

Neural responses to sanction threats in two-party economic exchange

Jian Li^{a,1}, Erte Xiao^b, Daniel Houser^c, and P. Read Montague^{a,d,2}

^bDepartment of Social and Decision Sciences, Carnegie Mellon University, Pittsburgh, PA 15213; ^cInterdisciplinary Center for Economic Science, Department of Economics, George Mason University, Fairfax, VA 22030; and ^aDepartment of Neuroscience and ^dMenninger Department of Psychiatry & Behavioral Sciences, Baylor College of Medicine, Houston, TX 77030

Communicated by Vernon L. Smith, George Mason University, Fairfax, VA, August 11, 2009 (received for review April 22, 2008)

Sanctions are used ubiquitously to enforce obedience to social norms. However, recent field studies and laboratory experiments have demonstrated that cooperation is sometimes reduced when incentives meant to promote prosocial decisions are added to the environment. Although various explanations for this effect have been suggested, the neural foundations of the effect have not been fully explored. Using a modified trust game, we found that trustees reciprocate relatively less when facing sanction threats, and that the presence of sanctions significantly reduces trustee's brain activities involved in social reward valuation [in the ventromedial prefrontal cortex (VMPFC), lateral orbitofrontal cortex, and amygdala] while it simultaneously increases brain activities in the parietal cortex, which has been implicated in rational decision making. Moreover, we found that neural activity in a trustee's VMPFC area predicts her future level of cooperation under both sanction and no-sanction conditions, and that this predictive activity can be dynamically modulated by the presence of a sanction threat.

cooperation | neuroimaging | perception shift | punishment | social norms

Sanctions are ubiquitous in modern human societies (1). The purpose of sanctions is to enforce norm obedience beyond the level that humans might achieve in the absence of punishment (2–4). However several recent field studies and laboratory experiments have established that adding monetary sanctions to an environment can reduce cooperation (5–7). Substantial speculation has arisen surrounding the source of this counterintuitive effect, including the possibility that the presence of sanctions might change individuals' perceptions of the environment, thus crowding out internal motivations for cooperation (5–8). The imposition of sanctions also might be perceived as a signal of distrust (9–11) and might create a hostile atmosphere (12, 13), leading to decreased cooperation.

Previous behavioral experiments have sought to distinguish these competing explanations. For example, a recent study (5) reported data from an experiment aimed at determining the relative importance of intentions and incentives in producing noncooperative behavior. Participants played a one-shot investment experiment in pairs. Investors sent a certain amount to trustees, requested a return on that investment, and, in some treatments, could threaten sanctions to enforce their requests. Decisions by trustees facing threats imposed (or not) by investors were compared with decisions by trustees facing threats imposed (or not) by nature. The main finding was that when not threatened, trustees typically decided to return a positive amount less than the investor requested, but when threatened, that decision was less common. This result is the same whether the sanction is imposed by a human investor or by nature, suggesting that the detrimental effect of sanctions on cooperation might not hinge specifically on trustees' perceptions of investor intentions. One explanation for such effects has been called the “perception shift” hypothesis, where a nonthreatened subject makes decisions directed by social norms and shifts to utility-driven choices in the presence of threats. In this paper, we pursue the neural substrates of such effects using an economic exchange game

equipped with the possibility that a player can threaten to sanction his or her partner.

The specific brain areas of interest to the perception shift hypothesis are reasonably well established. The parietal cortex has been shown to activate in self-interested economic decision making, especially expected utility calculations (14–16). Neural networks involved in social rewards also have been heavily researched (17–28). Of particular interest to us is the orbitofrontal cortex (OFC), which is known to be reliably involved in social reward evaluation and decision making processes (15, 17, 19, 28–31). But despite the substantial neuropsychology and psychiatry literature pointing to the importance of the prefrontal cortex and the OFC in social recognition and interaction (19, 21–25, 32, 33), ours are among the first experiments informing the OFC's role in perceiving and evaluating threats of sanctions. In particular, we investigate (i) how activation patterns in the OFC depend on whether one is threatened with sanctions and (ii) whether the activity of the medial area of the OFC, the ventromedial prefrontal cortex (VMPFC), a brain area that appears to be pivotal in human decision making (15, 17, 18, 34–38), also predicts subjects' social exchange decisions.

Our study used event-related functional magnetic resonance imaging (fMRI) and an investment game that has been used previously to reliably elicit detrimental sanction effects (5, 9) [Fig. 1; also see [supporting information \(SI\) Fig. S1](#)]. In this game, 2 mutually anonymous participants are paired together for 10 trials. One player is assigned the role of investor and the other is assigned the role of trustee, and both players are given 10 monetary units (MUs) at the beginning of each trial ([Figs. S1 and S2](#)). The subject pairs, as well as the subjects' roles within each pair, remain fixed for the entire 10 rounds. The investor moves first and makes 3 consecutive decisions: (i) the amount of money to send to the trustee (the amount of money was tripled on the way to the trustee), (ii) the amount of money to request back from the trustee, and (iii) whether or not to impose a threat (i.e., a monetary sanction). The sanction is a fixed loss—a 4-MU deduction from the trustee's final earnings should the trustee not send back the requested amount ([Fig. S1](#)). We collected blood oxygen level-dependent (BOLD) images from trustees while they made decisions in the investment game. Investor brain activity was not monitored. Because participants played the game in fixed pairs, reputation presumably could accumulate throughout the experiment. But this presents no difficulties for our analysis, because we focus on sanction–no-sanction contrasts

Author contributions: J.L., E.X., D.H., and P.R.M. designed research; J.L., E.X., D.H., and P.R.M. performed research; J.L. and P.R.M. analyzed data; and J.L., E.X., D.H., and P.R.M. wrote the paper.

The authors declare no conflict of interest.

Freely available online through the PNAS open access option.

¹Present address: Department of Psychology, New York University, New York, NY 10003.

²To whom correspondence should be addressed. E-mail: rmontague@hnl.bcm.edu.

This article contains supporting information online at www.pnas.org/cgi/content/full/0908855106/DCSupplemental.

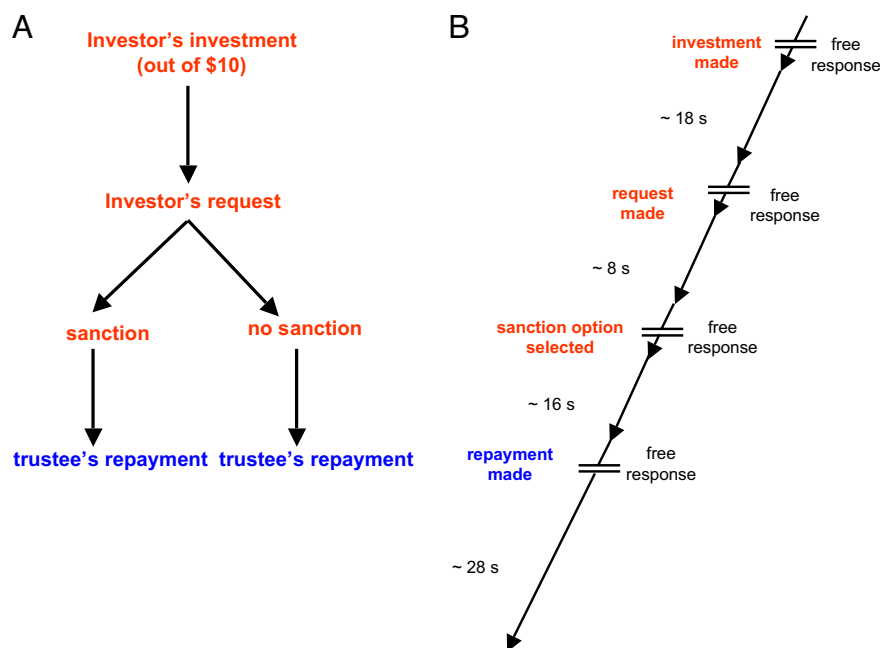


Fig. 1. Experiment task. The task involves 2 subjects sequentially exchanging MUs. Investors' choices are labeled in red; trustees' decisions, in blue. (A) The investor makes 3 decisions sequentially: investment amount, back-transfer request, and whether or not to threaten sanctions. Then the trustee makes the back-transfer decision. (B) Experiment timing. After each player makes her decision, the results are displayed simultaneously to both subjects. A total of 10 rounds are played, and at the end of each round each player's earnings are revealed to both players (also see Figs. S1 and S2).

across all rounds and subjects, thereby controlling for any reputation effects.

Results

Sanction Decisions and Their Effect on Trustees' Repayment Decisions.

On average, investors imposed threats of sanctions 49.3% of the time following a trustee's decision to defect and 46% of the time following a trustee's cooperation. Out of 52 investors, 8 imposed sanctions on every trial, while 11 never imposed a sanction. Overall, an investor's decision to impose a threat was uncorrelated with whether or not a trustee defected in the previous period ($P = .78$; two-sample χ^2 test); however, an investor was more likely to use sanctions in a given trial if (i) the trustee defected in the previous trial and (ii) a sanction had not been used in that previous trial ($\chi^2 = 23.38$; $P = .001$). Overall, investors chose the sanction option 46.3% of the time, ranging

from a high of 53.7% (round 9) to a low of 37.0% (round 1). Using a mixed-effect analysis including a one-sample t test and logistic regression, we found that the correlation between the use of sanctions and the round number did not survive statistical thresholds (average sigmoid slope, 1.64; $P = .053$). Three important variables—investor's investment (mean slope, -0.048 ; $P = .52$), investor's request (mean slope, -0.013 ; $P = .87$), and trustee's repayment (mean slope, -0.03 ; $P = .64$)—are not correlated with round numbers.

To assess trustees' behavioral responses to sanction threats, we first plot an "equal split" strategy as a baseline (Fig. 2B, dotted line). This strategy could emerge if a trustee were to treat the tripled investment amount as a common good and demand half of it. We compare this to trustees' mean real repayments when threatened and when not threatened with sanctions (Fig. 2B, blue and red lines, respectively). Each vertical line in the figures

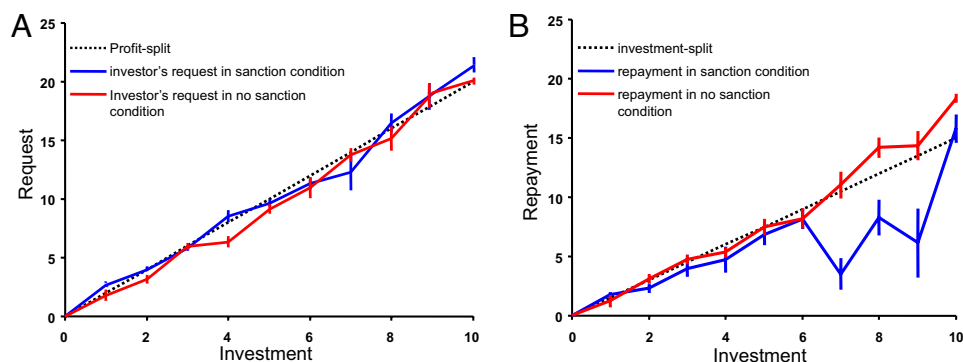


Fig. 2. Summary of players' decisions when sanctions are threatened versus not threatened. Error bars represent SEM. (A) The investor's request as a function of the investment amount. The dotted line indicates a request of two-thirds of the tripled investment amount, which implies equal earnings for investor and trustee. The blue and red curves indicate investors' requests under the threat and no-threat of sanctions condition, respectively. (B) The trustee's repayment as a function of investor's investment. The dotted line indicates a back-transfer amount of half of the tripled investment. The blue and red curves indicate trustee's back-transfer under the threat and no-threat of sanctions condition, respectively (also see Fig. S3).

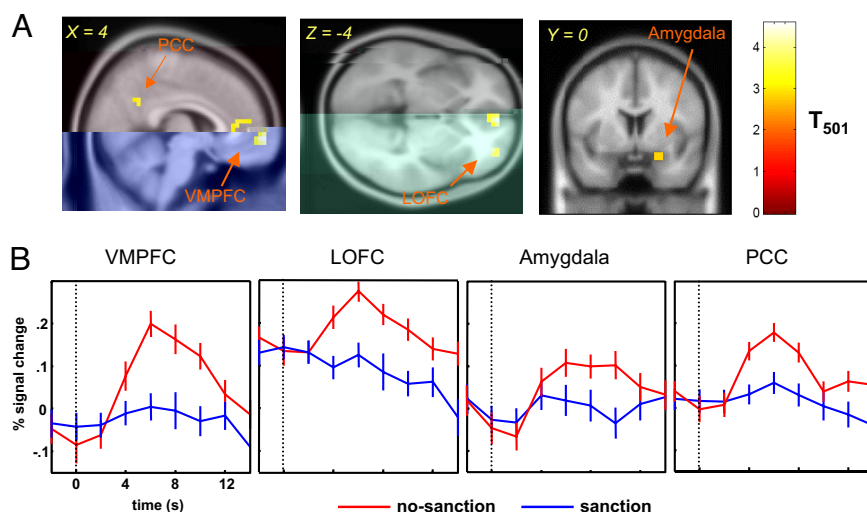


Fig. 3. The trustee's brain regions showing greater activation in the no-sanction condition than in the sanction condition ($P < .001$, uncorrected; cluster size $k > 5$ voxels). (A) A random-effects GLM analysis reveals several brain regions significantly more activated by the revelation of no sanction. These regions include the VMPFC (peak activation MNI coordinate [4 56 -4]), right amygdala (peak activation MNI coordinate [24 0 -20]), right LOFC (peak activation MNI coordinate [32 52 -4]), and PCC (peak activation MNI coordinate [4 -24 36]). (B) Mean event-related time courses of the 4 brain regions. The dashed line indicates the time onset; error bars are SEM. Bold signals in the VMPFC, LOFC, amygdala, and PCC are all significantly greater when the trustee is in the no-sanction condition (red traces) than when she is in the sanction condition (blue traces).

represents 1 SE of the trustees' mean repayment in both conditions. The trustee's repayment when threatened with sanctions is significantly different between sanction and no-sanction cases ($P < .05$; two-sample t test); see Fig. S3 and Table S1 for details. The difference is greater when the investments are larger (>6). Overall, trustees' average repayments are 6.05 MUs in sanction cases and 12.04 MUs in no-sanction cases (Table S1). Thus, the difference in repayment amounts cannot be explained solely by the possibility that trustees choose to keep 4 extra MUs in the sanction condition as compensation for the cost of the sanction.

Previous research suggests that trustees' repayments also might depend on whether the investor used the sanction to enforce an "unfair" request (5) (defined as a request for 2/3 of the tripled investment amount, which is the amount that equalizes investor and trustee earnings). To investigate unfair requests, we first explored investor behavior by plotting the back-transfer request against the investment decision for both the sanction and no-sanction conditions (Fig. 2A, blue and red lines, respectively). The dotted line in that figure indicates a request of 2/3 of the tripled investment. It is apparent that the investors' requests do not differ significantly between the sanction and no-sanction conditions ($P = .9$; t test), nor are the averages significantly different on average from equal-earnings requests ($P_{\text{no-sanction}} = .9$; $P_{\text{sanction}} = .9$). With respect to trustees' decisions, consistent with previous studies (5), we find that sanctions have a detrimental effect on trustees' returns both when the investor's back-transfer request is fair and when it is unfair, and that these detrimental effects are not statistically significantly different. In particular, a fair request results in a mean return equal to 53% of the tripled investment amount, while combining sanctions with a fair request reduces returns to 47% on average. When the request is unfair, the analogous change is from 59% to 47%; this between-condition difference (a 6% vs 12% reduction) is not statistically significant ($P > .15$, two-tailed Wilcoxon test). Previous reports suggest that subjects in repeated games might adopt sophisticated Nash equilibrium strategies (39,40), and we specifically tested those hypotheses (see *SI Text* for more details).

Trustees' Neural Responses to the Revelation of Sanctions. To gain insight into the neural underpinnings of this effect, we used a standard general linear model analysis (GLM) to compare trustees' brain responses in cases where sanctions were and were not threatened by the investor. The sanction–no-sanction contrast did not identify any prefrontal brain activities at $P < .001$ (uncorrected, 5 continuous voxels; see [Table S3](#)), but the no-sanction—sanction contrast did reveal differential activation in areas implicated in social reward processing (Fig. 3; [Table S2](#)). These brain areas include the VMPFC (peak activity at MNI [4 56 –4]), lateral OFC (LOFC; peak activity at MNI [32 52 –4]), posterior cingulate cortex (PCC peak activity at MNI [4 –24 36]), and right amygdala (peak activity at MNI [24 0 –20]). We conducted a region-of-interest (ROI) analysis to further investigate these results (Fig. 3B). In the figure, the vertical dotted line indicates the point at which either the sanction or no-sanction screen was revealed, and the red and blue curves represent brain activities in the no-sanction and sanction conditions, respectively (23–25, 31, 35).

Brain Activity Predicts Trustee's Repayment. We used standard parametric regression analysis to explore whether a trustee's neural activity at the revelation of the sanction screen might predict her subsequent back-transfer decision (which was made about 10 or 15 seconds later). Because the absolute back-transfer from a trustee does not inform a trustee's intention to cooperate, it is sensible to normalize the back-transfer by the maximum amount that the trustee could have sent (i.e., the tripled investment amount). The back-transfer-to-tripled-transfer amount ratio is a useful measure of a trustee's willingness to cooperate. Our analysis revealed a brain area at the superior frontal gyrus (DLPFC) (peak activity at MNI [24 52 20]; $P < .005$, uncorrected) (Fig. 4A). The activity of this area is negatively correlated with the back-transfer-to-investment amount ratio. Further ROI analysis demonstrated that as this back-transfer ratio increases, the BOLD signal at the DLPFC area decreases, and it returns to the baseline level when the trustee cooperates fully (Fig. 4A, Bottom; each vertical bar represents 1 SEM). Positive parametric regression analysis identified several brain areas, including the medial frontal gyrus (38), the inferior frontal

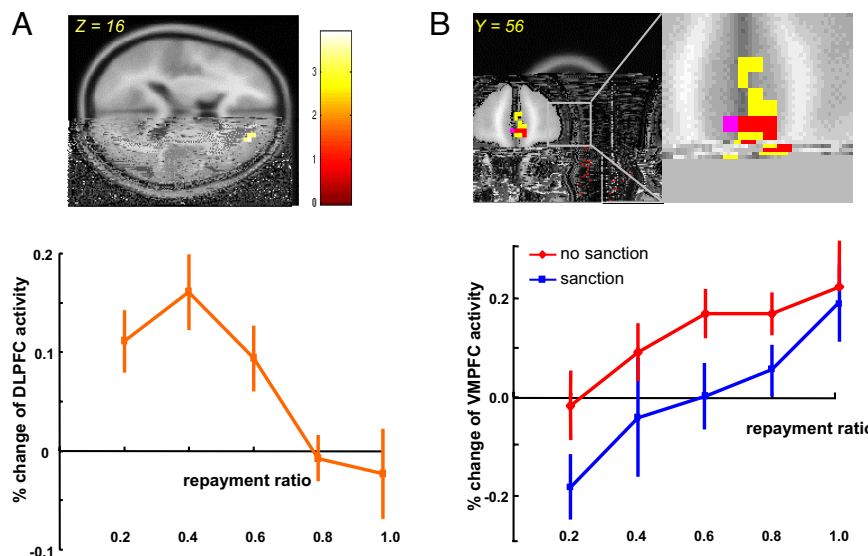


Fig. 4. Trustees' brain regions whose activations are parametrically correlated with trustees' normalized back-transfer (defined as the ratio of the back-transfer and the tripled investment amount). (A) Brain activity at dorsal lateral prefrontal cortex (DLPFC; peak activation MNI coordinate [24 52 20]) is negatively correlated with trustees' normalized back-transfers ($P < .001$, uncorrected; cluster size, $k > 5$ voxels). (B) A GLM ($P < .005$, uncorrected; cluster size, $k > 5$ voxels) showing that a subset of voxels (peak activation MNI coordinate [-4 56 -4]; purple) in the VMPFC area (yellow, with the overlap in orange) previously identified in Fig. 3A strongly and positively predicts trustees' normalized back-transfers. Further ROI analysis indicates that the VMPFC activity is correlated with trustees' normalized back-transfers in both sanction and no-sanction conditions. The slopes of the 2 curves (red and blue) do not differ significantly ($P = .1$, t test) while the intercept of the no-sanction curve (red) is significantly greater than that of the sanction curve (blue; $P < .01$, t test).

cortex, the middle temporal cortex, and the occipital cortex (Fig. 4B and Table S4; $P < .005$, uncorrected). Interestingly, one of those brain areas, the area in the VMPFC (peak activity at MNI [-4 56 -4]; Fig. 4B, purple) overlaps significantly with the VMPFC region identified in the previous sanction–no-sanction contrast (Fig. 4B, yellow; overlapping area indicated in orange).

The ROI analysis (Fig. 4B, Bottom) demonstrates this unique pattern of VMPFC activation. Although the VMPFC activity correlates with the repayment ratio in general, further separation of the VMPFC BOLD signal into sanction and no-sanction categories reveals a shift of the BOLD signal in both conditions (Fig. 4B; sanction in blue, no-sanction in red). Moreover, there is only weak evidence of differing slope coefficients ($P = .1$, two-sample t test); the intercepts are significantly different, however ($P < .01$, t test). It is also interesting to note that when the trustee plans to completely defect in the no-sanction situation, VMPFC activity remains at baseline, but when the trustee plans to defect under the sanction condition, VMPFC activity is well below baseline ($P < .05$, t test). The fact that brain activity at the VMPFC precedes the trustee's actual repayment choice by 10–15 seconds suggests that this brain area might be heavily involved in the trustee's final decision making and might generate a BOLD signal predicting the trustee's repayment ratio. This signal is thus responsive, in that it is susceptible to social cues (i.e., whether or not the trustee is threatened by sanctions), and also acts as a predictive signal, parametrically modulating the trustee's final repayment.

Discussion

Using an iterated version of the trust game with a sanction component, we have demonstrated an aversive effect of sanctions on human cooperation as measured by trustee's repayment in the investment game (5) (Fig. 2B). Recent theories that incorporate other preferences (particularly inequality aversion and kindness) shed light on motives for trustees' decisions in standard trust games (6, 41–49) but cannot explain the detrimental effect of punishment on reciprocity. We hypothesized

that this effect might be due to a “perception shift” from norm-sensitive choices to utility-based choices.

Differential Brain Activities in the Sanction–No-Sanction Contrast.

Our perception shift hypothesis suggests that trustees not threatened with sanctions make their reciprocity decision within a social context and are directed by social norms. Indeed, we found that when a trustee learns that he or she has not been threatened with sanctions, a neural network including the VMPFC, right amygdala, LOFC, and PCC is activated. Activation of these reward-related pathways supports our hypothesis for several reasons. Recent studies have found elevated brain activity in the LOFC area when subjects choose to comply with social norms (50, 52), while the medial part of the OFC (VMPFC) may be involved in preference generation and final decision making (17, 30, 33, 52–54). Although amygdala activation in humans has been associated with negative emotions and fear conditioning, emerging evidence suggests that the amygdala might be equally important to reward processing (22, 52, 53, 55–59). In addition, reciprocal connections between the amygdala and OFC have been studied extensively, and the functional interaction between these two regions is thought to be essential in goal-directed behaviors (53, 54, 56–59).

Differential Brain Activation in the Sanction–No-Sanction Contrast.

The sanction–no-sanction contrast did not reveal any differential brain responses in the prefrontal cortex. Instead, we observed bilateral parietal cortex (LIP) activation (Table S3). Parietal activity has been linked to the representation of expected utility in primate research and “rational” choices in both primates and humans (16, 60). Our finding of no differential activation of social or emotional systems under sanction threats seems to cast some doubt on the role of negative “intentions” in affecting behavior in this environment. Instead, this finding provides convergent support for the “cognitive shift” hypothesis that credible threats of sanctions generate a cognitive shift that diminishes social motivations and increases the likelihood of market-oriented earnings maximizing behavior (5–8).

49. Ellingsen T, Johannesson M (2008) Pride and prejudice: The human side of incentive theory. *Am Econ Rev* 98:990–1008.
50. Montague PR, Lohrenz T (2007) To detect and correct: Norm violations and their enforcement. *Neuron* 56:14–18.
51. Spitzer M, Fischbacher U, Herrnberger B, Gron G, Fehr E (2007) The neural signature of social norm compliance. *Neuron* 56:185–196.
52. Gottfried JA, O'Doherty J, Dolan RJ (2003) Encoding predictive reward value in human amygdala and orbitofrontal cortex. *Science* 301:1104–1107.
53. Winstanley CA, Theobald DE, Cardinal RN, Robbins TW (2004) Contrasting roles of basolateral amygdala and orbitofrontal cortex in impulsive choice. *J Neurosci* 24:4718–4722.
54. Arana FS, et al. (2003) Dissociable contributions of the human amygdala and orbitofrontal cortex to incentive motivation and goal selection. *J Neurosci* 23:9632–9638.
55. Schultz W (2000) Multiple reward signals in the brain. *Nat Rev Neurosci* 1:199–207.
56. Baxter MG, Parker A, Lindner CC, Izquierdo AD, Murray EA (2000) Control of response selection by reinforcer value requires interaction of amygdala and orbital prefrontal cortex. *J Neurosci* 20:4311–4319.
57. Moll J, et al. (2002) The neural correlates of moral sensitivity: A functional magnetic resonance imaging investigation of basic and moral emotions. *J Neurosci* 22:2730–2736.
58. Holland PC, Gallagher M (2004) Amygdala–frontal interactions and reward expectancy. *Curr Opin Neurobiol* 14:148–155.
59. Schoenbaum G, Chiba AA, Gallagher M (1998) Orbitofrontal cortex and basolateral amygdala encode expected outcomes during learning. *Nat Neurosci* 1:155–159.
60. McClure SM, Laibson DI, Loewenstein G, Cohen JD (2004) Separate neural systems value immediate and delayed monetary rewards. *Science* 306:503–507.
61. Hare TA, Camerer CF, Rangel A (2009) Self-control in decision-making involves modulation of the vmPFC valuation system. *Science* 324:46–48.

Supporting Information

Li et al. 10.1073/pnas.0908855106

SI Text

Image and Statistical Analysis

Image Acquisition and Preprocessing. High-resolution T1-weighted scans ($1 \times 1 \times 1$ mm) were acquired on Siemens 3T Allegra scanners using a Siemens MRPage sequence. Functional image details are as follows: echo-planar imaging, repetition time (TR) = 2000 ms; echo time (TE) = 40 ms; flip angle = 90 degrees; 64×64 matrix with 26 4-mm-thick axial slices, yielding functional $3.4 \times 3.4 \times 4$ mm³ voxels. To optimize functional sensitivity in the OFC, we acquired images using an oblique 30-degree angle to the AC–PC axis. All of the imaging data were processed and analyzed using SPM2 (<http://www.fil.ion.ucl.ac.uk/spm/software/spm2>) and xjView (<http://people.hnl.bc-m.tmc.edu/cuixu/xjView>). Functional images were realigned using a 6-parameter rigid-body transformation. Each individual's structural T1 image was co-registered to the average of the motion-corrected images using 12-parameter affine transformation. Individual T1 structural images were segmented into gray matter, white matter, and cerebrospinal fluid before the individual gray matter was nonlinearly warped into an MNI gray matter template. Functional images were then slice-timing artifact-corrected and normalized into MNI space by applying the transformation matrix adopted from previous T1 warping. The images were then smoothed with an 8-mm isotropic Gaussian kernel and high-pass-filtered in the temporal domain (filter width, 128 s).

Statistical Analysis. For GLM analysis, functional images were divided into separate rounds (10 rounds) that included all images preceding each round by 20 s and following the end of each round by 8 s. A separate GLM was specified and estimated for each round of the task for each subject. All visual cues and motor responses were constructed and estimated independently for each subject by convolving a delta function at the onset of those events with a canonical hemodynamic response function implemented within SPM2. The random-effects analyses depicted in Fig. 3 and Tables S2 and S3 were performed as follows. A fixed-effect analysis was performed for each round for each subject to estimate the brain activity of effects of interest. Beta images generated from previous analyses were further separated into 2 uneven groups of 281 (no-sanction condition) and 238 (sanction condition) contrast images of a single between-group factor (sanction or no-sanction), and a 2-sample *t* test was performed. Table S2 identifies brain regions with significantly greater activity ($T_{517} = 3.11$; $P < .001$, uncorrected) in response to the no-sanction screen relative to the sanction screen. Table S3 identifies brain regions with significant greater activity ($T_{517} = 3.11$; $P < .001$, uncorrected) in response to the sanction screen versus the no-sanction screen.

ROI Analysis. ROI analyses for 4 brain regions shown in Fig. 3 (VMPFC, LOFC, amygdala, and DLPFC) were performed on the 5 most significantly activated voxels from the *t* test, as listed in Table S2. The spatially averaged signal was linearly detrended within each round and time-locked to the display of “sanction/no-sanction” information to the trustee's brain. The correlation between brain activity of the VMPFC and DLPFC and the normalized repayment ratio illustrated in Fig. 4 are based on averages grouped by level of normalized repayment ratio (amount of repayment/ $3 \times$ investment), binned into 5 normalized repayment ratio levels: [0–0.2), [0.2–0.4), [0.4–0.6), [0.6–

0.8), and [0.8–1.0]. Trial event numbers for the 5 repayment ratio levels are [23, 28, 75, 146, 9] for the no sanction condition and [75, 18, 36, 84, 30] for the sanction condition. Brain activities at the VMPFC and DLPFC shown in Fig. 4 are the averages of peak hemodynamic activities (at 4–6 s after event onset) and 2 data points surrounding the peak.

Nash Equilibrium Strategies

Investigating Whether Subjects Use Sophisticated Nash Equilibrium Strategies. Although we repeated our game 10 times, we first derive standard Nash equilibrium (NE) predictions based on selfish preferences for the one-shot game. In this environment, trustees should not return any amount if the investor does not impose a sanction threat. Consequently, the investor should send nothing, meaning that both would earn their endowment of 10 MUs. But threatening a sanction of 4 MUs can enforce a back-transfer request of at most 4; thus, in this case an NE occurs when an investor sends 1 (or 2) MU, requests a back-transfer of 3 (or 4) MUs, and threatens a sanction of 4. The trustee then returns 3 (or 4) MUs to the investor. In both cases, the investor earns 12 MUs. The trustee earns 10 MUs when the investor sends 1 and 12 MUs when the investor sends 2. Thus, there are multiple Nash equilibria for the one-shot game, and trustees are predicted to return more under a punishment threat (return 3 or 4) than when punishment is not threatened (return 0).

One NE for the repeated game involves playing any one-shot equilibrium in every round. But previous have found that find investors send more, and trustees subsequently return more, than would be predicted by such a “naïve” equilibrium (1, 2). One possible explanation for this is that fully rational participants play a “sophisticated” subgame-perfect NE strategy that leads to greater amounts sent and returned than are predicted by one-shot equilibria. Such “as if” cooperative equilibria can exist when the one-shot game admits 2 or more NEs, and when one of these NEs is “worse” for the players than the other(s) (3, 4).

In our game, the NE payoff of [12,10] is worse for a trustee than [12,12], which can lead to “as if” cooperative behavior. To see how this occurs, suppose that our game were repeated only twice instead of 10 times. Then, for example, in the first round the investor could send 3, ask for a return of 6, and threaten to punish. If the trustee reciprocated and returned the full 6 (implying that both earned 13 in the first round) then the investor would play the “nice” NE in the final round, meaning that the trustee would earn 12, for a total payoff of 25. If the trustee instead defected in the first round, then the investor would play the [12,10] equilibrium in the final round. It follows that the trustee would have no incentive to defect in the first round, and “as if” cooperative equilibrium behavior would follow. By similar reasoning, we can see that “as if” cooperation can be supported as a sophisticated NE of our 10-round game. But, as we demonstrate later, the cooperative patterns found in our data are inconsistent with this explanation.

We note that investors send, and trustees return, substantial amounts. To shed light on whether this seemingly cooperative behavior might stem from sophisticated noncooperative NE strategies, observe that a one-shot NE *must* be played in the last period of any sophisticated equilibrium path (4). The average amount invested in the final round of our game is 5.9. This amount is statistically significantly larger than 2 ($P < .001$, 2-sided *t* test), which is the largest possible equilibrium investment amount in the one-shot game. Indeed, the vast majority (>67%) of investors send more than 2 in the final round, and of

those who send 2 or fewer, 41% make a punishment or back-transfer request decision that is inconsistent with a one-shot equilibrium. Thus, only about 15% of investors in the final round of our game make decisions consistent with sophisticated NE. Moreover, the average amount returned by trustees in the final round is 9.7. This amount is statistically significantly larger than 4 ($P < .001$, two-sided t test), which is the maximum return consistent with equilibrium in the stage game.

Finally, also note that sophisticated equilibria can involve “trigger strategies,” under which an investor reverts to a one-shot NE following a defection. But in our data, the average investment following a defection (returning less than the investor requested) is 5.4, which is more than half of the endowment and again is statistically significantly larger than 2 ($P < .001$, two-sided t test.)

In light of our evidence, we conclude that sophisticated NE play is not a plausible explanation for the cooperative patterns found in our data.

Exploring the “Once Commodity, Always Commodity” Hypothesis.

Previous research suggests the “once commodity, always commodity” (OCAC) hypothesis that perception shifts can persist even when the source of the shift is removed (5). We investigated, both behaviorally and at the neural level, whether exposure to a sanction creates a perception shift that persists in future exchanges that do not include a sanction. To do this, we focused on the 33 pairs in which the investor chose both sanction and no-sanction at least once during the 10 rounds. We categorized each round of each pair in 1 of 3 mutually exclusive ways: (i) nonsanction trials before the investors imposed sanctions for the first time, (ii) sanction trials, and (iii) nonsanction trials experienced after sanction trials. We obtained a total of 88 obser-

vations on 20 unique subjects in group 1, 163 observations from 33 subjects in group 2, and 83 observations from 26 subjects in group 3. The OCAC hypothesis predicts that groups 2 and 3 should exhibit similar return behavior, and that group 1 should return more on average than both the others.

We find that trustees’ returns (measured as percentage of tripled investment) are higher in group 1 (49.2%) than in group 2 (42.6%), and that the decrease is (marginally) significant ($P = .10$, 2-tailed t test). However, the back-transfer rate in group 2 (42.6%) is less than that in group 3 (51.6%), and the difference is statistically significant ($P = .03$, two-tailed t test). This is inconsistent with the OCAC hypothesis, and thus we do not find behavioral evidence supporting this hypothesis in our environment.

To explore the OCAC hypothesis at the neural level, we conducted our imaging analysis using only the restricted sample of 33 subjects and only those observations occurring in either group 2 (sanction observations) or group 3 (the no-sanction trials that occur subsequent to a sanction trial). The OCAC hypothesis would predict neural activations consistent with “market” decision making in both groups. In fact, however, for the no-sanction–sanction contrast, we found results again supporting social reward systems (VMPFC, PCC, OFC, and amygdala), but at a lower threshold; due to the substantially reduced sample size, we used $P = .01$. Fig. S4 reports the results of this analysis and shows that the activations are closely related to those revealed when using the full sample. Similarly, we investigated the sanction–no-sanction contrast with the restricted sample. Again at a lower threshold ($P = .01$), we found activations in lateral inferior parietal cortex (LIP) that line up well with our original findings. It follows that neither our behavioral nor neural evidence supports the OCAC hypothesis.

1. Houser D, Xiao E, McCabe K, Smith V (2008) When punishment fails: Research on sanctions, intentions and non-cooperation. *Games Econ Behav* 62:509–532.
2. Fehr E, Rockenbach B (2003) Detrimental effects of sanctions on human altruism. *Nature* 422:137–140.
3. Friedman J (1985) Cooperative equilibria in finite-horizon noncooperative supergames. *J Econ Theor* 35:390–398.

4. Benoit JP, Krishna V (1985) Finitely repeated games. *Econometrica* 53:905–922.
5. Andreoni J, Harbaugh W, Vesterlund L (2003) The carrot or the stick: Rewards, punishments, and cooperation. *Am Econ Rev* 93:893–902.

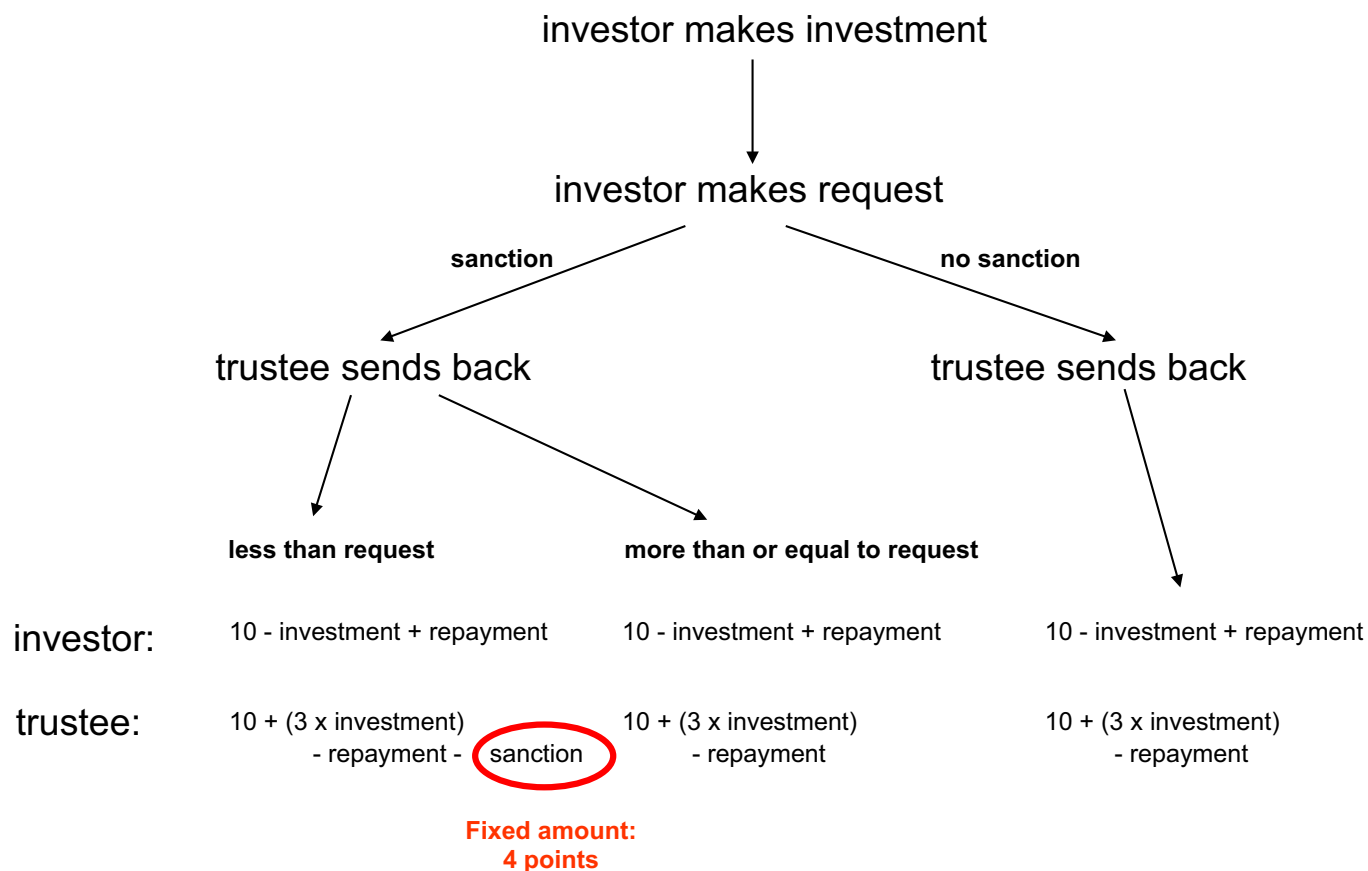


Fig. S1. The 2-player investment game. Two players are paired anonymously. Both the investor and the trustee are endowed with 10 points at the beginning of each round of the experiment (10 rounds total). The investor first decides how many points to invest, how many to request back and, whether or not to threaten punishment. The trustee observes these 3 pieces of information and then decides how many points to send back to the investor. If the trustee returns less than the investor requested, and if the investor chose the threat option, then a penalty of 4 points is deducted from the trustee's final earnings. If the threat was not chosen, then the trustee's and the investor's earnings depend only on the amounts sent and returned, respectively, as described above.

structure of an exchange

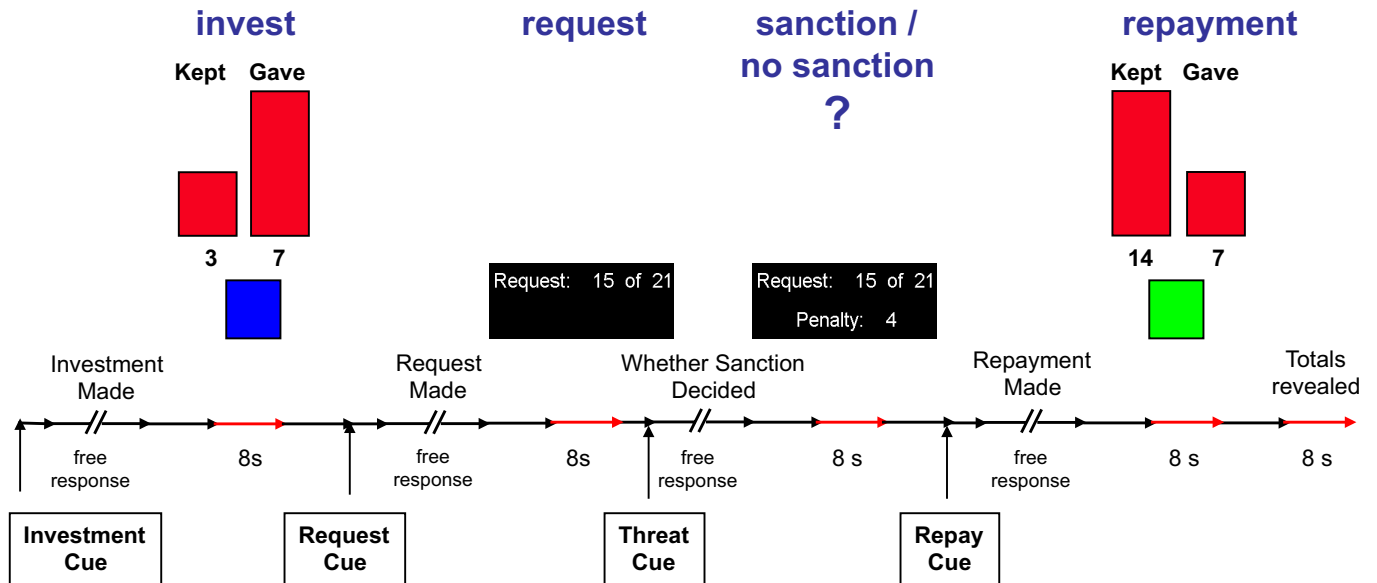
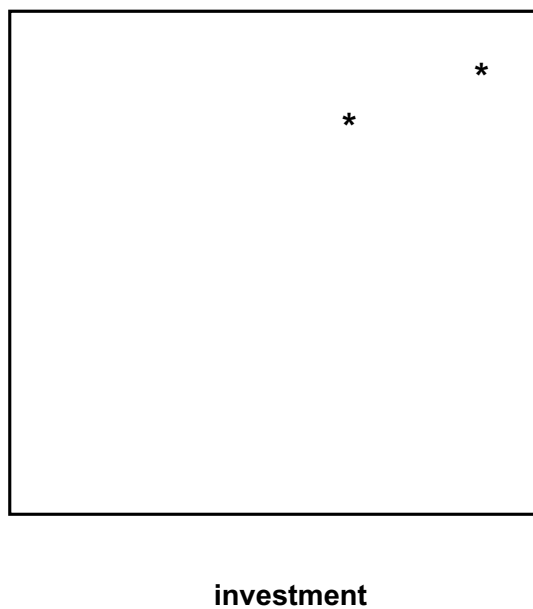


Fig. S2. Time line for the 2-player investment game. Each pair of subjects completed 10 consecutive exchanges. Each exchange began with a screen that indicated the beginning of the round, followed by a cue to invest. The investor then invested between 0 and 10 monetary units. After the investor's decision, the investment was displayed to both subjects for 8 seconds. The timing of the investor's next 2 decisions—the back-transfer request and whether or not to threaten a sanction—proceeded in an identical manner. After the investor completed 3 decisions, the trustee was prompted to return an amount (between 0 and triple the investment amount) back to the investor. The trustee's decision was revealed to both subjects for 8 seconds, followed by 8 seconds of a blank screen. That round's total earnings for both subjects was then displayed. Rounds were separated by a variable 12- to 42-second interval.



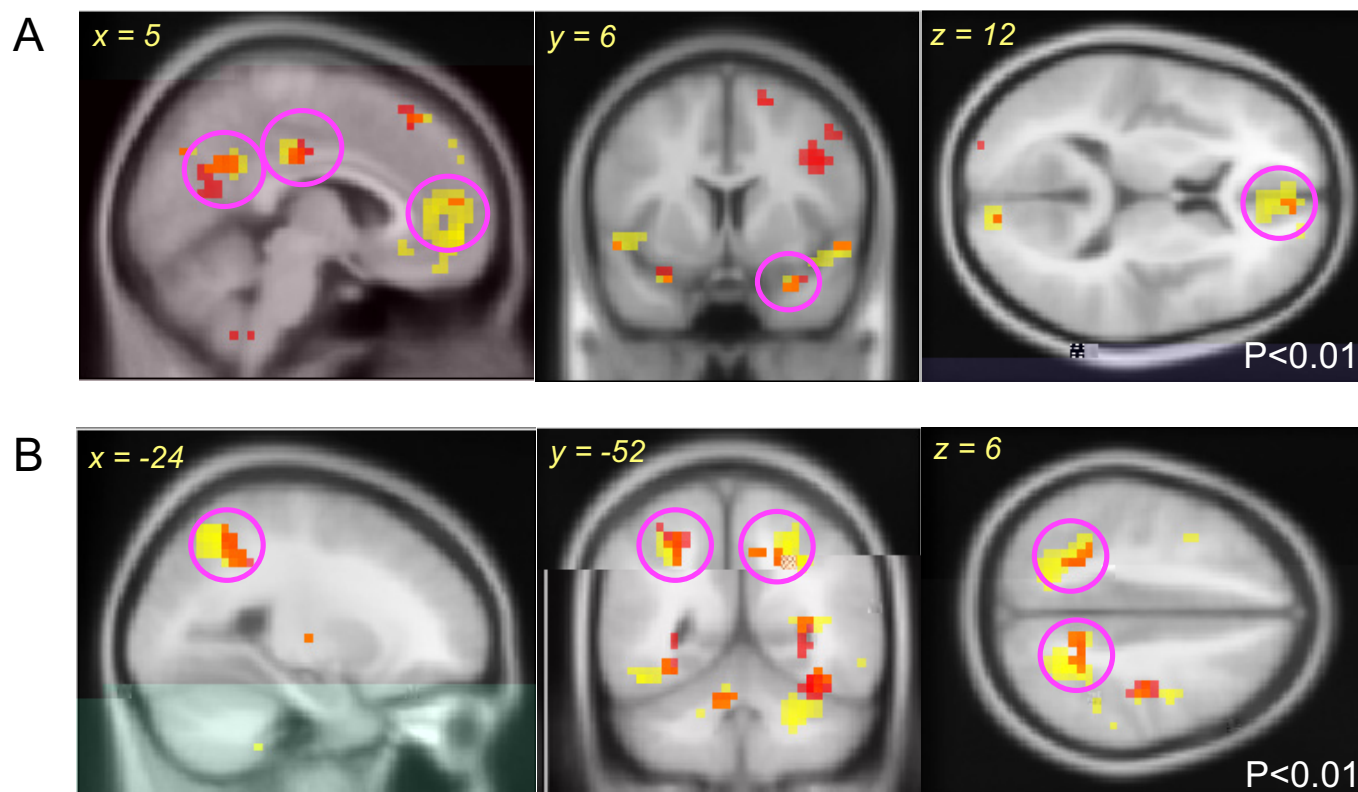


Fig. S4. Direct comparison between brain activation revealed by the full data set (52 subjects) and the OCOC hypothesis (33 subjects). (A) Brain areas, including the PCC, VMPFC, and amygdala (magenta circles), as revealed by no-sanction–sanction contrast, are represented in the overlapping pattern (orange) from the full 52-subject data set (yellow) and the restricted 33-subject data set (red). (B) Bilateral parietal cortex (magenta circles), as revealed by the no-sanction–sanction contrast, are represented in the overlapping pattern (orange) from the full (yellow) and restricted (red) data sets.

Table S1. Average behavior and payoff of investors and trustees

	Sanction	No sanction	Significance
Investment	4.89	7.09	*
Request	10.6	13.89	-
Request/(3 × investment)	0.72	0.64	*
Repayment	6.05	12.04	-
Repayment/(3 × investment)	0.46	0.55	*
Repayment/request	0.67	0.89	*
Investor's payoff	11.58	14.95	*
Trustee's payoff	17.01	19.22	-

*Statistical significance

Table S2. Brain responses differentially activated in no-sanction versus sanction conditions

Region of activation	Peak MNI coordinates			Voxels	z-value
	X	Y	Z		
Medial frontal gyrus (R)	4	56	−4	83	4.45
Superior temporal gyrus (R)	48	16	−16	52	4.52
Superior temporal gyrus (L)	−48	16	−12	31	3.76
Lateral frontal gyrus (R)	32	52	−4	15	4.03
Superior frontal gyrus (R)	20	40	48	35	3.78
Superior frontal gyrus (L)	−28	40	36	24	3.26
Occipital lobe (R)	12	−92	12	12	3.07
Occipital lobe (L)	−16	−88	−16	19	3.58
Precuneus (R)	4	−52	32	12	3.49
Posterior cingulate cortex	4	−24	36	11	3.41
Inferior frontal gyrus (R)	52	24	4	5	2.78
Amygdala (R)	24	0	−20	7	2.7

Regions with ≥ 5 significant voxels were identified at $P < .005$ (uncorrected).

Table S3. Brain responses differentially activated in sanction versus no-sanction conditions

Region of activation	Peak MNI coordinates			Voxels	z-value
	X	Y	Z		
Parietal lobe (L)	−24	−60	52	72	3.99
Parietal lobe (R)	28	−48	40	81	4.13
Inferior temporal gyrus	−44	−68	−4	67	4.1
Temporal lobe	28	−68	20	27	3.29
Precentral gyrus (R)	44	−4	36	68	3.97
Precentral gyrus (L)	−44	−8	36	80	3.79
Fusiform gyrus (R)	36	−48	−16	18	3.63
Medial frontal gyrus	−8	−24	68	17	3.3
Midbrain	4	−12	−12	59	4.17
Cerebellum	24	−48	−36	44	4.19

Regions with ≥ 5 significant voxels were identified at $P < .005$ (uncorrected).

Table S4. Brain responses positively correlated with repay ratio by trustees

Region of activation	Peak MNI coordinates			Voxels	z-value
	X	Y	Z		
Medial frontal gyrus	−4	56	−4	6	2.84
Inferior frontal gyrus	36	16	−20	18	3.89
Middle temporal gyrus	−60	−60	8	5	3.42
Temporal lobe	−52	−8	−28	7	3.4
Occipital lobe	−16	−96	−8	9	3.37

Regions with ≥ 5 significant voxels were thresholded at $P < .005$ (uncorrected).