = gang predictive; SUBURB-P = suburb predictive.

crimination between suburb cues was similar to that between gang cues; for this condition, the interaction was nonsignificant ( $F < 1$ ).

These results again indicate that, other things being equal, cues previously experienced as predictive of neutral outcomes will form evaluative stereotypes more rapidly than cues experienced as non-predictive. Thus Experiment 4A demonstrated that the predictiveness bias observed in earlier experiments could be replicated using this Web-based procedure.

### Experiment 4B

Experiment 4B investigated whether an equivalent bias would be seen in the learning of nonevaluative stereotypes, an observation that would be consistent with the general-purpose account offered by an associative mechanism. Failure to find such an effect would support suggestions that, in the stereotype formation process, evaluative information engages qualitatively different processes than nonevaluative information (Berndsen et al., 2002; Klauer & Meiser, 2000).

### Method

**Participants.** A total of 159 datasets were collected. Nineteen datasets deriving from IP addresses for which data had already been received were excluded from further analyses. Of the remainder, 104 participants entered their gender as female and 32 as male; four made no entry. Ages ranged from “under 18” to “over 70,” with the median being “25–30.” Eighty-three datasets were for condition GANG-P and 57 were for condition SUBURB-P. This difference was significant (binomial test  $p = .034$ ), presumably reflecting differential drop-out rates in the two conditions as discussed earlier.

**Procedure.** Stage 1 was as for Experiment 4A. In Stage 2, rather than reading valenced behavior statements, participants received nonevaluative information. They were told the type of car driven by target individuals, either a hatchback or a sedan, each illustrated by a grayscale picture. Suburb–gang combinations that were paired with positive behaviors in Experiment 4A were paired with the sedan; combinations paired with negative behaviors in Experiment 4A were paired with the hatchback.

On test, participants rated which type of car people from each suburb and each gang were likely to drive, using an 11-point scale with the hatchback at the leftmost end and the sedan at the rightmost end. Participants were instructed that, if they thought people from a group always drove a particular type of car, they should click at the appropriate end of the scale, whereas if they thought those people drove both cars the same amount, they should click toward the middle of the scale. This yielded car ratings that were analogous to the likeability ratings of earlier experiments.

**Results and Discussion**

Results were discarded for participants who did not reach criterion during Stage 1 (six in condition GANG-P, 17 in condition SUBURB-P). Remaining participants reached criterion after an average of 2.14 blocks in condition GANG-P and 3.10 blocks in condition SUBURB-P, implying significantly faster learning in condition GANG-P,  $t(115) = 3.45$ .

Figure 5B shows mean car ratings on test. High values indicate that a cue was associated strongly with the sedan (equivalent to high likeability in Experiment 4A); low values indicate a strong association with the hatchback (equivalent to low likeability in Experiment 4A). The pattern of results is very similar to that of Experiment 4A. An ANOVA with factors of category (gang vs. suburb), car (sedan vs. hatchback), and condition revealed a highly significant three-way interaction,  $F(1, 115) = 32.29, MSE = 9.74$ , indicating that the different Stage 1 training received by the two conditions selectively altered their patterns of stereotype formation during Stage 2. The Category  $\times$  Car interaction was significant,  $F(1, 115) = 4.88, MSE = 9.74$ , again supporting a general, preexisting bias toward learning about gang cues more rapidly than suburb cues. Aside from a main effect of car,  $F(1, 115) = 252.54, MSE = 9.16$ , no other effects reached significance,  $F_{max}(1, 115) = 1.16, p = .28$ .

Simple interaction effects revealed that, in condition GANG-P, car ratings were more extreme for gangs than for suburbs; Category  $\times$  Car interaction,  $F(1, 37) = 49.74, MSE = 8.91$ . In contrast, in condition SUBURB-P discrimination between suburbs was stronger than between gangs; the Category  $\times$  Car interaction approached significance,  $F(1, 39) = 3.94, MSE = 11.34, p = .054$ .

Thus we can safely conclude that the predictiveness bias observed in Experiments 1, 2, and 4A does not apply exclusively to evaluatively valenced information. Instead this bias seems to be a more general learning phenomenon, supporting the view that the mechanisms of stereotype formation can, to some extent at least, be described by general-purpose associative models.

Although the dependent variable in Experiments 4A and 4B is different (likeability ratings versus car ratings), the similarity between the two motivated a between-experiments comparison using a four-way ANOVA with factors of experiment, condition, category, and outcome (which distinguishes groups paired with positive and negative behaviors in Experiment 4A, and with the sedan and the hatchback in Experiment 4B). There were no significant effects involving the experiment factor,  $F_{max}(1, 246) = 2.93, ns$ , implying that the change from valenced behavior to cars had no significant influence on participants' responding, a finding consistent with the general-purpose approach offered by associative models.

**Experiment 5**

In an effort to establish that predictiveness bias does not apply only to situations in which participants must make global evaluations of groups (as in Experiments 1, 2, and 4A), Experiment 4B demonstrated a similar bias in formation of nonevaluative associations, namely to the type of car that members of a particular group tend to drive. However, it could be argued that even this latter experiment does not meet the criteria for demonstrating predictiveness bias in *stereotype* formation, because the neutral features involved have little social meaning and hence do not bear on theories regarding group differences. One view of stereotypes suggests that they are characterized by clear evaluative connotations (and hence have clear social meaning) but apply in descriptively circumscribed ways. For example, the recently influential stereotype content model (Fiske, Cuddy, Glick, & Xu, 2002) suggests that many group stereotypes vary according to two fundamental dimensions of competence and warmth (see also Abele, Cuddy, Judd, & Yzerbyt, 2008). According to this model, groups can differ stereotypically on these dimensions while sharing the same (positive or negative) valence (e.g., warm vs. competent; see Ford & Stangor, 1992). On this narrower definition, stereotyping (or stereotypic differentiation) occurs when people make inferences regarding a group that are tuned to a specific evaluative trait (e.g., "kind") but do not make other evaluatively similar inferences regarding that same group (e.g., "intelligent"). Hence according to this view, stereotyping would be demonstrated if a group that had been experienced as performing kind, but not intelligent, behaviors were judged as being kind, but not intelligent, and vice versa. Such judgments would reflect tuning to a specific and socially meaningful evaluative trait ("this group is kind, but not intelligent"), rather than a global evaluative judgment ("this group is generally positive, and therefore likely to be both kind and intelligent").

Experiment 5 made use of this distinction in a test of the idea that predictiveness bias applies to stereotyping in this more circumscribed sense, just as it does to formation of prejudice (global evaluations) and nonevaluative associations more generally. The design of Experiment 5 was similar to that of Experiment 1 and is shown in Table 6. Rather than groups being paired with positive or negative behaviors in Stage 2, statements instead described stereo-

Table 6  
*Design of Experiment 5*

Stage 1	Stage 2	Test
G1,G5 $\rightarrow$ gr	G1,G7 $\rightarrow$ kind	Kindness and intelligence ratings for G1–G8
G1,G6 $\rightarrow$ gr	G2,G8 $\rightarrow$ intelligent	
G2,G5 $\rightarrow$ ye	G3,G5 $\rightarrow$ kind	
G2,G6 $\rightarrow$ ye	G4,G6 $\rightarrow$ intelligent	
G3,G7 $\rightarrow$ ye		
G3,G8 $\rightarrow$ ye		
G4,G7 $\rightarrow$ gr		
G4,G8 $\rightarrow$ gr		

*Note.* Symbols G1–G8 represent different gangs to which target individuals were described as belonging. gr and ye represent different colors of clothing worn by these target individuals (green and yellow, respectively). "Kind" and "intelligent" refer to the nature of behavior statements that were attributed to target individuals. On test, participants rated the kindness and likeability of each group on two independent scales from 0 (*not kind/not intelligent*) to 10 (*very kind/very intelligent*).

typically kind or intelligent behaviors, and rather than being asked to provide a single, global evaluation of the groups on test, participants judged the kindness and intelligence of group members separately. To the extent that participants learned to discriminate appropriately between groups paired with kind behaviors and those paired with intelligent behaviors, this experiment therefore tests stereotype formation, even on this latter, rather specific view of stereotyping. Given our claim that predictiveness bias is a feature of a general-purpose learning system, we would once again expect better discrimination (i.e., stronger stereotype formation) for gangs that were predictive of clothing color in Stage 1 than for those that were nonpredictive.

## Method

**Participants, apparatus, and stimuli.** Twenty-five Cardiff University students (19 women, 6 men) took part in exchange for £5. The behavior statements used in Experiment 5 were taken from Fuhrman, Bodenhausen, and Lichtenstein (1989), who provided a list of 400 sentences describing behaviors, scored (and then ranked, on the basis of these scores) for the kindness and intelligence of those behaviors. The 20 sentences used to describe “kind” behaviors in Experiment 5 were selected on the basis of a high rank for kindness ( $M$  rank = 30.2 out of 400), combined with a moderate rank for intelligence ( $M = 141.3$ ). Examples include “He gave his coat to someone when it was cold,” and “He gave his balloon to a child who had let hers go.” The 20 “intelligent” sentences were selected on the basis of a high intelligence rank ( $M = 17.4$ ) combined with a moderate kindness rank ( $M = 147.8$ ). Examples include “He successfully defended himself in a court case” and “He has written five books.” Other stimuli, apparatus, and instructions were as for Experiment 1.

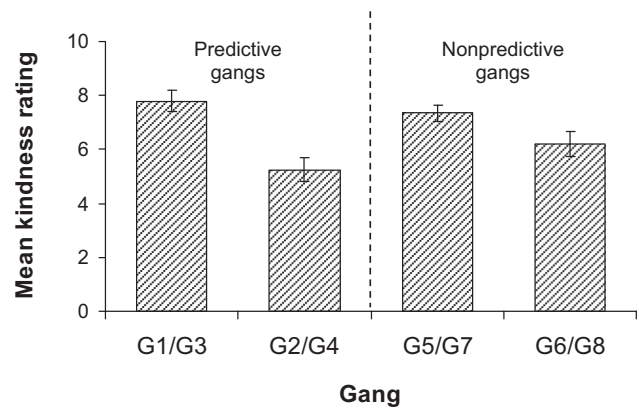
**Procedure.** Stage 1 was as for Experiment 1. Stage 2 was similar to that for Experiment 1, but gangs that were paired with positive behaviors in Experiment 1 were paired with kind behaviors in Experiment 5, and gangs that were paired with negative behaviors in Experiment 1 were paired with intelligent behaviors in Experiment 5. Stage 2 comprised 10 blocks, with each of the four trial types in Table 6 appearing once per block in random order.

On each test trial, participants rated each gang’s kindness and intelligence. For kindness ratings, the question “How KIND are people who are members of the [gang name]?” appeared, above a scale running from 1 (*not kind*) to 10 (*very kind*) on which participants entered their rating. For intelligence ratings, the question “How INTELLIGENT are people who are members of the [gang name]?” appeared, above a scale from 1 (*not intelligent*) to 10 (*very intelligent*). Whether the kindness scale or the intelligence scale appeared at the top of the screen was determined randomly for each participant but remained constant across all test trials. Participants entered their rating for the scale at the top of the screen and then clicked an “OK” button, which caused the second rating scale to appear at the bottom of the screen. Participants rated each of the eight gangs in random order.

## Results and Discussion

Mean percent correct rose steadily across Stage 1, reaching 84.5% in the final block. Figure 6A shows mean kindness ratings for

### A. Kindness ratings



### B. Intelligence ratings

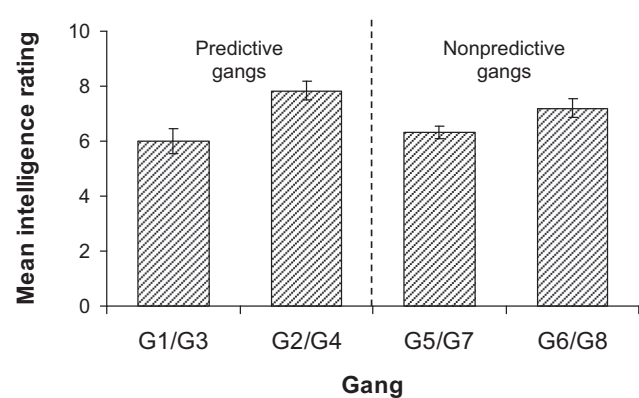


Figure 6. Mean kindness ratings ( $\pm$ SEM) (Panel A) and mean intelligence ratings ( $\pm$ SEM) (Panel B) for Experiment 5. Data were averaged over Gangs G1 and G3 (both predictive, and both paired with kind behaviors; labeled “G1/G3”), Gangs G2 and G4 (predictive, paired with intelligent behaviors; labeled “G2/G4”), Gangs G5 and G7 (nonpredictive, paired with kind behaviors; labeled “G5/G7”), and Gangs G6 and G8 (nonpredictive, paired with intelligent behaviors; labeled “G6/G8”).

the gangs on test, and Figure 6B shows mean intelligence ratings. These data have been averaged over Gangs G1 and G3 (both predictive, and both paired with kind behaviors), Gangs G2 and G4 (predictive, paired with intelligent behaviors), Gangs G5 and G7 (nonpredictive, paired with kind behaviors), and Gangs G6 and G8 (nonpredictive, paired with intelligent behaviors).

If participants had formed strong associations from predictive gangs to the behaviors with which they were paired, then we would expect Gangs G1/G3 to be rated as very kind but not very intelligent (recall that the “kind” sentences were selected on the basis of implying only moderate intelligence) and Gangs G2/G4 to be rated as very intelligent but not very kind. Moreover, if participants had formed weak associations from nonpredictive gangs to the behaviors with which they were paired, then they should be less willing to discriminate between G5/G7 and G6/G8 in terms of the kindness versus intelligence of these gangs.

With regard to kindness ratings, then, we would expect relatively high ratings for G1/G3, relatively low ratings for G2/G4,



and intermediate ratings for G5/G7 and G6/G8; this is the pattern observed in Figure 6A. These data were analyzed using an ANOVA with factors of predictiveness (predictive gangs [G1/G3 and G2/G4] vs. nonpredictive gangs [G5/G7 and G6/G8]) and behavior (gangs paired with kind behaviors in Stage 2 [G1/G3 and G5/G7] vs. gangs paired with intelligent behaviors in Stage 2 [G2/G4 and G6/G8]). Consistent with the anticipated pattern of results, this revealed a significant Predictiveness  $\times$  Behavior interaction,  $F(1, 24) = 9.34$ ,  $MSE = 1.39$ , indicating that kindness ratings discriminated more clearly between predictive gangs than between nonpredictive gangs. The main effect of behavior was significant,  $F(1, 24) = 12.05$ ,  $MSE = 7.03$ , and the main effect of predictiveness was nonsignificant,  $F(1, 24) = 1.06$ ,  $MSE = 1.60$ ,  $p = .31$ .

With regard to intelligence ratings, we would expect relatively low ratings for G1/G3, relatively high ratings for G2/G4, and intermediate ratings for G5/G7 and G6/G8; this is the pattern observed in Figure 6B. An ANOVA as described above again found a significant Predictiveness  $\times$  Behavior interaction,  $F(1, 24) = 4.48$ ,  $MSE = 1.29$ , indicating that intelligence ratings discriminated more clearly between predictive gangs than between nonpredictive gangs. The main effect of behavior was significant,  $F(1, 24) = 6.56$ ,  $MSE = 7.05$ , and the main effect of predictiveness was nonsignificant ( $F < 1$ ).

Experiment 5 used a design in which all behaviors were evaluatively positive but differed in the traits that they described (kindness or intelligence). It is clear that participants' inferences about the characteristics of the different gangs, as revealed by their ratings on test, are tuned to the specific stereotype content of the behaviors with which those gangs have been paired. The main effect of behavior in each of the analyses described above indicates that gangs paired with kind behaviors tended to be judged as more kind than those paired with intelligent behaviors, and vice versa. Hence these ratings cannot simply reflect global evaluations of the gangs (in which case kindness and intelligence would always go together, since both are evaluatively positive traits); instead they reflect specific, and selective, learned knowledge about traits of the gangs, and as such constitute stereotypes. Most importantly, and consistent with the results of the previous experiments, the extent of this selective learning about traits (i.e., the extent of stereotype formation) was influenced by the prior predictiveness of the groups involved, with stronger stereotypes formed for groups previously experienced as predictive of clothing color than those experienced as nonpredictive.

### General Discussion

There exists a long-running debate in the social psychology literature over whether stereotyping necessarily implies the potential for bias in perception and information processing. There is a long tradition that sees stereotyping as biased because it fails to do justice to the uniqueness and variation among individuals who might be tarred with the stereotype of their social category (Fiske, 1998; Hamilton, 1981; Lippmann, 1922). Previous support for this position has been taken from research on stereotype formation deriving from the illusory correlation effect (Hamilton & Gifford, 1976), demonstrating that systematically biased stereotypic beliefs can arise even when groups are described by behaviorally equivalent information. The dominant "cognitive miser" metaphor also

suggests a process of distortion and bias, implied by a view of stereotypes as cognitive shortcuts or energy-saving devices akin to otherwise adaptive heuristics (Fiske & Neuberg, 1990; Macrae, Milne, & Bodenhausen, 1994). However, other researchers have made the argument that the illusory correlation effect may in fact reflect real contingencies in the skewed contingency data (e.g., Fiedler, 1991, 1996; McGarty & de la Haye, 1987; Smith, 1991) and thus reflect more rational, and unbiased, information processing strategies. In line with this possibility, some have argued that stereotyping may reflect the reality of differences between groups (e.g., Eagly & Steffen, 1984, 1986; Oakes, Haslam, & Turner, 1994) and that the perception of stereotypic differences, when the intergroup dimension is salient and relevant, is functional and rational (Oakes et al., 1994). This has left the suggestion that the stereotype formation mechanism can be fundamentally biased at something of a theoretical impasse.

By making recourse to phenomena of learning first established in the field of animal conditioning, we believe that the current research injects new life into the bias debate by showing that the mechanism of stereotype formation can indeed be biased, specifically by a spurious influence of the previously experienced predictiveness of social groups. In Experiments 1, 2, and 4, groups were established as either predictive or nonpredictive of an evaluatively neutral property (e.g., clothing color), before being paired with statements describing valenced behaviors. The prior predictiveness of groups biased the extent to which participants developed evaluative stereotypes, with previously predictive groups supporting stronger stereotypes than previously nonpredictive groups despite both classes of groups being paired with evaluatively identical information. Unlike in the case of the "standard" illusory correlation effect, it is hard to see how an unbiased model could ever account for our findings: The fact that the statistical relationship between groups and behavior valence was identical for previously predictive and nonpredictive groups, coupled with the fact that prior predictiveness was established with respect to a property that was statistically independent of behavior valence, ensures that there is no objective basis in the relevant data for a stereotypic difference between these groups. As such we believe that these experiments represent an unequivocally illusory effect of bias in stereotype formation.

Although predictiveness bias lies beyond the scope of unbiased models of stereotyping, our findings fit well with the predictions of associative models incorporating a variable associability parameter that determines the extent to which a given cue engages the learning process, on the basis of that cue's predictive history. The notion of associability being determined by "predictiveness" allows this approach to account for a biasing influence of prior learning on stereotype formation (by altering the rate of learning about a group) even if the content of that prior learning is not directly related to the content of the stereotype-relevant information learned subsequently (e.g., clothing color and behavior valence are not directly related in our experiments).

Experiment 3 was a first attempt to address the question of the cognitive level at which predictiveness bias occurs. Predictiveness bias lies beyond the scope of unbiased reasoning-based or statistical models of stereotyping, but it could feasibly be explained by adding an additional mechanism to such accounts that acts to bias the reasoning process, for example via the concept of entitativity. Experiment 3, however, demonstrated that a manipulation demon-

strated in previous studies to increase the ease of reasoning (Arkes & Harkness, 1983; Le Pelley et al., 2005; Shanklee & Mims, 1982; Ward & Jenkins, 1965) in fact eradicated predictiveness bias. The most straightforward account of this finding is that the root of this bias is not in reasoning. More generally the dissociation provided by the comparison of Experiments 2 and 3 is consistent with a dual-process view, under which learning can be a product of either associative or higher level cognitive processes, with each tending to dominate under different conditions. Although this conclusion rests on a null result, the demonstration of a robust effect in all other experiments indicates that meaning should be attached to the failure to detect a predictiveness bias in Experiment 3 (in which, according to the higher level accounts advanced above, the effect should, if anything, be stronger).

That said, as noted in the Discussion of Experiment 3, our results cannot rule out a higher level account in which reasoning about predictiveness produces a bias in the *acquisition* of information on which subsequent, reasoned, judgments are based. Given, however, that related effects of predictiveness are observed in animals such as rats (e.g., Oswald et al., 2001) and pigeons (George & Pearce, 1999) to which we would not typically ascribe reasoning abilities, and that certain predictiveness effects in humans are inconsistent with reasoning-based accounts (Le Pelley et al., 2005), parsimony favors the associative account of the current findings. Moreover, this simple and formalized associative account offers a clear and testable model of such predictiveness effects.

The associative model of predictiveness bias constitutes a general-purpose approach; the rules governing learning about the attributes of social groups are the same as those governing formation of associations in any categorization task. Hence this model does not distinguish between formation of attitudes/prejudice and formation of stereotypes. The problem in substantiating this argument is that there exist different views of what constitutes a stereotype. On one view, to which we subscribe, we have stereotyping when our expectations about, and judgment of, an individual are based on information (attributes or traits) that does not necessarily derive from our experience with that individual but derives from knowledge we have about a group to which the individual belongs (cf. Stangor & Lange, 1994). On this account, whether or not that trait or disposition has evaluative content is irrelevant, and consequently many previous researchers have used nonevaluative dimensions to probe stereotype activation (e.g., Bargh, Chen, & Burrows, 1996; Blair & Banaji, 1996; Dijksterhuis, Spears, & Lepinasse, 2001; Macrae, Bodenhausen, Milne, & Calvini, 1999; Stangor, Carr, & Kiang, 1998), because any effects observed under such circumstances clearly cannot be attributed to global evaluations (prejudice). Following this approach, in Experiment 4 we demonstrated a predictiveness bias in stereotype formation when groups were paired with nonevaluative information (types of car driven by group members). This conflicts with prior studies indicating that biases in the formation of evaluations do not apply to learning of nonevaluative information (Klauer & Meiser, 2000), although differences in the learning preparation make the reason for this discrepancy difficult to interpret.

However, on an alternative account it could be argued that even Experiment 4 fails to meet the criteria for demonstrating stereotype formation, to the extent that the nonevaluative car information has little social meaning. We remain skeptical on this issue; it seems clear that cars do have social meaning and bear on causal theories

about group differences (and indeed, that social meaning can be imputed to almost any group differences; Spears, 2002). For example, we might perceive Germans as very efficient, which is why they drive efficient and reliable cars; in contrast Italians may be perceived as flamboyant and hence keen on flamboyant cars. Nevertheless, an alternative view of stereotyping might claim that stereotypes are characterized by clear evaluative content but are specifically “tuned” and distinct from the evaluation itself (e.g., kindness and intelligence are two specific traits that both have positive value). On this view stereotyping is demonstrated when people make inferences regarding a group that are tuned to a specific trait but do not make other evaluatively similar inferences regarding the same group (e.g., judging a group as kind but not intelligent). In order to conclusively establish that predictiveness bias does influence stereotyping, regardless of which definition one chooses to use, Experiment 5 used a design in which groups were paired with statements describing either stereotypically kind or intelligent behaviors and demonstrated a predictiveness bias in participants’ selective judgments of the kindness and intelligence of the groups involved.

In summary, although the literature has defined stereotyping in a variety of ways, sometimes to include an evaluative component (often called stereotypic prejudice when this component is negative), our theoretical account is equally able to explain stereotype formation from the level of minimal nonevaluative bases to more complex forms of this definition that include evaluative components.

### Observational Versus Feedback-Driven Learning

In Experiments 1, 2, 4, and 5, participants’ task during Stage 1 was of a different nature to that in Stage 2. Stage 1 involved predictive learning with explicit feedback, whereas in Stage 2 participants read statements—no response was required, and no feedback was provided. A predictive learning task was used in Stage 1 to maximize the probability of participants learning the cue–outcome contingencies (in the absence of feedback participants would have less incentive to learn about, or even attend to, the onscreen information), which in turn should maximize the associability difference between predictive and nonpredictive cues (although Experiment 3 did not use corrective feedback, the clothing color ratings [see Table 4] allowed us to verify that participants appreciated the difference in predictiveness of the different groups) and hence maximize our chances of detecting a predictiveness bias. It remains for future research to further clarify the necessary and sufficient conditions for this predictiveness bias; for example, will similar effects obtain if Stage 1 learning is purely observational, or must people actively process the Stage 1 information for it to influence subsequent learning?

### Predictiveness and Transfer

The occurrence of a predictiveness bias relies explicitly on experience of the groups’ predictive status during Stage 1 exerting an influence on learning about those groups during Stage 2. In other words, there must be some kind of transfer of information between the two stages. Perhaps the most peculiar aspect of the effect is that this transfer occurs, despite the objective indepen-

dence and subjective dissimilarity of the two phases of the experiments reported here. Why might this be?

According to the Mackintosh (1975) model, each stimulus has a parameter associated with it, which Mackintosh labeled alpha. The alpha value of a particular stimulus is learned on the basis of the experienced predictiveness of that stimulus. That is, the value of alpha is changed incrementally with experience of the relationship between a particular stimulus and other events of significance, such that stimuli that are predictive of other events will tend to develop a higher alpha than nonpredictive stimuli. This alpha value then itself influences the rate of future learning about that stimulus—stimuli with higher alphas will be learned about more rapidly than those with lower alphas. Thus alpha determines how “associable” the stimulus is in future, hence our description of it as the associability of a stimulus. Because alpha is a stimulus-specific parameter, such that a stimulus’s alpha value is “owned by” that stimulus, there is the potential for differences in alpha values learned to stimuli during a given phase of training to persist to a subsequent phase of training involving those same stimuli, and hence to foster differences in the rate of learning about those stimuli (see also Kruschke, 2001, 2003). The question then becomes, what is alpha, and why might it persist?

### Alpha as Attention

Both Mackintosh (1975) and Kruschke (2001, 2003) suggested that alpha represents the attention that is paid to a cue. This approach sees attention as being a learned response to a cue (see also Lubow, 1989) that is based on the experienced predictiveness of that cue. Thus participants learn to attend to particular cues, and to ignore others, on the basis of their predictiveness. It does not seem unreasonable to suggest that learned attentional responses might persist to a new training context. Suppose that participants have learned, during Stage 1, to ignore Gangs G5–G8 (i.e., if their gaze should fall on one of these gangs, they should immediately move on to look for another cue) and to attend to Gangs G1–G4 (i.e., if their gaze should fall on one of these gangs, they should maintain focus on it). Indeed, recent experiments in our laboratory in which we have monitored eye gaze during a category learning task have confirmed exactly this pattern (Le Pelley, Beesley, & Griffiths, 2009). With extended training, as provided in the experiments reported here, this attentional response may become relatively automatic, as demonstrated by Livesey et al. (2009) and Beesley and Le Pelley (2009). To the extent that this is the case, then the attentional response will persist to a subsequent task in which the same stimuli are presented. This account need then assume only that the amount of attention that is paid to a cue influences how much is learned about that cue in order to explain predictiveness bias.

It is important to note that this approach does not assume that there will be transfer of *predictiveness*. That is, it does not state that a stimulus that has been predictive of outcomes in Context X will transfer to be perceived as more predictive of outcomes in Context Y. What it does argue is that there will be transfer of attention, such that a stimulus that has been predictive in Context X will be learned about more rapidly in Context Y. It may be the case that this stimulus is actually nonpredictive in Context Y, in which case this transfer of attention will ensure that the organism is quick to learn about its nonpredictive status (and as a conse-

quence of rapid learning about its nonpredictive status, this approach anticipates that attention will rapidly be moved away from this stimulus and toward a stimulus that is more predictive in Context Y).

### Alpha as Memory

An alternative view exists in which alpha relates to memory processes, rather than attention. Learning of a stereotype involves encoding the relationship between the mental representation of a particular group label and a particular type of behavior. The extent to which this stereotype is learned, then, will depend in part on how well the group label is represented in memory, where a stimulus is represented in memory “well” if it is represented as very distinct from other, similar stimuli—that is, if it allows memories to be clearly addressed. Suppose we experience members of Group A performing positive behaviors and members of Group B performing negative behaviors. If Groups A and B are represented very distinctly in memory, then it will be easy to encode this stereotype information, as each piece of behavioral information will be addressed to (associated with) the correct group label. If, on the other hand, the mental representations of Groups A and B are very similar—such that the two groups are highly confusable—then information regarding members of Group A might be mistakenly addressed to the representation of Group B, and vice versa. Hence stereotype formation will proceed more slowly.

To reiterate, according to this account the rate of stereotype formation regarding a particular group will be determined, in part, by how well that group is represented in memory. And it is possible that prior experience of a group’s predictiveness might influence how well it is represented in memory, in a manner consistent with the general principles of the Mackintosh (1975) model. That is, consistent pairings of a particular group with the same color in Stage 1 might improve the mnemonic encoding of that group, producing a more distinct representation (and hence faster subsequent learning of stereotypes in Stage 2) than for a group that is inconsistently paired with two different outcomes during Stage 1.

Both of the general approaches described above agree on the fundamental idea that prior experience of a group’s predictiveness influences some aspect, alpha, of the processing of that group’s representation and that this difference in processing then subsequently influences the rate of stereotype formation regarding the group. Moreover, both agree that alpha increases for predictive groups and decreases for nonpredictive groups, in line with the Mackintosh (1975) model. The difference is that the former approach identifies alpha with the attention paid to a group, while the latter identifies alpha with the group’s representation in memory. The results of the current experiments do not allow us to decide between these alternatives, because we have measured only learning about the different groups, and the two accounts make the same predictions with regard to rate of learning. However, future work could address this issue by measuring other consequences of differences in predictiveness. For example, we noted above that the attentional view anticipates that predictiveness will influence other properties of a cue that are related to attention, such as the extent to which it commands overt and covert orienting, or its suscepti-

bility to the “attentional blink” (Livesey et al., 2009), and several such predictions have been confirmed in studies of human associative learning using nonsocial stimuli (see Le Pelley, in press, for a review). In a related vein, the memory-based view anticipates that previously predictive groups will show an advantage over previously nonpredictive groups when these group labels are used as targets in tests of memory. For example, a natural interpretation of this account suggests that speeded recognition of predictive groups should be faster than that of nonpredictive groups.

We should also acknowledge, however, that attention and memory might well be entangled: It seems plausible that we would form a better memory representation of those stimuli to which we are attending (see Griffiths & Mitchell, 2008; Mitchell, in press). Similarly, we might pay more attention to those stimuli that are represented in memory as being more distinct from others. Given this potential interaction between attention and memory, it may well be difficult to deconfound the two and hence deduce the locus at which predictiveness exerts its primary effect, even if additional measures of stimulus processing are taken as described above.

### Learned Predictiveness, Primitive Categories, and the Real World

Suppose that a perceiver were to read that a British male with dark hair had committed an antisocial act. Would this information lead them to think more negatively of British people, or of men, or of people with dark hair? Our findings indicate that, other things being equal, stronger stereotypes will form regarding those features that have previously been more predictive of other significant behavioral or physical features. If in the past this perceiver has experienced gender to be more predictive of a person’s attributes (their physical appearance, clothing, mannerisms, etc.) than nationality, then gender will have an advantage over nationality in stereotype formation.

The phrase “other things being equal” is essential here. Our demonstrations of predictiveness bias come from carefully controlled laboratory studies, using novel and fictitious groups, in which we can isolate the influence of predictiveness on learning. This is not so easy in the real world, however, where many factors will interact to determine the extent to which a given cue will engage in stereotype formation. For example, our own experiments indicate that preexperimental biases interact with experimentally defined predictiveness to produce the general advantage for gangs over suburbs observed in Experiments 2 and 4. Likewise, previous studies have demonstrated that people’s “folk theories” about groups are related to the extent to which those groups support stereotypes (Martin & Parker, 1995), and instructions regarding the entitativity of groups will influence the extent to which those groups will engage in stereotype formation (Crawford et al., 2002). It is possible that either or both of these mechanisms could contribute to the general advantage for gangs over suburbs.

More generally, the suggestion that certain features might be favored over others in stereotype formation is not without precedent. Hamilton and Sherman (1994) noted that certain categories (in particular gender, race, and age) play a greater role in person evaluation than others, and they offered three possibilities for the processing advantage maintained by these *primitive categories*.

The first was that these features provide broad categories, providing a good basis on which to divide up our experience with other people. The second suggestion was that primitive categories represent features that are salient in a person’s appearance and hence are immediately obvious to perceivers. The third suggestion corresponded to the principle of discounting discussed earlier: Primitive categories represent features that have in the past been predictive of behavior and hence will cause perceivers to discount the influence of other available categories when faced with occurrences of similar behavior in future. Similarly, Rothbart and Taylor (1992) argued that social categories that form *natural kinds* (those based on fundamental and unalterable distinctions, e.g., race or gender) will be more important in stereotyping than those that do not.

The fact that so many influences combine to determine the extent to which a cue engages in stereotype formation means that it is very difficult to identify a “pure,” unambiguous example of the influence of predictiveness on stereotyping in the real world; hence the value of controlled experiments in establishing the potential for such effects. Likewise, many of the factors influencing stereotyping that are listed in the preceding paragraphs fall beyond the scope of the simple associative model presented here. For example, given that this model is essentially blind to content, it will not distinguish between a social category that forms a natural kind and one that does not. We must stress that our argument is not that predictiveness is the sole factor that determines stereotype formation, merely that it is one factor that might make a contribution. Similarly, the associative theory presented here is not intended to provide a comprehensive model of all of the potential influences on stereotyping, merely to provide a parsimonious explanation of one such influence.

More speculatively, we could perhaps add predictiveness to Hamilton and Sherman’s (1994) list: It is conceivable that at least part of the advantage for features belonging to primitive categories lies in the fact that these features have previously been experienced as predictive of any of a wide range of properties, including aspects that are entirely unrelated to the behavior currently under consideration (e.g., gender is a reliable predictor of height, pitch of voice, body shape, etc.).

### Stereotype Activation and Expression, and Multiply Categorizable Objects

Spears (2002) has noted that stereotype formation is a “strangely neglected” topic within social psychology, that social cognition research tends to treat stereotypes as givens—“cognitive heuristics, which are part of our mental repertoire . . . that are activated and then applied” (p. 127). Although this approach has yielded much interesting data, it tells us little about how the stereotypes came into being in the first place. It is this latter issue to which the current article is addressed. As formal models of how knowledge structures are dynamically acquired, associative learning theories (such as Mackintosh’s [1975] model) are readily applied to the issue of stereotype formation, just as they apply to other nonsocial examples of category learning.

In support of Spears’s (2002) argument, although the current article represents, to the best of our knowledge, the first experi-

mental study of stereotype *formation* with regard to targets that simultaneously belong to more than one category, previous studies have investigated the *activation and expression* of existing stereotypes for such multiply categorizable targets (Gilbert & Hixon, 1991; Macrae et al., 1995; Shih et al., 2002; Smith et al., 1996; Zárate & Smith, 1990). That is, these studies have assumed the existence of stereotypes regarding particular groups and have looked at the circumstances under which those stereotypes might be activated and applied. Consider the target of a Black male: Under what circumstances will this target elicit stereotypes relating to race as opposed to stereotypes relating to gender? This question has been addressed by self-categorization theory, through the concept of *accessibility* (Oakes, 1987; see also Bruner, 1957): the relative readiness of a category to become activated in the perceiver. Thus a perceiver for whom the “Black” category is more accessible than the “male” category will ascribe more race-related than gender-related stereotypes to a Black male. But what influences the accessibility of a category? Oakes (1987) noted that the two primary determinants of accessibility are “the likelihood of particular types of objects or events occurring in the perceiver’s present environment” and “the current tasks, goals and purposes of the perceiver” (p. 127).

Smith and Zárate (1992) suggested a similar approach with their exemplar-based model of stereotype expression. They argued that a perceiver’s relative attention to different categorization dimensions would influence the extent to which an individual is categorized on that dimension, which would determine the stereotypic attributes ascribed to that individual. For example, a perceiver paying attention to race would ascribe to a Black male more Black-stereotypical attributes than male-stereotypical attributes. Smith and Zárate stated that “answering questions about why perceivers attend to some social dimensions (such as race and gender) and not to others requires going beyond the types of cognitive processes outlined earlier [low-level categorization mechanisms] to enter the realm of social, motivational, and contextual factors” (p. 12). They also suggested that “an activated social motive, such as motives for affiliation, power, or sex, will increase the perceiver’s attention to motive-relevant attributes of social stimulus persons” (Smith & Zárate, 1992, p. 12).

To reiterate, Oakes’s (1987) concept of accessibility and Smith and Zárate’s (1992) attentional exemplar model relate only to the activation or expression of existing stereotypes and do not address the question of how these stereotypes form in the first place. This latter issue is the focus of the current article, where we demonstrate that attention-like processes, based on the learned predictiveness of groups, modulate stereotype formation. Moreover, these attention-like processes are influenced by the predictiveness of groups with respect to information that is entirely independent of, and hence objectively irrelevant to, the perceiver’s “activated social motive” (in Smith & Zárate’s terms) or “current goal” (in Oakes’s)—that is, behavior evaluation.

It is possible, of course, that predictiveness might have a similar biasing effect on stereotype expression: The prior predictiveness of a category might be one determinant of the accessibility of that category (or, in Smith and Zárate’s [1992] terms, the attention paid to that category). Thus it is possible that attention-like processes in stereotype expression are influenced by the types of low-level associative processes that Smith and Zárate reject, and it is

straightforward to provide a simple, formal mechanism for how such processes might operate: We have used Mackintosh’s (1975) theory to model attentional influences on the *learning* of stereotypes, but a similar mechanism could also influence the expression of already-learned stereotypes. It is clearly of value to bring a level of formalization to the rather underspecified concepts of accessibility and attention outlined above—despite their suggestion of several abstract, high-level properties that might influence attention to social dimensions, Smith and Zárate do not provide any mechanism, formal or otherwise, for how attention might change. However, in the absence of experimental evidence, the suggestion that predictiveness might influence stereotype expression in the same way that it influences stereotype formation must remain speculative.

## Conclusions

Our findings support the idea that at least some aspects of complex and behaviorally significant examples of human learning, such as stereotype formation, can be understood in terms of established principles of associative learning. We have argued that differences in prior predictiveness can produce biases in stereotype formation, which might explain (in part) why certain cues are favored over others in stereotype formation. However, while the psychological mechanisms underlying stereotype formation seem to be subject to predictiveness bias (as demonstrated by experiments in which the cue–outcome contingencies during Stage 1 are statistically independent of those in Stage 2), this is not to say that such mechanisms will necessarily lead to learning of inappropriate or spurious relationships. In the real world, it is possible (indeed likely) that cues that have been experienced as predictive in the past will also be accurate predictors of other information in future. Under such circumstances it is of course advantageous for information processing to be focused on these cues at the expense of others, as this will speed the learning of potentially important predictive information. It is precisely this optimizing and focusing role that associability mechanisms might play.

## References

- Abele, A. E., Cuddy, A. J. C., Judd, C. M., & Yzerbyt, V. Y. (2008). Fundamental dimensions of social judgment. *European Journal of Social Psychology, 38*, 1063–1065.
- Allan, L. G. (1980). A note on measurement of contingency between two binary variables in judgment tasks. *Bulletin of the Psychonomic Society, 15*, 147–149.
- Arkes, H. R., & Harkness, A. R. (1983). Estimates of contingency between two dichotomous variables. *Journal of Experimental Psychology: General, 112*, 117–135.
- Baeyens, F., & De Houwer, J. (1995). Evaluative conditioning is a qualitatively distinct form of classical conditioning: A reply to Davey (1994). *Behaviour Research and Therapy, 33*, 825–831.
- Bargh, J. A., Chen, M., & Burrows, L. (1996). Automaticity of social behavior: Direct effects of trait construct and stereotype activation on action. *Journal of Personality and Social Psychology, 71*, 230–244.
- Beesley, T., & Le Pelley, M. E. (2009). The effect of predictive history on the learning of sub-sequence contingencies. *Quarterly Journal of Experimental Psychology*. Advance online publication. doi:10.1080/17470210902831767
- Berndsen, M., Spears, R., van der Pligt, J., & McGarty, C. (2002). Illusory

- correlation and stereotype formation: Making sense of group differences and cognitive biases. In C. McGarty, V. Y. Yzerbyt, & R. Spears (Eds.), *Stereotypes as explanations: The formation of meaningful beliefs about social groups* (pp. 90–110). Cambridge, England: Cambridge University Press.
- Blair, I. V., & Banaji, M. R. (1996). Automatic and controlled processes in stereotype priming. *Journal of Personality and Social Psychology, 70*, 1142–1163.
- Bonardi, C., Graham, S., Hall, G., & Mitchell, C. J. (2005). Acquired distinctiveness and equivalence in human discrimination learning: Evidence for an attentional process. *Psychonomic Bulletin & Review, 12*, 88–92.
- Bruner, J. S. (1957). On perceptual readiness. *Psychological Review, 64*, 123–152.
- Campbell, D. T. (1958). Common fate, similarity, and other indices of the status of aggregates of persons as social entities. *Behavioral Science, 3*, 14–25.
- Chapman, L. J., & Chapman, J. P. (1967). Genesis of popular but erroneous psychodiagnostic observations. *Journal of Abnormal and Social Psychology, 72*, 193–204.
- Crawford, M. T., Sherman, S. J., & Hamilton, D. L. (2002). Perceived entitativity, stereotype formation, and the interchangeability of group members. *Journal of Personality and Social Psychology, 83*, 1076–1094.
- Davey, G. C. L. (1994). Is evaluative conditioning a qualitatively distinct form of classical conditioning? *Behaviour Research and Therapy, 32*, 291–299.
- De Houwer, J., & Beckers, T. (2002). A review of recent developments in research and theories on human contingency learning. *Quarterly Journal of Experimental Psychology, 55B*, 289–310.
- De Houwer, J., Thomas, S., & Baeyens, F. (2001). Associative learning of likes and dislikes: A review of 25 years of research on human evaluative conditioning. *Psychological Bulletin, 127*, 853–869.
- De Houwer, J., Vandorpe, S., & Beckers, T. (2005). On the role of controlled cognitive processes in human associative learning. In A. J. Wills (Ed.), *New directions in human associative learning* (pp. 41–63). Mahwah, NJ: Erlbaum.
- Dijksterhuis, A., Spears, R., & Lepinasse, V. (2001). Reflecting and deflecting stereotypes: Assimilation and contrast in impression formation and automatic behavior. *Journal of Experimental Social Psychology, 37*, 286–299.
- Eagly, A. H., & Steffen, V. J. (1984). Gender stereotypes stem from the distribution of women and men into social roles. *Journal of Personality and Social Psychology, 46*, 735–754.
- Eagly, A. H., & Steffen, V. J. (1986). Gender and aggressive behavior: A meta-analytic review of the social psychological literature. *Psychological Bulletin, 100*, 309–330.
- Eiser, J. R., & Stroebe, W. (1972). *Categorization and social judgment*. London, England: Academic Press.
- Fiedler, K. (1991). The tricky nature of skewed frequency tables: An information loss account of distinctiveness-based illusory correlations. *Journal of Personality and Social Psychology, 60*, 24–36.
- Fiedler, K. (1996). Explaining and simulating judgment biases as an aggregation phenomenon in probabilistic, multiple-cue environments. *Psychological Review, 103*, 193–214.
- Field, A. P., & Davey, G. C. L. (1997). Conceptual conditioning: Evidence for an artifactual account of evaluative learning. *Learning and Motivation, 28*, 446–464.
- Fiske, S. T. (1998). Stereotyping, prejudice, and discrimination. In D. T. Gilbert, S. T. Fiske, & G. Lindzey (Eds.), *The handbook of social psychology* (Vol. 2, pp. 357–414). Boston, MA: McGraw-Hill.
- Fiske, S. T., Cuddy, A. J. C., Glick, P., & Xu, J. (2002). A model of (often mixed) stereotype content: Competence and warmth respectively follow from perceived status and competition. *Journal of Personality and Social Psychology, 82*, 878–902.
- Fiske, S. T., & Neuberg, S. L. (1990). A continuum of impression formation, from category-based to individuating processes: Influences of information and motivation on attention and interpretation. In M. P. Zanna (Ed.), *Advances in experimental social psychology* (Vol. 23, pp. 1–74). New York, NY: Academic Press.
- Ford, T. E., & Stangor, C. (1992). The role of diagnosticity in stereotype formation: Perceiving group means and variances. *Journal of Personality and Social Psychology, 63*, 356–367.
- Fuhrman, R. W., Bodenhausen, G. V., & Lichtenstein, M. (1989). On the trait implications of social behaviors: Kindness, intelligence, goodness, and normality ratings for 400 behavior statements. *Behavior Research Methods, Instruments, & Computers, 21*, 587–597.
- George, D. N., & Pearce, J. M. (1999). Acquired distinctiveness is controlled by stimulus relevance not correlation with reward. *Journal of Experimental Psychology: Animal Behavior Processes, 25*, 363–373.
- Gilbert, D. T., & Hixon, J. G. (1991). The trouble of thinking: Activation and application of stereotypic beliefs. *Journal of Personality and Social Psychology, 60*, 509–517.
- Gopnik, A., & Schulz, L. (2007). *Causal learning: Psychology, philosophy, and computation*. Oxford, England: Oxford University Press.
- Griffiths, O., & Le Pelley, M. E. (2009). Attentional changes in blocking are not a consequence of lateral inhibition. *Learning & Behavior, 37*, 27–41.
- Griffiths, O., & Mitchell, C. J. (2008). Selective attention in human associative learning and recognition memory. *Journal of Experimental Psychology: General, 137*, 626–648.
- Hamilton, D. L. (1981). Stereotyping and intergroup behavior: Some thoughts on the cognitive approach. In D. L. Hamilton (Ed.), *Cognitive processes in stereotyping and intergroup behavior* (pp. 333–354). Hillsdale, NJ: Erlbaum.
- Hamilton, D. L., & Gifford, R. K. (1976). Illusory correlation in interpersonal perception: A cognitive basis of stereotypic judgments. *Journal of Experimental Social Psychology, 12*, 392–407.
- Hamilton, D. L., & Rose, R. L. (1980). Illusory correlation and the maintenance of stereotypic beliefs. *Journal of Personality and Social Psychology, 39*, 832–845.
- Hamilton, D. L., & Sherman, S. J. (1994). Stereotypes. In R. S. Wyer & T. K. Srull (Eds.), *Handbook of social cognition* (Vol. 2, pp. 1–68). Hillsdale, NJ: Erlbaum.
- Hilton, J. L., & von Hippel, W. (1996). Stereotypes. *Annual Review of Psychology, 47*, 237–271.
- Kelley, H. H. (1972). Attribution in social interaction. In E. E. Jones et al. (Eds.), *Attribution: Perceiving the causes of behavior* (pp. 1–26). Morristown, NJ: General Learning Press.
- Klauer, K. C., & Meiser, T. (2000). A source-monitoring analysis of illusory correlations. *Personality and Social Psychology Bulletin, 26*, 1074–1093.
- Kruschke, J. K. (2001). Towards a unified model of attention in associative learning. *Journal of Mathematical Psychology, 45*, 812–863.
- Kruschke, J. K. (2003). Attention in learning. *Current Directions in Psychological Science, 12*, 171–175.
- Le Pelley, M. E. (2004). The role of associative history in models of associative learning: A selective review and a hybrid model. *Quarterly Journal of Experimental Psychology, 57B*, 193–243.
- Le Pelley, M. E. (in press). Attention and human associative learning. In C. J. Mitchell & M. E. L. Pelley (Eds.), *Attention and learning*. Oxford, England: Oxford University Press.
- Le Pelley, M. E., Beesley, T., & Griffiths, O. (2009). *Overt attention and predictive validity in human associative learning*. Manuscript in preparation.
- Le Pelley, M. E., & McLaren, I. P. L. (2003). Learned associability and

- associative change in human causal learning. *Quarterly Journal of Experimental Psychology*, 56B, 68–79.
- Le Pelley, M. E., Oakeshott, S. M., & McLaren, I. P. L. (2005). Blocking and unblocking in human causal learning. *Journal of Experimental Psychology: Animal Behavior Processes*, 31, 56–70.
- Le Pelley, M. E., Suret, M. B., & Beesley, T. (2009). Learned predictive-ness effects in humans: A function of learning, performance, or both? *Journal of Experimental Psychology: Animal Behavior Processes*, 35, 312–327.
- Lippmann, W. (1922). *Public opinion*. New York, NY: Harcourt Brace.
- Livesey, E. J., Harris, I. M., & Harris, J. A. (2009). Attentional changes during implicit learning: Signal validity protects a target stimulus from the attentional blink. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 35, 408–422.
- Lubow, R. E. (1989). *Latent inhibition and conditioned attention theory*. Cambridge, England: Cambridge University Press.
- Lubow, R. E., & Gewirtz, J. C. (1995). Latent inhibition in humans: Data, theory, and implications for schizophrenia. *Psychological Bulletin*, 117, 87–103.
- Mackintosh, N. J. (1973). Stimulus selection: Learning to ignore stimuli that predict no change in reinforcement. In R. A. Hinde & J. S. Hinde (Eds.), *Constraints on learning* (pp. 75–96). London, England: Academic Press.
- Mackintosh, N. J. (1975). A theory of attention: Variations in the associability of stimuli with reinforcement. *Psychological Review*, 82, 276–298.
- Macrae, C. N., Bodenhausen, G. V., & Milne, A. B. (1995). The dissection of selection in person perception: Inhibitory processes in social stereotyping. *Journal of Personality and Social Psychology*, 69, 397–407.
- Macrae, C. N., Bodenhausen, G. V., Milne, A. B., & Calvini, G. (1999). Seeing more than we can know: Visual attention and category activation. *Journal of Experimental Social Psychology*, 35, 590–602.
- Macrae, C. N., Milne, A. B., & Bodenhausen, G. V. (1994). Stereotypes as energy saving devices: A peek inside the cognitive toolbox. *Journal of Personality and Social Psychology*, 66, 37–47.
- Martin, C. L., & Parker, S. (1995). Folk theories about sex and race differences. *Personality and Social Psychology Bulletin*, 21, 45–57.
- McGarty, C., & de la Haye, A.-M. (1997). Stereotype formation: Beyond illusory correlation. In R. Spears, P. J. Oakes, N. Ellemers, & S. A. Haslam (Eds.), *The social psychology of stereotyping and group life* (pp. 144–170). Oxford, England: Blackwell.
- McGarty, C., Haslam, S. A., Turner, J. C., & Oakes, P. J. (1993). Illusory correlation as accentuation of actual inter-category difference: Evidence for the effect with minimal stimulus information. *European Journal of Social Psychology*, 23, 391–410.
- Mitchell, C. J. (in press). Attention and memory in human learning. In C. J. Mitchell & M. E. L. Pelley (Eds.), *Attention and learning*. Oxford, England: Oxford University Press.
- Murphy, R. A., Schmeer, S., Mondragon, E., Vallee-Tourangeau, F., & Hilton, D. (2009). *Making the illusory correlation effect appear and then disappear: The effects of increased learning*. Manuscript submitted for publication.
- Oakes, P. J. (1987). The salience of social categories. In J. C. Turner, M. A. Hogg, P. J. Oakes, S. D. Reicher, & M. S. Wetherell (Eds.), *Rediscovering the social group: A self-categorization theory* (pp. 117–141). Oxford, England: Blackwell.
- Oakes, P. J., Haslam, S. A., & Turner, J. C. (1994). *Stereotyping and social reality*. Oxford, England: Blackwell.
- Oswald, C. J. P., Yee, B. K., Rawlins, J. N. P., Bannerman, D. B., Good, M., & Honey, R. C. (2001). Involvement of the entorhinal cortex in a process of attentional modulation: Evidence from a novel variant of an IDS/EDS procedure. *Behavioral Neuroscience*, 115, 841–849.
- Pearce, J. M., & Hall, G. (1980). A model for Pavlovian conditioning: Variations in the effectiveness of conditioned but not of unconditioned stimuli. *Psychological Review*, 87, 532–552.
- Pryor, J. B. (1986). The influence of different encoding sets upon the formation of illusory correlations and group impressions. *Personality and Social Psychology Bulletin*, 12, 216–226.
- Rothbart, M., & Taylor, M. (1992). Category labels and social reality: Do we view social categories as natural kinds? In G. R. Semin & K. Fiedler (Eds.), *Language and social cognition* (pp. 11–36). London, England: Sage.
- Sartre, J. P. (1948). *Anti-Semite and Jew*. New York, NY: Schocken Press.
- Schaller, M. (1994). The role of statistical reasoning in the formation, preservation and prevention of group stereotypes. *British Journal of Social Psychology*, 33, 47–61.
- Shanklee, H., & Mims, M. (1982). Sources of error in judging event covariations: Effects of memory demands. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 8, 208–224.
- Shanks, D. R. (2007). Associationism and cognition: Human contingency learning at 25. *Quarterly Journal of Experimental Psychology*, 60, 291–309.
- Shanks, D. R., & Darby, R. J. (1998). Feature- and rule-based generalization in human associative learning. *Journal of Experimental Psychology: Animal Behavior Processes*, 24, 405–415.
- Sherif, M. (1967). *Group conflict and co-operation: Their social psychology*. London, England: Routledge and Kegan Paul.
- Sherman, J. W., Kruschke, J. K., Sherman, S. J., Percy, E. J., Petrocelli, J. V., & Conrey, F. R. (2009). Attentional processes in stereotype formation: A common model for category accentuation and illusory correlation. *Journal of Personality and Social Psychology*, 96, 305–323.
- Shih, M., Ambady, N., Richeson, J. A., Fujita, K., & Gray, H. M. (2002). Stereotype performance boosts: The impact of self-relevance and the manner of stereotype activation. *Journal of Personality and Social Psychology*, 83, 638–647.
- Sloman, S. A. (1996). The empirical case of two systems of reasoning. *Psychological Bulletin*, 119, 3–22.
- Smith, E. R. (1991). Illusory correlation in a simulated exemplar-based memory. *Journal of Experimental Social Psychology*, 27, 107–123.
- Smith, E. R., & DeCoster, J. (1998). Knowledge acquisition, accessibility, and use in person perception and stereotyping: Simulation with a recurrent connectionist network. *Journal of Personality and Social Psychology*, 74, 21–35.
- Smith, E. R., Fazio, R. H., & Cejka, M. A. (1996). Accessible attitudes influence categorization of multiply categorizable objects. *Journal of Personality and Social Psychology*, 71, 888–898.
- Smith, E. R., & Zárate, M. A. (1992). Exemplar-based model of social judgment. *Psychological Review*, 99, 3–21.
- Spears, R. (2002). Four degrees of stereotype formation: Differentiation by any means necessary. In C. McGarty, V. Y. Yzerbyt, & R. Spears (Eds.), *Stereotypes as explanations* (pp. 127–156). Cambridge, England: Cambridge University Press.
- Stangor, C., Carr, C., & Kiang, L. (1998). Activating stereotypes undermines task performance expectations. *Journal of Personality and Social Psychology*, 75, 1191–1197.
- Stangor, C., & Lange, J. (1994). Mental representations of social groups: Advances in understanding stereotypes and stereotyping. *Advances in Experimental Social Psychology*, 26, 357–416.
- Stroessner, S. J., Hamilton, D. L., & Mackie, D. M. (1992). Affect and stereotyping: The effect of induced mood on distinctiveness-based illusory correlations. *Journal of Personality and Social Psychology*, 62, 564–576.
- Tajfel, H. (1957). Value and the perceptual judgment of magnitude. *Psychological Review*, 64, 192–204.
- Tajfel, H. (1982). Social psychology of intergroup relations. *Annual Review of Psychology*, 33, 1–39.

- Tajfel, H., & Wilkes, A. L. (1963). Classification and quantitative judgment. *British Journal of Psychology*, *54*, 101–114.
- Van Rooy, D., Van Overwalle, F., Vanhooymissen, T., Labiouse, C., & French, R. (2003). A recurrent connectionist model of group biases. *Psychological Review*, *110*, 536–563.
- Ward, W. D., & Jenkins, H. M. (1965). The display of information and the judgment of contingency. *Canadian Journal of Psychology*, *19*, 231–241.
- Zarate, M. A., & Smith, E. R. (1990). Person categorization and stereotyping. *Social Cognition*, *8*, 161–185.

## Appendix

### Instructions to Participants

#### On-Screen Instructions Presented to Participants at the Outset of Experiment 1A

In this experiment we are interested in how people retain and process information that is presented to them visually. You will be studying information about a number of people who belong to certain gangs. For the sake of anonymity, these people are identified by their initials. Each person belongs to two gangs. For example, you might be told that:

“K.F. is a member of gang X, and a member of gang Y.”

On each trial in the first stage of this experiment you will see a statement describing a particular person, along with pictures of two people. Your job is to decide which of these pictures shows the person referred to in the statement at the top of the screen.

To enter your decision, click on the picture that you think is the person described at the top of the screen. When you have made your decision, click the OK button. The computer will then tell you whether your decision was correct or incorrect. A blue box will appear, indicating the picture that actually shows the person described in the statement. If you make an incorrect decision the computer will beep.

You will have to guess at first, but with the aid of the feedback your predictions should soon start to become more accurate. Your reaction times are not important: You may take as long as you like on each trial. Please do not write anything down at any point during the experiment.

#### Instructions Presented Prior to Stage 2

In the second stage of this experiment we are again interested in how people retain and process information that is presented to them visually. On each trial you will again be given information about which gangs a particular person belongs to. You will also now see a sentence describing a behavior performed by that person. For example, you might be told that a particular person “tried not to take sides when two of his friends had an argument,” or that a person “attempted to push into the middle of a queue.”

The gangs involved in this stage of the experiment will be the same as for the previous stage. In collecting behavior descriptions of people for this experiment we tried to draw a random sample from the population.

You will now be shown a rather large number of screens providing information about people along with statements about their behavior. On each trial, please read carefully all of the information presented on the screen (both the information describing the person and the information describing their behavior). You may find it helpful to read the information out loud.

After a fixed period of time, a button will appear. When you are ready, click this button to move on to the next trial.

#### Instructions Presented Prior to Test Phase

You will now be asked to rate your opinions of people who are members of the different gangs. Specifically, you are asked to rate how much you like these people, based on the behaviors that you experienced in the previous stage.

On each trial you will be asked how much you like people who are members of a particular gang. When deciding on how much you like a particular group of people, think about how much you would like a person who belongs to that gang to be a friend of yours.

Your ratings will be entered on a scale from 0 to 10, where a rating of 0 indicates that you **STRONGLY DISLIKE** people from the particular gang mentioned at the top of the screen, and a rating of 10 indicates that you **STRONGLY LIKE** people from the particular gang mentioned at the top of the screen. You may use any value from 0 to 10 to indicate your opinion.

To enter your rating, click on the appropriate option button. When you have entered your rating, click the OK button to continue.

Received February 15, 2009

Revision received October 14, 2009

Accepted October 15, 2009 ■