

Behavioral and Neural Properties of Social Reinforcement Learning

Rebecca M. Jones,¹ Leah H. Somerville,¹ Jian Li,² Erika J. Ruberry,¹ Victoria Libby,¹ Gary Glover,³ Henning U. Voss,⁴ Douglas J. Ballon,⁴ and B. J. Casey¹

¹Sackler Institute for Developmental Psychobiology, Weill Cornell Medical College, New York, New York 10065, ²Department of Psychology, New York University, New York, New York 10003, ³Lucas Center for Imaging, Department of Radiology, Stanford University, Stanford, California 94305, and ⁴Citigroup Biomedical Imaging Center, Department of Radiology, Weill Cornell Medical College, New York, New York, 10065

Social learning is critical for engaging in complex interactions with other individuals. Learning from positive social exchanges, such as acceptance from peers, may be similar to basic reinforcement learning. We formally test this hypothesis by developing a novel paradigm that is based on work in nonhuman primates and human imaging studies of reinforcement learning. The probability of receiving positive social reinforcement from three distinct peers was parametrically manipulated while brain activity was recorded in healthy adults using event-related functional magnetic resonance imaging. Over the course of the experiment, participants responded more quickly to faces of peers who provided more frequent positive social reinforcement, and rated them as more likeable. Modeling trial-by-trial learning showed ventral striatum and orbital frontal cortex activity correlated positively with forming expectations about receiving social reinforcement. Rostral anterior cingulate cortex activity tracked positively with modulations of expected value of the cues (peers). Together, the findings across three levels of analysis—social preferences, response latencies, and modeling neural responses—are consistent with reinforcement learning theory and nonhuman primate electrophysiological studies of reward. This work highlights the fundamental influence of acceptance by one's peers in altering subsequent behavior.

Introduction

Successfully navigating our social environment depends on learning from positive and negative encounters with others and shaping future behavior toward those individuals. Psychologists have proposed that positive social exchanges are fundamentally rewarding for humans (Bandura and Walters, 1963; Baumeister and Leary, 1995; Steinberg, 2008), suggesting that learning from social interactions may draw on basic reinforcement learning mechanisms. The present study was designed to test this hypothesis by building on reinforcement learning studies in nonhuman primates and human imaging studies (Schultz et al., 1997; Fiorillo et al., 2003; McClure et al., 2003; D'Ardenne et al., 2008).

Reinforcement learning from primary (e.g., food) and secondary (e.g., money) reinforcers has been shown to engage specific neural circuitry. In its simplest form, it is explained by the classic Rescorla–Wagner model (Rescorla and Wagner, 1972). According to this model, learning to associate arbitrary cues with

positive outcomes results in expectations of future positive outcomes in the presence of these cues. If there are discrepancies between the expected outcome to the cue and the actual outcome, a prediction error signal is generated. Nonhuman primate and human imaging studies have implicated the ventral striatum and orbital frontal cortex (OFC) in prediction error signaling (Schultz et al., 1997; Berns et al., 2001; Fiorillo et al., 2003; McClure et al., 2003). Studies have shown that as cues become reliably associated with receipt of a reward, manual responses to these cues quicken over time (O'Doherty et al., 2006; Spicer et al., 2007), while others demonstrate changes in choice behaviors based on reinforcement manipulations (Daw et al., 2006; Li and Daw, 2011). The learned association generates a neural signal to the cue that previously was associated with the reward itself (Schultz et al., 1997; O'Doherty et al., 2006). The current study examines whether similar changes in behavior (response latencies) and neural circuitry engaged during basic reinforcement learning are involved during learning within a social context.

This study tests the extent to which social reinforcement learning relies on similar learning mechanisms as those used in basic reinforcement learning. To do so, we created a task in which participants learned to differentiate three peers, each of whom was associated with a unique probability of social reinforcement (i.e., providing socially accepting feedback). Social reinforcement learning processes were evaluated at three levels of analysis—preference ratings, response latencies, and neural responses to expected cue values and prediction errors. We hypothesized that social preference ratings would become more favorable and response latencies would become faster toward the peer with the

Received June 13, 2011; revised July 12, 2011; accepted July 19, 2011.

Author contributions: R.M.J., L.H.S., and B.J.C. designed research; R.M.J., E.J.R., and V.L. performed research; J.L., G.G., H.U.V., D.J.B., and B.J.C. contributed unpublished reagents/analytic tools; R.M.J., L.H.S., and J.L. analyzed data; R.M.J., L.H.S., and B.J.C. wrote the paper.

This work was supported by NIDA R01 DA018879, NIDA T-32 Training Grant DA007274, the Mortimer D. Sackler family, and the Dewitt-Wallace fund. We gratefully acknowledge the assistance of the resources and staff at the Biomedical Imaging Core Facility of the Citigroup Biomedical Imaging Center at Weill Cornell Medical College.

The authors report no conflicts of interest.

Correspondence should be addressed to either Rebecca M. Jones or B. J. Casey, Sackler Institute for Developmental Psychobiology, Weill Cornell Medical College, 1300 York Avenue, Box 140, New York, NY 10065. E-mail: rej2004@med.cornell.edu or bj2002@med.cornell.edu.

DOI:10.1523/JNEUROSCI.2972-11.2011

Copyright © 2011 the authors 0270-6474/11/3113039-07\$15.00/0

greatest probability of providing social acceptance to the participant. We applied a simple Rescorla–Wagner rule in behavioral and functional imaging analyses to target the neural bases of these behavioral changes, hypothesizing that the ventral striatum and OFC would code prediction error signals (Schultz et al., 1997; O’Doherty, 2007). Thus, the current study elucidates neurobiological mechanisms for key learning processes during social exchanges that shape behavior through positive interactions.

Materials and Methods

Participants. Forty-six adults (aged 18–28 years; 22 females) participated in the experiment. Thirty-six completed the task during functional magnetic resonance imaging (fMRI) (aged 18–28 years; all right-handed; 19 females). Three individuals in the fMRI group were excluded due to insufficient number of correct trials in any condition ($n = 2$; 1 male) or noncompliance with the task ($n = 1$, male). Participants reported no history of neurological and/or psychiatric disorders in a standard screening or on the Structured Clinical Interview for DSM-IV Axis I Disorders (SCID) (First et al., 2007) and imaging participants reported no contraindications for an MRI. Two participants did not complete the SCID due to time constraints. All participants provided written consent approved by the Institutional Review Board at Weill Cornell Medical College and were debriefed and compensated following their participation.

Experiment cover story. The experiment was conducted during two separate sessions. The first session introduced the cover story, leading participants to believe they would receive actual social feedback during a task that would be completed on the second visit. Participants were shown up to five photographs of gender- and ethnicity-matched peers. They then selected three with whom they would like to interact, and rated the three peers for how likeable and attractive they looked on a scale from 1 (not very) to 10 (very). Participants also completed a personal survey where they listed information about themselves (birthday; hometown; and favorite music, TV shows, books, quotes, and activities). Participants were told that each of the three selected peers would see their survey over the next few days as well as the surveys of two other supposed participants. These three peers would write notes indicating a positive interest in the participant’s survey or in one of the other two surveys. Participants were told that each of these individuals could write a small number of notes, emphasizing their limited number and enhancing the positive value of receiving a note. Participants were then scheduled for a second session.

At the second session, participants were told that the experimenters had compiled the notes from the three selected peers. During the experiment, participants would be shown how often each of the peers decided to write notes to them (positive social reinforcement) or to one of the other supposed participants (no positive social reinforcement). Although it is possible that participants experienced the no positive social reinforcement trials as mildly rejecting, we have chosen not to adopt this interpretation because we do not have conclusive data supporting this possibility. Rather, these operational definitions were selected for consistency with studies of basic reward learning. At the beginning of the second session, participants were also reminded that receiving a note symbolized that the peer was interested in something written in their personal survey.

Unbeknownst to the participants, peer interaction (i.e., delivery of notes) was experimentally manipulated such that each of the three peers was associated with a distinct probability of social reinforcement (Fig. 1A) with Rare interaction defined by positive social reinforcement on 33% of the trials and no positive social reinforcement on 66% of the

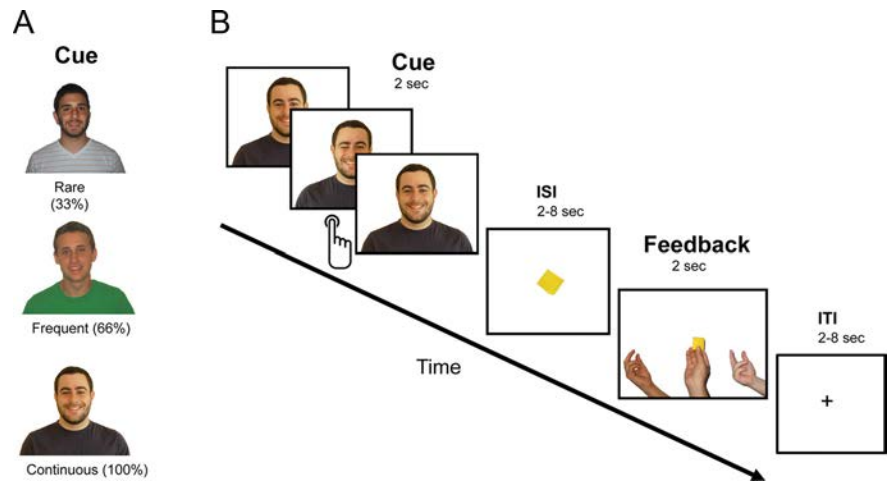


Figure 1. Task parameters. **A**, Three peers chosen by the participant were associated with distinct probabilities of positive reinforcement. **B**, Schematic of one trial within a run. The face of one peer (Cue) was displayed for 2 s, during which the face stimulus winked (500 ms) and participants pressed one of two buttons indicating in which eye the wink occurred, followed by a variable interstimulus interval (ISI), followed by the note outcome (Feedback). In this example, the participant received the note (positive social reinforcement) because it appeared in the middle hand. If the note appeared in one of the hands to the left or to the right of the middle hand, the participant did not receive the note (no positive social reinforcement). A variable intertrial interval (ITI) followed.

trials; Frequent interaction defined by positive social reinforcement on 66% of the trials and no positive social reinforcement on 33% of the trials; and Continuous interaction defined by positive social reinforcement on all trials (100%). The probability of reinforcement associated with each of the face stimuli was counterbalanced across participants to equate for low-level stimulus features across conditions.

Task parameters. At the start of each trial (Fig. 1B), a picture of one of the three peers was presented for two seconds (Cue). During the two seconds, the stimulus would wink for 500 ms in either the left or right eye, indicating that a note was ready to be passed. Participants signaled that they were ready to receive the note by pressing one of two buttons indicating whether the wink was in the left or the right eye. This behavioral component was included to ensure attention and to collect reaction time data as an index of learning about the reinforcement contingencies for each of the three peers across the experiment. After a jittered interstimulus interval of a picture of a folded note (2, 4, 6, or 8 s), three hands appeared at the bottom of the screen with one hand holding a note for 2 s (Feedback). Participants had been instructed that if the middle hand held the note, this signified that the participant had received a note from that peer (positive social reinforcement). If the note appeared in one of the hands to the left or right of the middle hand, this signified that the note was given to someone else (no positive social reinforcement). If the participant pressed incorrectly or did not respond during the cue, no feedback was given. A jittered intertrial interval (2, 4, 6, or 8 s) followed in which participants rested while viewing a fixation crosshair. Participants viewed 18 trials per run in a pseudorandomized order with six trials per condition (Rare, Frequent, Continuous) for six runs, for a total of 108 trials, 36 trials per condition. To enhance the believability of the cover story and keep participants engaged, one of the notes was shown between each run; these notes were generated by the experimenters and always indicated positive interest in the participant’s personal survey (e.g., “I love playing soccer too, and I am part of a weekend league”, “Where did you go when you visited Hawaii?”, “I also have a golden retriever”).

To further index learning with the reaction time data at the end of the experiment, after the six experimental runs, participants completed a reversal run (18 trials) during which reaction times were recorded. Contingencies were reversed for the Rare and Continuous conditions such that the Rare peer now provided 100% reinforcement to the participant and the Continuous peer now provided 33% reinforcement to the participant. The Frequent peer’s probability (66%) did not change.

The task was presented using E-Prime software, and the participants who completed the task during fMRI viewed images on an overhead

liquid crystal display panel with the Integrated Functional Imaging System-Stand Alone (IFIS-SA; fMRI Devices). E-Prime software, integrated with IFIS-SA, recorded button responses and reaction times using the Fiber Optic Button Response System (Psychology Software Tools).

At the end of the experiment, participants completed posttest ratings of attractiveness and likeability for each peer on the same scale used at the beginning of the experiment. To assess whether participants held explicit knowledge of the social reinforcement contingencies associated with each peer, they were asked whether any of the three peers provided positive reinforcement more often than any others. If the participant said yes, they were asked to describe what pattern they noticed, and descriptions were scored based on whether the participant accurately stated which peer provided the most, intermediate, and least positive social feedback. Three of the 43 participants correctly ranked the three peers in this way and were thus considered explicitly aware of the social reinforcement contingencies. Participants were then debriefed regarding the cover story and the rationale of the experiment.

Image acquisition. Participants were scanned with a Signa HDx 3.0T MRI scanner (General Electric Medical Systems) with a quadrature head coil. A high-resolution, 3D magnetization prepared rapid acquisition gradient echo anatomical scan (MPRAGE) was acquired (256×256 in-plane resolution, FOV = 240 mm; 124 1.5 mm sagittal slices). Functional scans were acquired with a spiral in and out sequence (TR = 2000 ms, TE = 30 ms, flip angle = 90°) (Glover and Thomason, 2004). Twenty-nine 5-mm-thick contiguous coronal slices were acquired per TR, for a total of 129 TRs per functional run with a resolution of 3.125×3.125 mm (64×64 matrix, FOV = 200 mm) covering the entire brain except for the posterior portion of the occipital lobe.

Behavioral analysis. Change in attractiveness and likeability of the peers before and after the task was tested with a 3 (probability: Rare, Frequent, Continuous) \times 2 (time: before task, after task) repeated-measures ANOVA using PASW Statistics 18 software (SPSS). Attractiveness and likeability ratings for three of the 43 participants were lost due to technical error.

Reaction times were analyzed in response to the cue after the wink occurred. Reaction times were z -score transformed for each individual after removing outliers (defined as reaction times 3 SDs above or below the individual's mean reaction time). Changes in reaction times and accuracy for the three conditions during the early and late trials were each tested with a 3 (probability: Rare, Frequent, Continuous) \times 2 [time: first half of trials (early), second half of trials (late)] repeated-measures ANOVA.

To test for reaction time modulation as a function of contingency reversal, we compared reaction times from the sixth run of the experiment to the reversal run with a 2 (probability: Rare, Continuous) \times 2 (time: sixth run, reversal run) repeated-measures ANOVA.

Prior research has demonstrated that not receiving reinforcement on a given trial modulates behavioral responses on the next trial (Liu et al., 2007). To determine whether reinforcement outcome influenced response latencies on the subsequent trial, we compared reaction times from trials when the participant had received positive social reinforcement on the preceding trial to trials when they had not received positive reinforcement using a paired samples t test.

Reinforcement learning model. We used a simple reinforcement learning algorithm (Rescorla–Wagner) to model the trial-by-trial variance in participants' reaction times (Rescorla and Wagner, 1972). The Rescorla–Wagner rule probes learning through a prediction error (PE) signal δ , which is the difference between the experienced outcome (R ; positive social feedback or no positive social feedback) and expected outcome (V) for each trial. PE takes the form of $\delta = R - V$ and can be used to subsequently update expected outcome weighted by a fixed learning rate α : $V_{t+1} = V_t + \alpha\delta_t$ for given trial t . Reaction time has been shown in previous studies to be a reliable indicator of learning contingencies and speeding or slowing in reaction times has been associated with conditioning as predicted by reinforcement learning models (Seymour et al., 2004; Bray and O'Doherty, 2007). We thus fitted the Rescorla–Wagner model to participants' trial-by-trial z -score transformed reaction times using a linear regression model to derive the best-fitting model parameters (α and V_0). We tested the rate of learning for each subject based on his or her

individual reaction time history, which yielded an average learning rate (α) of 0.15 across participants, suggesting learning effects on reaction time measures (one-sample t test of learning rate vs null hypothesis of 0; $p < 0.001$). The average learning rate of participants who completed the behavioral version of the experiment was comparable to the imaging sample ($p > 0.3$), suggesting consistency in our model.

Imaging analysis. The fMRI data analyses were performed with Analysis of Functional Neuroimages (AFNI) software (Cox, 1996). Functional data were slice-time corrected, realigned within and across runs to correct for head movement, coregistered with the high-resolution anatomical scan, scaled to percentage signal change units, and smoothed with a 6 mm full-width at half maximum Gaussian kernel. Images with movement >2 mm along the x , y , or z planes were excluded from the analysis. Functional data were transformed into standard Talairach coordinate space (Talairach and Tournoux, 1988) by using the warping parameters obtained from the Talairach transformation of the high-resolution anatomical scan. Talairach-transformed functional data were resampled to a resolution of $3 \times 3 \times 3$ mm.

For imaging analysis, we generated a linear reinforcement learning model with linear regression using reaction times of all participants to obtain a single set of signed model parameters (α and V_0) that best fit participants' behavior ($r = 0.19$, $p < 0.001$). This approach has been suggested to be less susceptible to extreme parameter value estimation for individual participants and tends to more stable (Daw et al., 2006; Bray and O'Doherty, 2007; Li et al., 2011). The learning rate ($\alpha = 0.07$) defined from modeling of the behavioral data was used to generate the PE and expected outcome values that were included as parametric regressors with signed numbers in individual-subject general linear models.

A general linear model analysis was performed to estimate neural responses to stimuli as a function of reinforcement learning. Each participant's GLM contained five task regressors: (1) cue onset times, defined as the time points at which peer faces were presented; (2) a parametric regressor paired with cue timings containing expected value estimates for each trial (V_t); (3) feedback onset times, containing values corresponding to the time points at which the note feedback was presented; (4) a parametric regressor paired with feedback onset time representing prediction error values (δ_t); and (5) incorrect trial onset times. Task regressors were convolved with a gamma-variate hemodynamic response function. Regressors of noninterest included motion parameters and linear and quadratic trends for each run. Separate random effects group analyses were conducted on individual participant beta estimates for the parametric regressor representing prediction error values (δ_t) and individual participant beta estimates for the parametric regressor representing expected values to the cues (V_t).

To test for basic effects of prediction error during the feedback presentation, a within-subjects voxelwise one-sample t test was performed to identify regions demonstrating activity that positively correlated with prediction error learning signals. To identify neural responses to expected values during the cue presentation of the trials, a within-subjects voxelwise one-sample t test was performed to identify regions showing activity that positively correlated with expected values to the cues. Results of all whole-brain analyses were considered significant by exceeding a p value/cluster size combination ($p < 0.005/50$ voxels) that corresponded to whole-brain $p < 0.05$, corrected for multiple comparisons as calculated with Monte Carlo simulations in AFNI.

As the OFC has been implicated in prior studies of prediction error learning (Berns et al., 2001; Takahashi et al., 2009), we selected this region as an a priori structure in which to search for learning signals. Specifically, an anatomical mask of the bilateral OFC was created, encompassing Brodmann areas 11 and 47 ($x = -50/51$, $y = 10/57$, $z = -3/-23$; $22,086$ mm³ voxels). Group analyses conducted within this mask applied $p < 0.05$ small volume corrected statistical thresholding.

Results

Behavioral

Likeability and attractiveness ratings

Enhanced social preference for peers was modulated by the probability of reinforcement experienced during the experiment, as indicated by a significant interaction between reinforcement

probability and time (pre-interaction, post-interaction) on likeability ratings ($F_{(2,78)} = 5.48, p < 0.01$; Fig. 2A). *Post hoc* analyses indicated that post-task ratings decreased linearly with decreasing interaction probability, such that peers who interacted less with the participant were rated as less likeable (linear term: $F_{(1,39)} = 7.17, p < 0.02$). Whereas pre-task likeability ratings were equivalent for all three peers ($ps > 0.48$), after the task the Frequent ($t_{(39)} = -2.26, p < 0.03$) and Continuous ($t_{(39)} = -2.68, p < 0.02$) peers were rated as more likeable than the rarely reinforcing peer, though there was not a significant difference in likeability ratings after the task between the Frequent and Continuous peers ($t_{(39)} = -0.69, p > 0.49$). Attractiveness ratings were not significantly modulated by task conditions (main effects of reinforcement probability, time, and interactions: $ps > 0.09$).

Accuracy

Participants responded correctly to 95.63% of trials (SD = 3.54%). Response accuracy was not significantly modulated by the task conditions (main effects of reinforcement probability, time, and interactions: $ps > 0.29$).

Reaction time

Response latencies to the cue varied as participants learned the reinforcement contingency outcomes associated with each peer, as indicated by a significant interaction between probability of reinforcement and time ($F_{(2,84)} = 3.98, p < 0.03$; Fig. 2B). *Post hoc* *t* tests showed that, whereas there was no difference in reaction times in the early trials (all $ps > 0.41$), individuals were faster during the late trials for the Frequent reinforcement condition ($t_{(42)} = 2.49, p < 0.02$), compared with the Rare reinforcement condition. There was a trend for responses to be faster in the Continuous reinforcement condition ($t_{(42)} = 2.01, p < 0.06$) than in the Rare reinforcement condition. Overall, participants were faster during the late versus early trials ($F_{(1,42)} = 15.21, p < 0.01$) and there was no main effect of probability of social reinforcement when collapsing across time ($F_{(2,84)} = 1.43, p > 0.25$).

To further test for the effects of learning, we compared reaction times for Rare and Continuous reinforcement before and after reinforcement contingencies were reversed at the end of the experiment. Evidence that participants had implicitly learned the contingencies was further supported by the interaction between time (sixth run vs reversal) and reinforcement probability (rarely reinforcing vs continuously reinforcing) on reaction times ($F_{(1,42)} = 10.15, p < 0.01$; Fig. 2C). *Post hoc* tests showed a significant reaction time speeding when the Rare condition switched to delivering Continuous reinforcement ($t_{(42)} = 3.13, p < 0.01$). There were no main effects of reinforcement probability ($F_{(1,42)} = 2.77, p = 0.1$) or time ($F_{(1,42)} = 1.35, p > 0.25$). There was also no difference in the Frequent condition (unchanged during reversal) reaction times between the last run and the reversal run ($t_{(42)} = -1.48, p > 0.15$).

Additionally, we examined how reaction times (RT) changed based on feedback from the preceding trial as another index of how the reinforcement contingencies altered behavior. We found that participants were faster on the subsequent trial after not receiving positive reinforcement (mean *z*-score RT: $-0.01, SD: 0.13$) versus when they had received positive reinforcement (mean *z*-score RT: $0.07, SD: 0.1$; $t_{(42)} = 2.86, p < 0.01$).

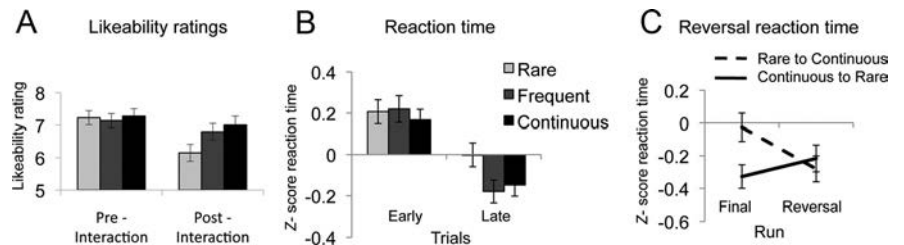


Figure 2. Behavioral responses to cues. **A**, Likeability ratings for the three peers before engaging in the task (pre-interaction) and after the task (post-interaction). **B**, Reaction times to the wink for the three peers, broken down by early and late trials of the experiment. **C**, Reaction times during the final (sixth) run of the experiment and during the run when the contingencies were reversed for the Rare and Continuous conditions.

Table 1. Brain regions reflecting positive correlations with prediction errors

Region	x	y	z	Number of voxels	t
R rostral ACC	2	29	11	316	4.23
R anterior insula	32	11	-10	56	3.55
L ventral striatum	-4	5	5	53	3.89
L orbital frontal cortex+	-32	26	-7	12	3.47

Coordinates represent activation clusters exceeding $p < 0.05$, whole-brain corrected, thresholding, and are listed in Talairach and Tournoux coordinate space. R, Right; L, left; +, small volume corrected.

Imaging

Prediction errors

As indexed by the prediction error parametric regressor, prediction error signals (δ_t) were positively associated with activity in the rostral anterior cingulate cortex, ventral striatum, anterior insula, and OFC (Table 1; Fig. 3). The parametric values in the general linear model encompassed positive and negative prediction errors, demonstrating that the BOLD fluctuations in these regions tracked learning signals reflecting reinforcement expectancies. Together, these findings delineate an orbital frontostriatal circuit showing significantly greater activity associated with the unexpected outcomes of either receiving or not receiving positive social reinforcement.

Expected cue values

We also examined regions of the brain that positively correlated with learning to distinguish the faces of the peers based on their differential rates of positive social reinforcement (learned cue value). Specifically, group analysis of the cue phase of trials that tracked positively with modulations of expected value (V_t) identified greater activity in the rostral anterior cingulate cortex with larger expected value (Fig. 4). No other regions survived whole-brain correction. No regions within the frontostriatal circuitry of interest demonstrated negative correlations with expected value (i.e., brain regions sensitive to lower expected values) at corrected thresholding.

Discussion

Repeated social exchanges shape our behavior toward others. In this experiment, we examined how different probabilities of positive interaction from distinct peers rapidly influence social learning. Within a reinforcement learning framework, we developed a novel social paradigm and demonstrated that the neural systems engaged while forming social expectations are similar to those involved in basic reward learning. This overlap in neural circuitry and function is consistent with prediction error-related learning and with our hypothesis that positive social interactions can serve as secondary reinforcers, taking on the attributes of primary rewards essential to survival (e.g., food).

Over the course of the experiment, participants learned to differentiate each of the cues (peers) by their distinct reinforce-

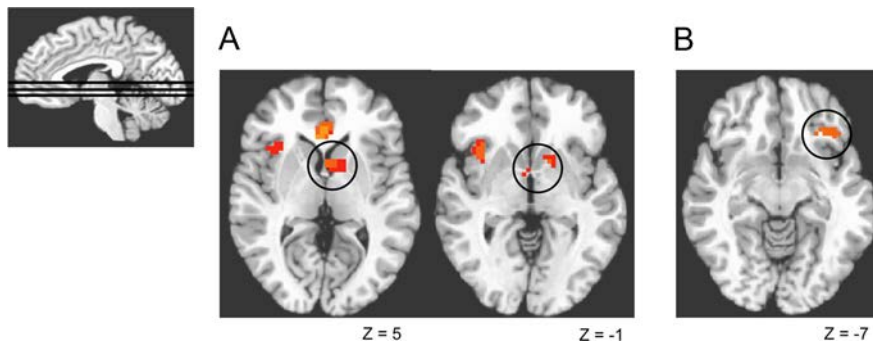


Figure 3. Brain regions reflecting positive correlations with prediction errors. **A**, Circles denote activity in the ventral striatum. Image threshold $p < 0.05$, whole-brain corrected. **B**, Circle denotes activity in the lateral orbital frontal cortex. Image threshold $p < 0.05$, small volume corrected (see Materials and Methods). All statistical activations are displayed on a representative high-resolution axial image. The left side of the image corresponds to the right side of the brain.

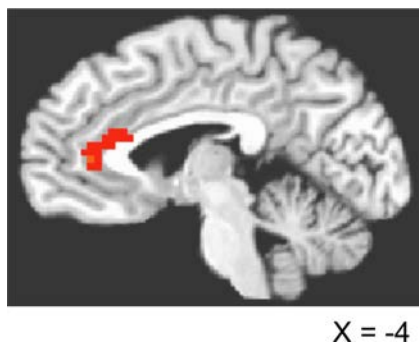


Figure 4. Neural activity with positive correlations with learned cue value. Activity in the rostral anterior cingulate cortex reflects a positive correlation with expected values for the cues. Image threshold $p < 0.05$, whole-brain corrected. Statistical activations are displayed on a representative high-resolution sagittal image.

ment outcomes. Specifically, one peer always provided positive social reinforcement, another one frequently provided positive social reinforcement, and the third rarely provided positive social reinforcement. Ratings of likeability changed from the beginning to the end of the experiment, with less reinforcing peers becoming less likeable, and more reinforcing peers yielding higher ratings of likeability by the end of the task. By asking participants to make a simple button response during the cue presentation, we tested whether speeding of response latencies (action tendencies) indexed learned associations between a given peer and their probability of providing positive social reinforcement. As expected due to the simplicity of the task, accuracy was at ceiling and there were no statistical differences in accuracy for the three peers.

In the current study, we observed faster responses to peers who provided positive social reinforcement more often, similar to studies where participants responded more quickly to cues that reliably predicted receiving a primary or secondary reward (O'Doherty et al., 2006; Spicer et al., 2007). Measuring differences in reaction times to cues to index learning differs from reinforcement studies that use modulated choice behavior as an indicator of learning (Tanaka et al., 2004; Daw and Doya, 2006; Schönberg et al., 2007). Choice tasks index changes in explicit preferences or a participant's strategy in maximizing reinforcement, while in the current study, changes in responses are thought to index differences in approach behaviors that are based on learning from a prior history of social feedback. In addition, participants showed faster reaction times after trials that did not

provide positive social reinforcement. This finding is similar to studies that demonstrate improved performance on a trial that follows receiving punishment (Hester et al., 2010) or choosing to make a bet more often after losing money than after winning money (Liu et al., 2007), though the present study did not assess strategic behavior directly. Together, the behavioral findings demonstrate that participants learned the reinforcement contingencies and thus provide an objective index of social learning.

The changes in likeability ratings and response latencies did not appear to be conscious behavioral choices. The majority (93%) of participants were unable to articulate the reinforcement patterns,

suggesting little if any explicit awareness of the reinforcement contingencies. These findings demonstrate that social preferences and actions can be influenced after only brief encounters with peers and without conscious awareness. Such rapid changes highlight the influence of positive social interactions on effectively altering subsequent behavior.

The neural correlates of these behavioral changes draw upon the same neural circuitry as that implicated in reinforcement learning (Alexander et al., 1986; Haber and Knutson, 2010). Prediction error (δ_t) learning engaged the ventral striatum and orbital frontal cortex, similar to previous studies using single-cell recordings (Schultz et al., 1997; Fiorillo et al., 2003; Sul et al., 2010) and human imaging studies with primary reinforcers such as juice (McClure et al., 2003; O'Doherty et al., 2003; D'Ardenne et al., 2008) and secondary reinforcers such as money or attractive or smiling faces (Bray and O'Doherty, 2007; Valentin and O'Doherty, 2009; Lin et al., 2011). Together, these findings support a role for the orbital frontostriatal circuit in generating learning signals from positive social reinforcement and provide a neural basis for how feedback during a social interaction is flexibly updated to inform subsequent social expectations.

The present study is distinct in its capacity to test whether registering violations in expectations of social acceptance draws on basic mechanisms that support prediction error learning. Using a simple Rescorla–Wagner learning model, we show that violations in expected social interaction are tightly coupled with changes in ventral striatal activity. No prior studies, to our knowledge, have applied a classic reinforcement learning model in the examination of learning from social reinforcers. In the social domain, studies have modeled trial-by-trial decisions about charitable donations (Hare et al., 2010) or intentions to trust a partner during economic exchanges (King-Casas et al., 2005). Furthermore, the current paradigm is distinct from previous studies that compare social acceptance to rejection (Eisenberger et al., 2003; Somerville et al., 2006; Guyer et al., 2009), as it targets the process of learning from the social feedback, rather than comparing acceptance to rejection. Therefore, the present study offers a unique explanation for how we learn from positive social interactions.

Our finding that the striatum is sensitive to expectations about receiving social feedback converges with other work targeting the neural mechanisms of social learning. Recently, Harris and Fiske (2010) showed sensitivity in this region to violations in expectations about personality trait information, and others have shown that the striatum is sensitive to violations of social group norms (Klucharev et al., 2009) as well as forming predictions

about investors' decisions (King-Casas et al., 2005; Phan et al., 2010). Our results complement these studies by demonstrating a neural mechanism for how prior positive interactions with others shape our expectations for future interactions. Given the increased sensitivity in the ventral striatum to appetitive stimuli during adolescence (Galvan et al., 2006; Somerville et al., 2011), as well as the greater influence of peers during adolescence (Spear, 2000; Gardner and Steinberg, 2005), this work clearly raises the question of how peer interaction differentially impacts learning and behavior across development and how this may be differentially represented in the brain. Accordingly, it would be interesting to explore whether adolescents show increased sensitivity during social learning relative to children and adults.

The expected values (V_t) to the cues corresponded with greater activity in the rostral anterior cingulate cortex. Previous studies have shown the rostral anterior cingulate cortex/medial prefrontal cortex is sensitive to cues that predict reward receipt (Tanaka et al., 2004; Knutson et al., 2005; Palminteri et al., 2009) and may play a role in general learning about the value of information and using this information for future decisions (Rushworth and Behrens, 2008). Lesion studies in nonhuman primates have shown this region is important for establishing patterns of social interest in other individual male or female macaques (Rudebeck et al., 2006). Human imaging studies have shown this region is sensitive when choosing to approach peers relative to celebrities (Güroğlu et al., 2008) and when engaging in a series of actions during live relative to recorded interactions (Redcay et al., 2010). Given these studies examining social value in the anterior cingulate cortex, and the extensive literature showing a general sensitivity of this region in monitoring response conflict (Botvinick et al., 1999, 2004), our findings suggest that learning social cue values drives changes in behavior that may differ or conflict with the cognitive demands of the situation (e.g., task demands). Over the course of the experiment, this conflict may increase as behavior is modulated in response to changing expected values. Although the current study did not find that the orbital frontostriatal circuit was sensitive to expected values, the findings in the anterior cingulate cortex may suggest a role for this region in processing behavioral tendencies toward learned social cues. These findings thus offer insight into the neural processing of quick social decisions.

Conclusions

Our findings provide direct evidence for how brief, positive, social interactions can significantly shape social learning across three discrete measures: social preferences, behavioral actions, and neural activity. After short interactions with others, social preferences and actions can be altered, highlighting the significance of social acceptance in biasing behavior. Moreover, we show that formal computational models of reinforcement learning apply to secondary reinforcement learning in the social domain. We demonstrate that the neural circuitry involved in forming prediction error signals about receiving social reinforcement, including the ventral striatum and orbital frontal cortex, overlaps with circuitry that subserves learning about other types of rewards (e.g., food or money). Overall, the findings suggest that similar mechanisms underlie basic reinforcement learning and our ability to rapidly and flexibly update our expectations during interactions with others, which enables us to effectively navigate the social environment.

References

- Alexander GE, DeLong MR, Strick PL (1986) Parallel organization of functionally segregated circuits linking basal ganglia and cortex. *Annu Rev Neurosci* 9:357–381.
- Bandura A, Walters RH (1963) *Social learning and personality development*. New York: Holt, Rinehart and Winston.
- Baumeister RF, Leary MR (1995) The need to belong: desire for interpersonal attachments as a fundamental human motivation. *Psychol Bull* 117:497–529.
- Berns GS, McClure SM, Pagnoni G, Montague PR (2001) Predictability modulates human brain response to reward. *J Neurosci* 21:2793–2798.
- Botvinick M, Nystrom LE, Fissell K, Carter CS, Cohen JD (1999) Conflict monitoring versus selection-for-action in anterior cingulate cortex. *Nature* 402:179–181.
- Botvinick MM, Cohen JD, Carter CS (2004) Conflict monitoring and anterior cingulate cortex: an update. *Trends Cogn Sci* 8:539–546.
- Bray S, O'Doherty J (2007) Neural coding of reward-prediction error signals during classical conditioning with attractive faces. *J Neurophysiol* 97:3036–3045.
- Cox RW (1996) AFNI: Software for analysis and visualization of functional magnetic resonance neuroimages. *Comput Biomed Res* 29:162–173.
- D'Ardenne K, McClure SM, Nystrom LE, Cohen JD (2008) BOLD responses reflecting dopaminergic signals in the human ventral tegmental area. *Science* 319:1264–1267.
- Daw ND, Doya K (2006) The computational neurobiology of learning and reward. *Curr Opin Neurobiol* 16:199–204.
- Daw ND, O'Doherty JP, Dayan P, Seymour B, Dolan RJ (2006) Cortical substrates for exploratory decisions in humans. *Nature* 441:876–879.
- Eisenberger NI, Lieberman MD, Williams KD (2003) Does rejection hurt? An fMRI study of social exclusion. *Science* 302:290–292.
- Fiorillo CD, Tobler PN, Schultz W (2003) Discrete coding of reward probability and uncertainty by dopamine neurons. *Science* 299:1898–1902.
- First MB, Spitzer RL, Gibbon M, Williams JB (2007) *Structured Clinical Interview for DSM-IV-TR Axis I Disorders: non-patient edition (SCID-I/NP, 1/2007 revision)*. New York: Biometrics Research, New York State Psychiatric Institute.
- Galvan A, Hare TA, Parra CE, Penn J, Voss H, Glover G, Casey BJ (2006) Earlier development of the accumbens relative to orbitofrontal cortex might underlie risk-taking behavior in adolescents. *J Neurosci* 26:6885–6892.
- Gardner M, Steinberg L (2005) Peer influence on risk taking, risk preference, and risky decision making in adolescence and adulthood: an experimental study. *Dev Psychol* 41:625–635.
- Glover GH, Thomason ME (2004) Improved combination of spiral-in/out images for BOLD fMRI. *Magn Reson Med* 51:863–868.
- Güroğlu B, Haselager GJ, van Lieshout CF, Takashima A, Rijpkema M, Fernández G (2008) Why are friends special? Implementing a social interaction simulation task to probe the neural correlates of friendship. *Neuroimage* 39:903–910.
- Guyer AE, McClure-Tone EB, Shiffrin ND, Pine DS, Nelson EE (2009) Probing the neural correlates of anticipated peer evaluation in adolescence. *Child Dev* 80:1000–1015.
- Haber SN, Knutson B (2010) The reward circuit: linking primate anatomy and human imaging. *Neuropsychopharmacology* 35:4–26.
- Hare TA, Camerer CF, Knoepfle DT, Rangel A (2010) Value computations in ventral medial prefrontal cortex during charitable decision making incorporate input from regions involved in social cognition. *J Neurosci* 30:583–590.
- Harris LT, Fiske ST (2010) Neural regions that underlie reinforcement learning are also active for social expectancy violations. *Soc Neurosci* 5:76–91.
- Hester R, Murphy K, Brown FL, Skilleter AJ (2010) Punishing an error improves learning: the influence of punishment magnitude on error-related neural activity and subsequent learning. *J Neurosci* 30:15600–15607.
- King-Casas B, Tomlin D, Anen C, Camerer CF, Quartz SR, Montague PR (2005) Getting to know you: reputation and trust in a two-person economic exchange. *Science* 308:78–83.
- Klucharev V, Hytönen K, Rijpkema M, Smidts A, Fernández G (2009) Reinforcement learning signal predicts social conformity. *Neuron* 61:140–151.
- Knutson B, Taylor J, Kaufman M, Peterson R, Glover G (2005) Distributed neural representation of expected value. *J Neurosci* 25:4806–4812.

- Li J, Daw ND (2011) Signals in human striatum are appropriate for policy update rather than value prediction. *J Neurosci* 31:5504–5511.
- Li J, Delgado MR, Phelps EA (2011) How instructed knowledge modulates the neural systems of reward learning. *Proc Natl Acad Sci U S A* 108:55–60.
- Lin A, Adolphs R, Rangel A (2011) Social and monetary reward learning engage overlapping neural substrates. *Soc Cogn Affect Neurosci*. Advance online publication. Retrieved March 22, 2011. doi:10.1093/scan/nsr006.
- Liu X, Powell DK, Wang H, Gold BT, Corbly CR, Joseph JE (2007) Functional dissociation in frontal and striatal areas for processing of positive and negative reward information. *J Neurosci* 27:4587–4597.
- McClure SM, Berns GS, Montague PR (2003) Temporal prediction errors in a passive learning task activate human striatum. *Neuron* 38:339–346.
- O'Doherty JP (2007) Lights, camembert, action! The role of human orbitofrontal cortex in encoding stimuli, rewards, and choices. *Ann N Y Acad Sci* 1121:254–272.
- O'Doherty JP, Dayan P, Friston K, Critchley H, Dolan RJ (2003) Temporal difference models and reward-related learning in the human brain. *Neuron* 38:329–337.
- O'Doherty JP, Buchanan TW, Seymour B, Dolan RJ (2006) Predictive neural coding of reward preference involves dissociable responses in human ventral midbrain and ventral striatum. *Neuron* 49:157–166.
- Palminteri S, Boraud T, Lafargue G, Dubois B, Pessiglione M (2009) Brain hemispheres selectively track the expected value of contralateral options. *J Neurosci* 29:13465–13472.
- Phan KL, Sripada CS, Angstadt M, McCabe K (2010) Reputation for reciprocity engages the brain reward center. *Proc Natl Acad Sci U S A* 107:13099–13104.
- Redcay E, Dodell-Feder D, Pearrow MJ, Mavros PL, Kleiner M, Gabrieli JD, Saxe R (2010) Live face-to-face interaction during fMRI: a new tool for social cognitive neuroscience. *Neuroimage* 50:1639–1647.
- Rescorla RA, Wagner AR (1972) A theory of Pavlovian conditioning: variations in the effectiveness of reinforcement and nonreinforcement. In: *Classical conditioning II: current research and theory* (Black AH, Prokasy WF, eds), pp 64–99. New York: Appleton Century Crofts.
- Rudebeck PH, Buckley MJ, Walton ME, Rushworth MF (2006) A role for the macaque anterior cingulate gyrus in social valuation. *Science* 313:1310–1312.
- Rushworth MF, Behrens TE (2008) Choice, uncertainty and value in prefrontal and cingulate cortex. *Nat Neurosci* 11:389–397.
- Schönberg T, Daw ND, Joel D, O'Doherty JP (2007) Reinforcement learning signals in the human striatum distinguish learners from nonlearners during reward-based decision making. *J Neurosci* 27:12860–12867.
- Schultz W, Dayan P, Montague PR (1997) A neural substrate of prediction and reward. *Science* 275:1593–1599.
- Seymour B, O'Doherty JP, Dayan P, Koltzenburg M, Jones AK, Dolan RJ, Friston KJ, Frackowiak RS (2004) Temporal difference models describe higher-order learning in humans. *Nature* 429:664–667.
- Somerville LH, Heatherton TF, Kelley WM (2006) Anterior cingulate cortex responds differentially to expectancy violation and social rejection. *Nat Neurosci* 9:1007–1008.
- Somerville LH, Hare T, Casey BJ (2011) Frontostriatal maturation predicts cognitive control failure to appetitive cues in adolescents. *J Cogn Neurosci* 23:2123–2134.
- Spear LP (2000) The adolescent brain and age-related behavioral manifestations. *Neurosci Biobehav Rev* 24:417–463.
- Spicer J, Galvan A, Hare TA, Voss H, Glover G, Casey B (2007) Sensitivity of the nucleus accumbens to violations in expectation of reward. *Neuroimage* 34:455–461.
- Steinberg L (2008) A social neuroscience perspective on adolescent risk-taking. *Dev Rev* 28:78–106.
- Sul JH, Kim H, Huh N, Lee D, Jung MW (2010) Distinct roles of rodent orbitofrontal and medial prefrontal cortex in decision making. *Neuron* 66:449–460.
- Takahashi YK, Roesch MR, Stalnaker TA, Haney RZ, Calu DJ, Taylor AR, Burke KA, Schoenbaum G (2009) The orbitofrontal cortex and ventral tegmental area are necessary for learning from unexpected outcomes. *Neuron* 62:269–280.
- Talairach J, Tournoux P (1988) Co-planar stereotaxic atlas of the human brain. New York: Thieme Medical Publishers.
- Tanaka SC, Doya K, Okada G, Ueda K, Okamoto Y, Yamawaki S (2004) Prediction of immediate and future rewards differentially recruits cortico-basal ganglia loops. *Nat Neurosci* 7:887–893.
- Valentin VV, O'Doherty JP (2009) Overlapping prediction errors in dorsal striatum during instrumental learning with juice and money reward in the human brain. *J Neurophysiol* 102:3384–3391.