

Computing Social Value Conversion in the Human Brain

 Haruaki Fukuda,^{1,2,3*}  Ning Ma,^{1,2*}  Shinsuke Suzuki,^{1,4,5}  Kenichi Ueno,^{2,6}
 Justin L. Gardner,^{7,8}  Noritaka Ichinohe,⁹  Masahiko Haruno,¹⁰  Kang Cheng,^{6,11†} and  Hiroyuki Nakahara^{1,2,12}

¹Laboratory for Integrated Theoretical Neuroscience, RIKEN Brain Science Institute, Wako, Saitama 351-0198, Japan, ²RIKEN Center for Brain Science, Wako, Saitama 351-0198, Japan, ³Department of General System Studies, The University of Tokyo, Meguro, Tokyo 153-8902, Japan, ⁴Frontier Research Institute for Interdisciplinary Sciences, Tohoku University, Sendai, Miyagi 980-8578, Japan, ⁵Institute of Development, Aging and Cancer, Tohoku University, Sendai, Miyagi 980-8575, Japan, ⁶Support Unit for Functional Magnetic Resonance Imaging, RIKEN Brain Science Institute, Wako, Saitama 351-0198, Japan, ⁷Department of Psychology, Stanford University, Stanford, California 94305, ⁸Laboratory for Human Systems Neuroscience, RIKEN Brain Science Institute, Wako, Saitama 351-0198, Japan, ⁹Department of Ultrastructural Research, National Institute of Neuroscience, National Center of Neurology and Psychiatry, Kodaira, Tokyo 187-8502, Japan, ¹⁰Center for Information and Neural Networks, National Institute of Information and Communication Technology, Suita, Osaka 565-0871, Japan, ¹¹Laboratory for Cognitive Brain Mapping, RIKEN Brain Science Institute, Wako, Saitama 351-0198, Japan, and ¹²Department of Intelligence Science and Technology, Kyoto University, Kyoto, Kyoto 606-8501, Japan

Social signals play powerful roles in shaping self-oriented reward valuation and decision making. These signals activate social and valuation/decision areas, but the core computation for their integration into the self-oriented decision machinery remains unclear. Here, we study how a fundamental social signal, social value (others' reward value), is converted into self-oriented decision making in the human brain. Using behavioral analysis, modeling, and neuroimaging, we show three-stage processing of social value conversion from the offer to the effective value and then to the final decision value. First, a value of others' bonus on offer, called offered value, was encoded uniquely in the right temporoparietal junction (rTPJ) and also in the left dorsolateral prefrontal cortex (ldlPFC), which is commonly activated by offered self-bonus value. The effective value, an intermediate value representing the effective influence of the offer on the decision, was represented in the right anterior insula (rAI), and the final decision value was encoded in the medial prefrontal cortex (mPFC). Second, using psychophysiological interaction and dynamic causal modeling analyses, we demonstrated three-stage feedforward processing from the rTPJ and ldlPFC to the rAI and then from rAI to the mPFC. Further, we showed that these characteristics of social conversion underlie distinct sociobehavioral phenotypes. We demonstrate that the variability in the conversion underlies the difference between prosocial and selfish subjects, as seen from the differential strength of the rAI and ldlPFC coupling to the mPFC responses, respectively. Together, these findings identified fundamental neural computation processes for social value conversion underlying complex social decision making behaviors.

Key words: computational; decision making; fMRI; social preference; social value orientation; value

Significance Statement

In daily life, we make decisions based on self-interest, but also in consideration for others' status. These social influences modulate valuation and decision signals in the brain, suggesting a fundamental process called value conversion that translates social information into self-referenced decisions. However, little is known about the conversion process and its underlying brain mechanisms. We investigated value conversion using human fMRI with computational modeling and found three essential stages in a progressive brain circuit from social to empathic and decision areas. Interestingly, the brain mechanism of conversion differed between prosocial and individualistic subjects. These findings reveal how the brain processes and merges social information into the elemental flow of self-interested decision making.

Introduction

The reinforcement learning framework has provided a foundation for understanding the computations and associated neural

mechanisms involved in self-regarding valuation (Schultz et al., 1997; Daw and Doya, 2006; Dayan and Nakahara, 2018). It provides a simple, rigorous account of choice behavior that maxi-

Received Dec. 5, 2018; revised March 30, 2019; accepted April 14, 2019.

Author contributions: H.F., N.M., S.S., and H.N. designed research; H.F., N.H., K.U., and H.N. performed research; H.F. and N.M. analyzed data; H.F., N.M., and H.N. wrote the first draft of the paper; N.M. and H.N. wrote the paper;

S.S., K.U., J.L.G., N.I., M.H., K.C., and H.N. edited the paper; J.L.G., N.I., and M.H. contributed unpublished reagents/analytic tools.

This work was supported by Japan Society for the Promotion of Science (KAKENHI Grants 26120732 and 16H06570 to H.N.). We thank C. Yokoyama and Y. Mochizuki for helpful discussions and comments on the manuscript and T. Asamizuya and C. Suzuki for technical assistance with the fMRI experiments.

mizes one's own reward. Considerable empirical evidence supports the behavioral quantification and its underlying computations and neural signals in the human brain (Rangel et al., 2008; Rushworth et al., 2011).

Human decision-making behavior is also often influenced by social factors (Montague et al., 2006; Rilling and Sanfey, 2011; Glimcher and Fehr, 2014), including concerns for others' welfare, mentalizing, social norms, and social interactions (Stanley and Adolphs, 2013; Ruff and Fehr, 2014; Lee and Seo, 2016). Previous studies using value-based decision-making frameworks and quantitative approaches such as economic games indicated that social factors can modulate brain signals in valuation and decision-related areas in addition to social areas such as the temporoparietal junction (TPJ) (Hsu et al., 2008; Behrens et al., 2009; Fehr and Krajbich, 2014; Hutcherson et al., 2015). Social value, or reward to others, is the simplest of the various social factors, but an important one. Previous studies have examined the neural correlates of social value in various altruistic behaviors (Zaki and Ochsner, 2012; Gospic et al., 2014; Kuss et al., 2015; Strombach et al., 2015; Crockett et al., 2017). In the studies, subjects make choices typically when rewards to the self and others are comparable and antagonistic to each other. The studies showed that the subject often makes decisions considering the balance of reward allocation to the self and others, and have eloquently demonstrated that brain regions related to reward, cognitive control, and emotion are involved in the different components of decision making under different reward allocation balances. Although variety in reward allocation balance is important in human social behavior, a more fundamental question, at a lower level than the various balances, remains elusive: How are brain signals for the social value of others' reward converted into downstream brain signals for decision making? We term such a neurocomputational process "social value conversion." Examination of social value conversion would require clearly and quantitatively separating the computations and brain signals for social value in reference to those for self-regarding decision making and then tracing the integration of social value into self-regarding decision making.

In this study, we investigated social value conversion by concurrently tracking the computational and neural stages by which social value signals migrate into the self-regarding decision-making process. In the behavioral task, the subject chose one of two options, each with a probabilistic reward to the self (standard reward) attached. We examined social value conversion as the behavior-modifying effect on choice behavior of an additional reward to others (other-bonus reward) versus an additional reward to the self (self-bonus reward). We used computational modeling and quantitative fMRI analysis to derive and track three major stages of social value conversion. The first stage is other-regarding detection involving the offered other-bonus value. The second stage is an intermediate value, called the effective value, which links the offer and final decision and represents the actual impact of the offer on choice. The third stage involves the final decision value to make choices. Further, we predicted that social value conversion is a primitive computation that may be essential for different forms of social behavior including inter-

individual variation. We then examined individual differences in social value conversion to probe the neural basis of social preference (Bogaert et al., 2008; Haruno and Frith, 2010). Social preference collectively refers to social concerns regarding the allocation of rewards between the self and others given an individual's predispositions such as inequity aversion (Frith and Frith, 2012; Fehr and Krajbich, 2014). Here, we examined one type of social preference called social value orientation (SVO) (Van Lange et al., 1997) and found that distinct prosocial and individualistic SVO phenotypes in subjects showed differences in brain signals and processes of social value conversion.

Materials and Methods

Subjects

Forty-seven healthy, normal subjects (17 female; 46 right-handed; age, 20–32 years; mean \pm SD, 22.0 \pm 2.41 years) participated in our main experiment. All subjects were screened to exclude those with a previous history of neurological or psychiatric illness. Subjects provided written informed consent to participate in this study, which was approved by RIKEN'S Third Research Ethics Committee. Because of large head motion during fMRI scans (>2 mm in any direction), 4 subjects were excluded from the analysis. The remaining 43 subjects were included in our main behavioral analysis. In the behavioral model fitting, we found that the choice behavior of seven subjects (despite the purpose of this study and the control of the experimental procedure) was not positively modulated by other-bonus rewards (for details, see "Model selection"). Therefore, the results of 36 subjects were included in the blood-oxygen-level-dependent (BOLD) signal analyses.

To increase their sample size and to facilitate further analysis of the individualistic group (for details, see "Analysis of a larger sample of individualistic subjects" section), we also conducted an additional round of the experiment in which individualistic subjects were specifically recruited (seven new subjects).

Experimental task

Each subject performed two tasks, called the control task and main task (Fig. 1); the main and control task were composed of other-bonus trials and standard trials and self-bonus trials and standard trials, respectively (Fig. 1B). Standard trials provided the reference for the choice behavior, in which one of the two options was selected, each associated with a probabilistic reward outcome. The other-bonus trial was designed to probe modification of the reference choice behavior. In addition to the probabilistic reward outcome, one of the two options had an additional reward designated as the other-bonus. When the subject chose the option with bonus, it was endowed to a specific charity (for details, see "Others' reward and charity organizations" section) regardless of the probabilistic outcome. The self-bonus trial served as the reference for the modified choice behavior in other-bonus trials and was designed to probe the modification of the choice behavior when the bonus was for the subject and not for others.

Each trial consisted of four phases (Fig. 1A). For trials in the fMRI experiment, at the beginning of each trial, a pair of options with a fixation point between them was presented to the subject for 1.5–4.5 s (CUE phase). The fixation point then changed to a question mark, and the subject could make their choice by pressing a button with their right hand within 1.5 s (RESPONSE phase). After the response phase, the chosen option was indicated by a gray frame, initiating the CONFIRM phase (0.5–1.5 s). Then, to show whether the subject's chosen option was rewarded was revealed by displaying a circle or cross, respectively, in the center of the screen for 1.0 s (OUTCOME phase). This was followed by a jittered intertrial interval (ITI; 1.5–5.5 s) before the next trial started. All timing jitters for the CUE, CONFIRM, and ITI phases were randomly generated by sampling from a uniform distribution. To maintain the subjects' motivation (Behrens et al., 2008), the accumulation of earned points over the trials was shown by a horizontal bar at the bottom of the monitor throughout the trials (except during ITI phases) in the control task, and the earned points for both the self and others were shown in the main task.

This study is dedicated to the memory of Kang Cheng, who passed away on Nov. 8, 2016.

The authors declare no competing financial interests.

*H.F. and N.M. contributed equally to this work.

†Deceased Nov. 8, 2016.

Correspondence should be addressed to Hiroyuki Nakahara at hiro@brain.riken.jp.

<https://doi.org/10.1523/JNEUROSCI.3117-18.2019>

Copyright © 2019 the authors

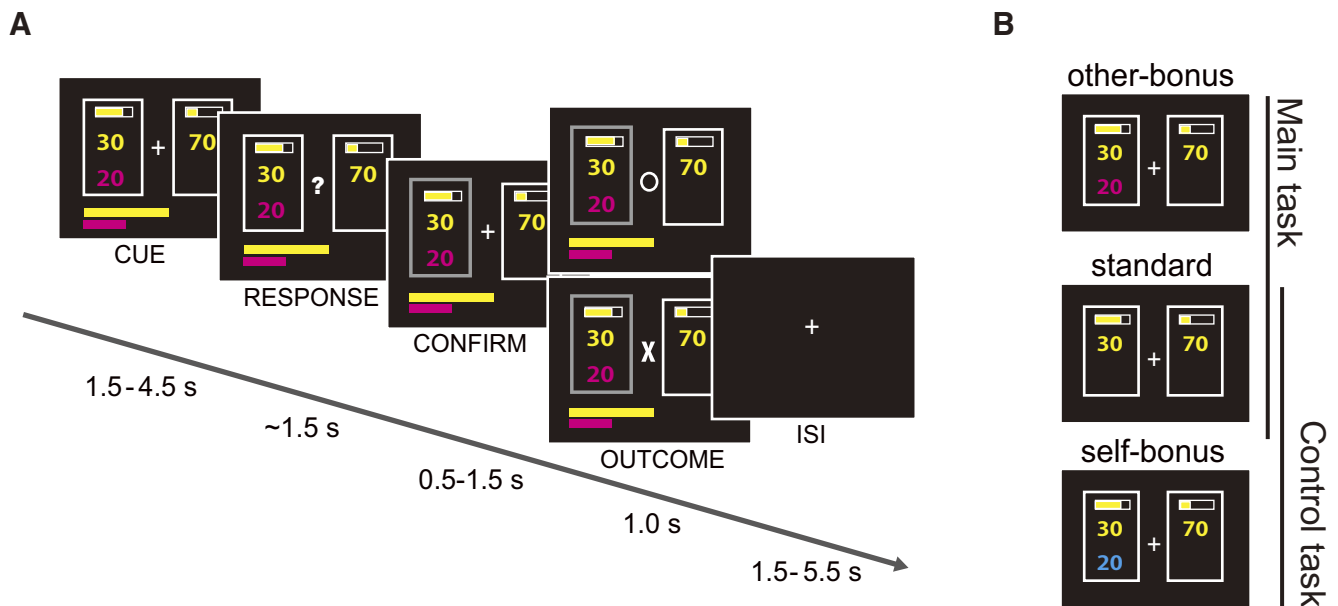


Figure 1. Experimental task. **A**, Example of an other-bonus trial in the fMRI experiment. The subject was asked to choose between two options to maximize their reward gains. Here, the left option was chosen so the standard reward (i.e., the yellow number indicating the reward magnitude, 30 points) would be given with a probability indicated by the yellow bar at the top, and the other-bonus (20 points), indicated by the magenta number at the bottom, would be given to a charity. The accumulated payments for self and others were indicated by the yellow and magenta horizontal bars at the bottom for each trial, respectively. **B**, The main task was composed of other-bonus (magenta) and standard trials and the control task was composed of self-bonus (cyan) and standard trials. For display purpose, the color of numbers indicating other-bonus and self-bonus are shown in magenta and cyan, respectively, in the above figure; the actual experiment used red and green, respectively.

Setting of rewards

The standard trial is a one-armed bandit task in which subjects were instructed to choose the option that would maximize the number of points earned among two options presented in the trial. Each option consisted of reward magnitude indicated by a number and a reward probability indicated by a filled yellow horizontal bar above the number, respectively. In every trial, the reward probability (p) for one option was uniformly randomly sampled from 0.1, 0.2, . . . 0.9 and the probability for the other option was set to $(1 - p)$. The reward magnitude for one option (m) was uniformly randomly sampled from 1, 2, . . . 99 points and the reward for the other option was set to $(100 - m)$. The probability and magnitude were first sampled independently and then assigned randomly to each option but with one constraint: larger probability and larger magnitude were paired randomly with probability 0.2 (i.e., not with probability 0.5, which would occur by completely random assignment) to avoid too many “easy” trials. In the other-bonus and self-bonus trials, additional reward (bonus) was randomly assigned to one of the two options in each trial. It was indicated by a colored number below the standard reward magnitude. For example, for some subjects, the number was shown in red and green in the self-bonus and other-bonus trials, respectively, and these colors were counterbalanced across the subjects. In contrast to the probabilistic nature of the standard reward, the bonus was deterministic, meaning that it was always given to the subject or a charity if the option was chosen.

The magnitude of the bonus was drawn randomly from 0, 10, and 20. Zero magnitude was included to exclude the possibility that the mere display of a bonus value affected the choice behavior. This setting was chosen based on the results of preliminary behavioral experiments using approximately the same distribution of standard rewards as in the main experiment (data not shown). This setting was decided based on the following considerations for the other-bonus. First, for most subjects, a detectable change in the choice behavior should be observed, but the magnitude of the bonus should not be excessively large compared with the standard reward because additional cognitive/behavioral factors such as envy, which should be avoided in this study, might come into play. Second, we aimed to include a smallest number of different magnitudes to have a sufficient number of trials to estimate other-bonus values in behavior given that each magnitude of trials should be examined with different

standard rewards. After determining the magnitudes of the other-bonus, we chose to set the same magnitudes for the self-bonus to directly compare the behavioral changes between the other- and self-bonus.

Others' reward and charity organizations

The others' reward was set to be given to the charity organization selected by the subject before the behavioral experiments. We chose this setting to make sure that the others' reward would matter to the subject in a positive way and that the relationship between the subject and others would instead be neutral, neither too close (in-group) nor too remote (out-group).

The subject chose one charity organization from a list of six. To ensure that the subject correctly knew about the organizations' activities, the subject read descriptions of all the charity organizations before choosing. All six of the listed organizations are well known in Japan: the Japanese Red Cross Society, the UN Refugee Agency, Médecins Sans Frontières (Doctors Without Borders), UNICEF, Central Community Chest of Japan (*Akai Hana Kyōdōbokin*), and Ashinaga (a charity that focuses on education and support for disadvantaged children).

After becoming familiarized with the tasks, the subject was instructed to choose a charity. Two slightly different instructions were given across subjects. One said to choose the organization to which their earnings from the others' reward would be given (16 subjects) and one said to choose the organization to which they would most prefer it be given (27 subjects). We merged them together in the main analysis because there was no significant difference between the two groups with respect to behavior (behavioral weights: self-bonus, $t_{(41)} = 1.244$, $p = 0.220$; other-bonus, $t_{(41)} = 1.877$, $p = 0.068$; proportion choosing options with the other-bonus: $t_{(41)} = 0.855$, $p = 0.398$; SVO classification: $\chi^2 = 0.216$, $p = 0.642$). Also, by questionnaire, we asked the subjects whether they already knew about each organization, to what extent (9-point scale) they were willing to support the activities of each organization, and whether they had ever directly participated in the organization's activities beyond donation. Prior participation in the organization's activities was a criterion for exclusion from the experimental analysis, but we found that no subjects had done so. All of the subjects knew beforehand about the organization that they chose, expressed a positive attitude toward it, and had not directly participated in its activities.

Experiment procedure

All subjects participated in both the behavioral and fMRI experiments (outside and inside the scanner, respectively). For all three types of trials (standard, self-bonus, and other-bonus), the subject was instructed to choose which of two options they preferred. Before the behavioral experiment, the subject was familiarized with the tasks, performing 35 trials for each trial type. The subject sat in front of a computer screen to do the behavioral experiment (three control and three main blocks), which together took about 120 min. Each block in the behavioral experiment was composed of 15 standard trials and 90 bonus trials (90 other-bonus or self-bonus trials in the main or control task, with 30 trials for each of the following bonus magnitudes: 0, 10, and 20). The order of trials was permuted randomly within a block. We used slightly shorter CUE and ITI phases in the behavioral experiment than in the fMRI experiment (1.0–3.0 s and 1.0–3.0 s, respectively) to enable a larger number of trials. The subject's behavior (decision: choosing the left or right option) was recorded by button pressing on the keyboard of the computer. For the behavioral results, we combined the trials of the fMRI and behavioral experiments because no difference in behavior was found between the behavioral experiments and all experiments (fMRI + behavioral experiments), as indicated by the behavior weight for the other-bonus ($t_{(42)} = 0.283, p = 0.779$).

After the behavioral experiment, the subject was put into the scanner for the fMRI experiment. Before the experiment, the subject performed a short practice session (20 standard trials). Each of six blocks in the fMRI experiment was composed of nine standard trials and 42 bonus trials (42 other-bonus or self-bonus trials, with 14 trials for each of the bonus magnitudes). The fMRI experiment took about 120 min, including preparation time, anatomic scanning, and functional scanning, during which the subject's behavior and brain responses were recorded.

At the beginning of the experiment, the subject was informed about the monetary reward that they would receive, which was based on the average of all the points they earned in the behavioral and fMRI experiments (i.e., the points they gained from both standard and self-bonus from 12 blocks in total). In brief, the total monetary compensation in yen was approximately equal to $200 \times (\text{average points} - 14) + 6000$, where 6000 yen was the base participation fee. For the other-bonus, the monetary conversion did not have the base fee: $200 \times (\text{average points of others' reward} - 14)$. Ultimately, the average monetary gain of the subjects was 9283 ± 410 yen (mean \pm SD; range, 8200–9800 Japanese yen) and that of the others was 1031 ± 262 yen (range, 1000–2430 Japanese yen).

Postexperiment questionnaire on SVO

After the experiment, the subjects were asked to complete an SVO questionnaire composed of nine items, which was adopted from a questionnaire used in a previous study (Van Lange et al., 1997). The SVO questionnaire was introduced to the subjects only after the main experiment (except the additional sampling of seven new subjects for whom SVO was used for screening before the main experiment). Therefore, the subjects performed the main experiment without knowing that they would also perform the SVO questionnaire. Therefore, their behavior in the main task was not directly affected by their answering the questionnaire.

The items were all forced choice questions with three options, each indicating a point gain distribution between the self and another person. The instructions were also adopted from the previous study (Van Lange et al., 1997). Briefly, the subject was instructed to think of an anonymous person to be paired with, whom the subject had never met and would never meet. The subject was told to consider that their choices would determine the points that they would earn and also that the other person would earn. They were told to assume that the other person would do the same. The subject was also informed that the point gains in this questionnaire would not be used for actual payment to themselves or others. Each option contained two numbers, one indicating the points that the subject would gain and one indicating the points that the other person would gain (e.g., [50, 40]). The subject was asked to choose the option that they most preferred. The three options in each question, with the distribution of the subject's and the others' gains, corresponded to the designations of

prosocial, individualistic, and competitive (e.g., [45, 45], [50, 40], and [45, 0], respectively).

Each subject was classified by their questionnaire responses into one of the three SVO phenotypes (prosocial, individualistic, and competitive), following the criterion of the original study: a subject who chose one characteristic (e.g., prosocial) six times or more times was classified as having that characteristic. When there was no such dominant SVO characteristic in their responses, subjects were unclassified. Among the 43 subjects in our main analysis, no one was classified as competitive. Therefore, in this study, we focused on only prosocial and individualistic subjects. Among the 36 subjects in the main experiment, we found 21 subjects were prosocial, 12 were individualistic, and three were unclassified.

Behavioral analysis and computational models

Choice behavior plot. We plotted the subjects' choice behavior (Fig. 2A), applying smoothing by Gaussian filter (variance = 10) to the standard value difference (ΔV_S) axis for the choice behavior (right option, 1; left option, 0) of individual subject in each type of trial separately for each condition (for a specific magnitude of either the self/other-bonus, attached to the left or right option). Then, in each case, the mean and SE were calculated across all subjects. Based on the sigmoidal choice curves, we quantified the effects of the bonus on the choices, thus obtaining the extent of behavioral change (Fig. 2B). For each bonus magnitude in the self-bonus and other-bonus trials of each subject, we took the difference in choice probability between each sigmoidal choice curve and that of the corresponding trial with zero magnitude bonus at the point of indifference ($\Delta V_S = 0$). Then, we used Page's test (Page, 1963) (a nonparametric, repeated-measure test for monotonic order) to examine our hypothesis that the extent of behavioral change increases with increasing bonus magnitude (other- and self-bonus in the corresponding trials), wherein the null hypothesis was that there was no monotonic order between the two factors.

Behavior modeling. Subjects' choice behaviors were modeled as a logit function of decision value (DV, the value difference between the two options). The logit function is given by the following:

$$\text{Logit}(q) \equiv \log\left(\frac{q}{1-q}\right)$$

where q indicates the choice probability. As shown later, this can be rewritten as $q = f(Z)$, where $f(Z) \equiv \frac{1}{1 + e^{-Z}}$, $Z = \beta \times DV$, and β is a free parameter often called inverse temperature. Without loss of generality, we write q as the probability of choosing the option on the right side. Note, however, that in many of the fMRI analyses that follow, we often used this function in the alignment of the choice made; in other words, we used the function of choice probability that would correspond to $\text{Logit}(q)$ and $\text{Logit}(1 - q)$ in the equation above for the right and left option chosen, respectively; it was used to fit the model to behavior by the maximum likelihood approach, and to analyze behavior and the variables in BOLD signals.

DV is modeled in our analysis as follows:

$$DV = \Delta V_S + w_S \Delta S + w_O \Delta O \quad (1)$$

Where Δ indicates the difference between options (right minus left option), and V_S , S , and O are variables relating to the standard, self-bonus, and other-bonus, respectively. More precisely, $\Delta V_S = V_{SR} - V_{SL}$ indicates the standard value difference, where V_{SR} and V_{SL} are the standard reward values of the right and left options, respectively. V_{SR} (or V_{SL}) was computed as the product of reward magnitude and probability, wherein the probability was allowed to be distorted by risk (Platt and Huettel, 2008). We chose to use the Prelec-2-parameter model (Prelec, 1998), given by $\bar{p} = \exp(-\delta \log^{\gamma} p)$, where the left side is the risk-distorted reward probability, p is a given reward probability, and δ and γ are two free parameters fitted by the logit function (see below). For ΔS and ΔO , note that only one option had an additional deterministic reward in each of the other-bonus and self-bonus trials, so the difference was the magnitude but with a sign (direction), either S_R or $-S_L$ (O_R or $-O_L$) when it

Table 1. Results of model selection

| Comparison of the main model with models dropping one variable each | | | | |
|---|-------------------------|-------------------------|-------------------------|-------------------------|
| | Drop ΔV_s | Drop ΔS | Drop ΔO | Main |
| AIC | 1.302 ± 0.050 | 0.692 ± 0.123 | 0.540 ± 0.143 | 0.509 ± 0.140 |
| Comparison with Main | $t_{(42)} = 36.149, *1$ | $t_{(42)} = 14.401, *1$ | $t_{(42)} = 5.746, *1$ | – |
| Selecting the best model for standard value | | | | |
| | $V_s = mp$ | $V_s = V_s = \bar{p}$ | $V_s = p$ | $V_s = m$ |
| AIC | 0.819 ± 0.202 | 0.810 ± 0.181 | 0.829 ± 0.167 | 1.343 ± 0.066 |
| Comparison with Main | $t_{(42)} = 9.187, *1$ | $t_{(42)} = 12.899, *1$ | $t_{(42)} = 13.495, *1$ | $t_{(42)} = 34.672, *1$ |

The Main model ($DV = \Delta V_s + w_s \Delta S + w_o \Delta O$, where $V_s = \bar{p}$ and $\bar{p} = \exp(-\delta \log^2 p)$) was selected based on AIC value. The main model was chosen based on consideration of DV (dropping each term) for $V_s = mp, V_s = V_s = \bar{p}, V_s = p$, and $V_s = m$ and other models (see Materials and Methods). We found that the main model was better than any of those models, which justified the use of the main model in this study. *1; $p < 0.001$.

was attached to the right or left option, respectively. Therefore, their respective weights, w_s and w_o , indicated the extents of modifications in the choice behavior by ΔS and ΔO , respectively.

Model fitting and selection. The maximum likelihood approach was used to fit the models to the behavioral data. For individual subjects, we minimized the sum of the negative log-likelihood of choice probabilities (MATLAB command `fminsearch`; MATLAB R2012b, The MathWorks) for the options chosen by each subject (i.e., the probabilities q and $1 - q$ when the right and left options were chosen, respectively). Each minimization was repeated 50 times using randomly generated initial values. Free parameters resulting in a minimal summed negative log-likelihood were then selected, corresponding to the best fit of the model.

To select the best model at the group level while taking into account the different number of free parameters between the models, a paired t test was used to compare differences in the distribution of Akaike’s information criterion (AIC) values obtained for the different models (Suzuki et al., 2012). First, to use the expected reward as a probabilistic outcome and include risk parameters (i.e., using rather than p), we examined the fit of our model ($V_s = m\bar{p}$) for standard trials in comparison with that of four reference models in which V_s was replaced with m, p, \bar{p} , or mp . Among the four reference models, the first three were used to test whether the reward magnitude or probability contributed to the expected reward, whereas the last reference model, mp , was used to test the inclusion of risk parameters. In the group-level comparison, the original model was the best model with statistical significance (Table 1). Therefore, we adopted $V_s = m\bar{p}$ in our analyses. Second, to further ensure our main model ($DV = \Delta V_s + w_s \Delta S + w_o \Delta O$), we compared the fit of our main model to that of three reference models, each of which lacks one of the three variables of the main model: $\Delta V_s, \Delta S$, or ΔO . By the group-level comparison of AIC value distributions, we found that the main model was superior to any of the three reference models. These results were also confirmed by using the Bayesian information criterion (BIC) for model selection (Table 1).

We also examined the models’ fit to the behavior of each individual subject to ensure that each of the main variables in the behavior model had a considerable effect on the individual behavior. First, for standard trials, the AIC values of the main model ($V_s = m\bar{p}$) were always the smallest compared with those of the m, p, \bar{p} , and mp models for all individuals, indicating that all individuals used both reward magnitude and probability in their valuation. Second, we compared the fit of the main model with that of the models omitting one of $\Delta V_s, \Delta S$, or ΔO for all the trials of each subject. We found that the main model’s fit was better for all individuals in comparison with the models omitting ΔV_s or ΔS . However, for the case of ΔO , we found that, for seven subjects, ΔO in the main model had no significantly positive effect compared with the model omitting ΔO and these subjects were excluded in our main analysis of BOLD signals.

We conducted several other verifications of the main model, including interaction terms, offset term, and accumulated payment for other’s reward, and also testing the effect of mere presence of a bonus. All of these verifications supported using the main model in this study, and these results are briefly summarized here. A model with an additional interaction term on DV ($w_{V_s \Delta S} \Delta V_s \times \Delta S + w_{V_s \Delta O} \Delta V_s \times \Delta O$) was significantly

worse than the main model in fitting the behavior ($t_{(42)} = 5.276, p < 0.001$ by AIC). A model with an additional constant (offset) was not superior to the main model ($t_{(42)} = 1.321, p = 0.192$ by AIC), but was significantly worse than the main model ($t_{(42)} = 5.522, p < 0.001$ by BIC), suggesting that the main model was at least not inferior. A model with a term of accumulated payment, $-w_{A_O A_O} \times \Delta O$ was significantly worse ($t_{(42)} = 3.117, p = 0.003$ by AIC, where A_O indicates the accumulated payments for other). We tested this because this term might affect the behavior, decreasing the selection of the option with the other-bonus as the payment accumulated. Investigating whether the mere presence of bonus (i.e., zero magnitude of reward) affects the behavior revealed that, in the trials with additional zero rewards, the main model was better ($t_{(42)} = 4.664, p < 0.001, t_{(42)} = 14.805, p < 0.001$ by AIC, respectively) than two models with the effect of the mere presence of bonus (adding either the term $w_{[0] \Delta[0]}$ or the term $w_{[0]_S} \Delta[0]_S + w_{[0]_O} \Delta[0]_O$ to the DV), where $\Delta[0] = 1$ or 0 if the zero reward appeared on the right or left option and the subscripts S and O indicate the self- and other-bonus, respectively.

Signed effective value. We use signed effective value to indicate the link between an offered bonus value and a final decision for our BOLD analyses. Even for the same offer, the effect of the offered value on the choice (i.e., the amount of change in choice probability due to addition of the offer) may vary depending on the standard value difference (ΔV_s) in the particular trial. This is because of the nonlinearity of the logit function connecting a decision value to choice probability (q). For this reason, the variable signed effective value was defined as a second-stage variable that represents the effective impact of the offer (the first stage) on the choice (the third stage).

First, given a trial with a particular standard value difference (ΔV_s), we assessed the effect of an offered value on choices. This was done using the (choice) effectivity, which is defined as the derivative of the choice probability (q) with respect to ΔV_s and quantifies how the offered value affects the choice as follows:

$$\text{effectivity} = \left. \frac{\partial q(Z)}{\partial Z} \right|_{Z=\Delta V_s} \tag{2}$$

Then, the signed effective other-bonus value is formulated as the product of the signed other-bonus value (signed with respect to choice, e.g., $w_o \Delta O$ for the other-bonus) and the effectivity as follows:

$$\text{signed effective other-bonus value} = w_o \Delta O \times \left. \frac{\partial q(Z)}{\partial Z} \right|_{Z=\Delta V_s} \tag{3}$$

Therefore, the signed effective other-bonus value corresponds to the first-order Taylor expansion term for the variable ΔO at $Z = \Delta V_s$ in the logit function. Similar calculations have been used in previous model-based fMRI studies (Bornstein and Daw, 2012, 2013), which also estimated the dependent variables in nonlinear equations by using the first-order Taylor expansion term. Analogous formulations and terminology were also used for the self-bonus value, such as the signed effective self-bonus value as follows:

$$w_s \Delta S \times \left. \frac{\partial q(Z)}{\partial Z} \right|_{Z=\Delta V_s}$$

Further confirmations for similar behaviors in SVO groups. Even though the two groups of subjects responded differently to the SVO questionnaires, they behaved similarly in our tasks, as indicated by our finding (see Results) that the behavioral weight for the other-bonus was not significantly different between the two groups. For further confirmation, we conducted two additional analyses (Table 2). The first analysis examined the proportion choosing the option with the other-bonus in all other-bonus trials to determine whether the behavior in our task was different between the two groups. This examination was also conducted for smaller subsets of the trials (subsets I and II). Subset I was trials with a conflict between reward gain to the self and others: the trials in which the standard value of the option attached with other-bonus was smaller than that of the option without the other-bonus. Subset II contains the trials in which the two options are similar in nature to those in the SVO

Table 2. Prosocial and individualistic groups behaved similarly within our task

| | Prosocial subjects | Individualistic subjects | Prosocial versus individualistic |
|---|--------------------|--------------------------|---|
| Proportion choosing the option with the other-bonus reward | | | |
| All trials | 0.550 ± 0.049 | 0.553 ± 0.028 | $t_{(31)} = 1.572, p = 0.126$ |
| Subset I | 0.388 ± 0.053 | 0.379 ± 0.032 | $t_{(31)} = 0.333, p = 0.741$ |
| Subset II | 0.722 ± 0.1666 | 0.707 ± 0.176 | $t_{(31)} = 1.341, p = 0.190$ |
| Negative log-likelihood of the responses of the SVO questionnaire | | | |
| SVO prosocial model | 0.637 ± 0.104 | 0.707 ± 0.090 | Better fitted to prosocial subjects ($t_{(20)} = 8.879, p < 0.001$) |
| SVO individualistic model | 0.812 ± 0.112 | 0.604 ± 0.060 | Better fitted to individualistic subjects ($t_{(11)} = 3.036, p = 0.039$) |

Two additional tests that confirmed similar behaviors in our task between the prosocial and individualistic subjects, but distinct behavior in SVO questionnaire responses (for details, see Materials and Methods). In the first test for the behavior in our task, there was no difference between the two groups of subjects in the proportion choosing the option with other-bonus in all (other-bonus) trials, and even two types of subset trials (subsets I and II). In the second test for distinct behavior in the questionnaire responses, we employed a logistic regression analysis to the responses in the SVO questionnaire, devising corresponding SVO prosocial or individualistic model (see Materials and Methods). We found that the SVO choices by the prosocial and individualistic subjects were better fitted by the corresponding SVO model than by the other model. Therefore, whereas this confirms that their SVO choice behavior matched with the corresponding SVO model, it clearly indicates that the individualistic subjects behaved differently regarding other-bonus between our task and the questionnaire. Although they appreciated other-bonus in our task (as shown by their w_o being significantly larger than zero), they instead ignored them in their SVO choices.

questionnaire. These trials were selected as follows: (1) all other-bonus trials in our task and all pairs of options (prosocial vs individualistic option) in the SVO questionnaire were projected onto a 2D map with the self-regarding difference on the x -axis and the other-regarding difference on the y -axis; (2) in this map, a wider range was covered by our task than by the SVO questionnaire, so (3) we selected the other-bonus trials in the range of the SVO questionnaire on both the x - and y -axes.

In the second analysis, we investigated whether the two groups of subjects behaved differently in the responses to the SVO questionnaire using models based on the behavior weights estimated from our task. For this, we set two models that would approximately correspond to two original SVO classifications: an “SVO prosocial” model that used the weights of both w_s and w_o , and an “SVO individualistic” model that used only w_s (setting $w_o = 0$). Using logistical function analysis, the negative log likelihoods of the fit by each model to the behavior in the prosocial and individualistic group are shown in Table 2.

fMRI

Data acquisition and preprocessing. The fMRI images were collected using a 4 T whole-body MRI system (Agilent Technologies) with a transverse electromagnetic volume coil as the transmitter (Takashima Seisakusho) and a 16-array head-shaped coil as the receiver (Nova Medical). For subjects positioned in the scanner, visual input was provided via a fiberoptic goggle system (Avotec) that subtended $25^\circ \times 19^\circ$ of visual angle, and the subject used a button box to make their responses. The BOLD signal was measured using a two-shot T2*-weighted echo planar imaging (EPI) sequence (230 volumes for each block, TR = 2202 ms, TE = 20.5 ms, FA = 64°). Twenty-five axial slices (thickness = 3.0 mm, gap = 1 mm, FOV = 192×192 mm, matrix = 64×64 , thus resulting in voxel size = $3 \times 3 \times 3$ mm) parallel to the AC-PC plane (AC: anterior commissure; PC: posterior commissure) were acquired per volume. The start of an experimental task was synchronized with the first EPI acquisition timing. Before, after, and between the functional runs, a set of high-resolution (1 mm₃) and a set of low-resolution (1.72 mm₃) whole-brain anatomical images were acquired using a T1-weighted 3D MPRAGE pulse sequence (TI = 500 ms, FA = 15° , TR = 9.5 ms, TE = 3.7 ms for the high-resolution scans or 2.5 ms for the low-resolution scans, TR = 7.3 ms). The low-resolution anatomical imaging slices were parallel to the functional imaging slices and were used to aid in coregistering the functional data to the high-resolution anatomical data. A pressure sensor was used to monitor and measure the respiration signal, and a pulse oximeter was used to measure the cardiac signal. The respiratory and cardiac signals were used in postprocessing to remove physiological fluctuations from functional images (Hu et al., 1995).

Functional and anatomical images were mainly analyzed using Brain Voyager QX 2.8 (Brain Innovation; RRID:SCR_013057). Functional images for each subject were preprocessed, which included slice time correction, 3D motion correction, spatial smoothing with a Gaussian kernel (FWHM = 8 mm), and high-pass temporal filtering (three cycles per run length). Anatomical images of each subject were transformed into the standard Talairach space (Talairach and Tournoux, 1988). Functional images were then normalized and resized according to the transformed

structural images and thus transformed into the standard Talairach space. Then, images from all scanning sessions were connected.

Generalized linear model analysis. We used a so-called model-based analysis (O’Doherty et al., 2007) to analyze the BOLD signals in both control and main tasks, using generalized linear model (GLM) regression with two levels of analysis as follows. At the first level, we submitted the BOLD signals of an individual subject into the GLM. When generating the GLM regressors, we used the behavioral weights of our model, which were estimated by fitting to the trials in fMRI experiment of each individual subjects and the estimated behavioral weights were also used as a covariate at the second-level analysis.

Using model-based analysis, we examined computations of social value conversion in three essential stages. The first and third stage correspond to the input and output stages of social value conversion, respectively, whereas the second stage links the two stages. A variable in the first stage, called offered value (e.g., $w_o|\Delta O|$, offered other-bonus value) indicates the value of the possible outcome ($|\Delta O|$) to others before one’s choice that would affect the others’ outcome. The variable in the second stage is called signed effective value (see “Signed effective value” section). The variable in the third stage is the DV (Eq. 1), which is assumed to determine the choice (or choice probability) given by the value of the chosen option minus the value of the unchosen option.

Therefore, we created subject-specific design matrices containing the following regressors for the GLM. There were 12 regressors for the variables of no interest. Six regressors encoded the average BOLD responses for the onset and period of the DECISION, ISI, and OUTCOME phases, where the DECISION phase was defined as the period from the onset of CUE until the subject responded in the RESPONSE period and the other two phases were defined as shown Figure 1A. Six motion correction parameters to account for motion effects, and the outcome for the standard reward in each trial (1 when rewarded and 0 when not, in the OUTCOME phase) were also included. There were five regressors for the variables of interest, (1) the decision value (DV, on the axis of the chosen minus unchosen option); (2) offered other-bonus value, $w_o^i O$, and (3) signed effective other-bonus value, $w_o^i(\partial q / \partial \Delta V_s)\Delta O$ (Eq. 3), only for other-bonus trials; (4) offered self-bonus value, $w_s^i S$, and (5) signed effective self-bonus value, only for self-bonus trials. All regressors for the variables of interest were modeled in the DECISION phase. All regressors of interest were mean-corrected, normalized, and convolved with a canonical hemodynamic response function before being entered into GLM analysis.

The estimated effect sizes or their contrasts from the first-level GLM were entered into a whole-brain random-effects analysis to extract significant brain activations for each regressor in group-level statistics (i.e., groupwise one-sample t test). In the random-effects analysis, we also included a covariate for examining relationships with individual variability in the behavioral weights. Specifically, we included w_o^i for offered and signed effective other-bonus values and w_s^i for offered and signed effective self-bonus values. The significance of these covariates was tested by t test (against a null hypothesis of zero estimated effect size). The brain regions with significant effect sizes for each regressor were reported based

Table 3. Areas exhibiting significant changes in BOLD signals by GLM analyses

| Activated clusters | LR | x | y | z | BA | k | corrected <i>p</i> |
|---|----|-----|-----|-----|----|------|--------------------|
| DV | | | | | | | |
| mPFC/medial frontal gyrus | — | 0 | 59 | 4 | 10 | 450 | 0.035 |
| Middle temporal gyrus | R | −60 | −4 | −2 | 21 | 411 | 0.034 |
| Self-bonus | | | | | | | |
| Precuneus/inferior parietal lobes | R | 9 | −70 | 37 | 7 | 7035 | 0.000 |
| rdlPFC/dmPFC/medial frontal gyrus | R | 9 | 26 | 37 | 6 | 2041 | 0.000 |
| Cingulate gyrus | LR | 6 | −31 | 28 | 23 | 1459 | 0.000 |
| ldlPFC/middle frontal gyrus | L | −27 | 50 | 16 | 10 | 231 | 0.036 |
| Other-bonus | | | | | | | |
| Precuneus/TPJ/middle temporal gyrus/inferior parietal lobes | R | 33 | −64 | 28 | 39 | 1473 | 0.016 |
| ldlPFC/middle frontal gyrus | L | −39 | 44 | 10 | 10 | 954 | 0.022 |
| dmPFC/medial frontal gyrus | R | 9 | 17 | 49 | 6 | 564 | 0.035 |
| Signed effective other value | | | | | | | |
| Insula/inferior frontal gyrus* | R | 36 | −4 | −11 | 21 | 59 | 0.037 |
| Conjunction of self and other-bonus value | | | | | | | |
| Precuneus/inferior parietal lobes | R | 28 | −65 | 30 | 7 | 1418 | |
| ldlPFC/middle frontal gyrus | L | −39 | 44 | 16 | 10 | 270 | |
| dmPFC/medial frontal gyrus | R | 11 | 16 | 45 | 6 | 535 | |
| Contrast of self and other-bonus value | | | | | | | |
| TPJ | R | 60 | −58 | 31 | 10 | 179 | |

Activated clusters of the variable were observed by voxelwise GLM analysis of BOLD signals, $p < 0.05$, FWE-corrected by permutation test (cluster-defining threshold: $p < 0.005$) (Winkler et al., 2014). For the brain regions marked with asterisks, activated clusters were observed by the corresponding covariate of the variable in the second-level analysis. Activations in a whole-brain voxelwise conjunction (Nichols et al., 2005) was extracted by taking logical AND over the two activations (obtained by thresholding at uncorrected $p < 0.005$) with the same statistical criterion. Activation by a whole-brain voxelwise contrast for offered other-bonus > self-bonus was also obtained using the same statistical criterion. The stereotaxic coordinates are in Talairach space (for details, see “GLM analysis” section). BA, Brodmann’s area; BOLD, blood-oxygen-level dependent.

on corrected p -values ($p < 0.05$), using familywise error (FWE) correction for multiple comparisons; we first thresholded contrast maps at $p < 0.005$ (i.e., uncorrected) and then estimated corrected p -values by permutation test (5000 permutations) using PALM software (FMRIB, University of Oxford, <https://fsl.fmrib.ox.ac.uk/fsl/fslwiki/PALM>) (Winkler et al., 2014). The coordinates reported in Tables 3 and 4 for a given cluster were those of the highest activated voxel within it. To report the coordinates for a given anatomic brain region, such as the rTPJ and ldlPFC, we further used a higher-values-first watershed searching algorithm by BVQX tools (version 0.8d) of Brain Voyager to identify the local maxima peak in the respective brain region (thus, the peaks reported in the text were not necessarily the same as reported in the tables).

Trial-by-trial GLM analysis. Trial-by-trial GLM analysis (whole-brain) was performed to examine the BOLD signal changes on a trial-by-trial basis for each subject (Fig. 3B) (Gläscher, 2009). This GLM analysis focused on the DECISION phase in each trial; thus, the regressor of interest encodes 1 for the DECISION phase of the trial, whereas any other periods or trials were encoded by 0. To conduct the trial-by-trial GLM analysis for the results shown in main Figure 3B, we used the leave-one-out, cross-validated region of interest (ROI). For each subject, we obtained the cross-validated ROI and then extracted the BOLD signal changes in each trial. For each subject, we used z transformation to normalize the BOLD signal changes across all trials (subtracting the mean and dividing by the SD) and binned them separately for each trial type as low, medium, or high (corresponding to the 33rd, 66th, and 100th percentiles, respectively) to obtain the individual’s binwise mean BOLD signal changes. Then, the mean and SEM of the normalized BOLD signal changes in each bin were computed across subjects. We investigated whether the BOLD signal changes increased with the order of the bins by Page’s test.

Extraction of ROIs. The ROIs were extracted as 8-mm-diameter spheres with centers at the local maxima peaks within the clusters of GLM activations. The local maxima peaks were obtained by a higher-values-first watershed searching algorithm in BVQX tools. To extract cross-validated BOLD signals, we used a previously reported leave-one-out procedure (Gläscher et al., 2009) to provide an independent criterion for ROI selection and thus ensure statistical validity (Kriegeskorte et al., 2009). We reestimated our second-level analysis n times (where n equals

Table 4. Areas exhibiting significant changes in BOLD signals by voxelwise PPI analyses

| Brain region | BA | Hemi | x | y | z | Cluster size | Corrected <i>p</i> |
|---|----|------|-----|-----|-----|--------------|--------------------|
| ALL subjects | | | | | | | |
| rAI × signed effectivity of other-bonus value on mPFC responses | | | | | | | |
| | 10 | − | 1 | 57 | 4 | 23 | 0.018 |
| rTPJ × effectivity of other-bonus value on rAI responses | | | | | | | |
| | 10 | R | 30 | −10 | −8 | 27 | 0.002 |
| ldlPFC × effectivity of other-bonus value on rAI responses | | | | | | | |
| | 13 | R | 48 | 5 | −11 | 31 | 0.008 |
| ldlPFC × offered self-bonus value on mPFC responses | | | | | | | |
| | 10 | R | 9 | 47 | −5 | 17 | 0.009 |
| dmPFC × offered self-bonus value on mPFC responses | | | | | | | |
| | 10 | R | 9 | 62 | 4 | 32 | 0.008 |
| Prosocial subjects | | | | | | | |
| ldlPFC × offered other-bonus value on mPFC responses | | | | | | | |
| | 10 | R | 3 | 56 | 1 | 56 | 0.042 |
| ldlPFC × offered self-bonus value on mPFC responses | | | | | | | |
| | 10 | L | −12 | 41 | 10 | 73 | 0.002 |
| Individualistic subjects | | | | | | | |
| rAI × signed effectivity of other-bonus value on mPFC responses | | | | | | | |
| | 10 | L | −3 | 65 | 10 | 12 | 0.009 |
| ldlPFC × offered self-bonus value on mPFC responses | | | | | | | |
| | 9 | R | 12 | 41 | 10 | 23 | 0.003 |

Activated clusters significant at the $p < 0.05$ level (FWE-corrected), using voxelwise PPI with SVC at the cluster level by permutation test within the corresponding target ROIs (e.g., indicated as “on XX responses,” XX = ROI; for details, see Materials and Methods). The rest of table format is the same as for Tables 3. BOLD, Blood-oxygen-level dependent. Side remarks: no cluster survived for the case of, [rTPJ × offered other-bonus value] in all subjects or each group, [ldlPFC × offered other-bonus value] in all subjects or the individualistic group, or for the case of [rAI × signed effectivity of other-bonus value] in the prosocial group (none of them was shown in the table above).

the number of subjects), leaving out a different subject each time. Starting at the peak voxels for the focal signal, we selected the nearest activated voxels ($p < 0.005$, uncorrected) in these cross-validation second-level analyses. The selected voxel was defined as a cross-validated ROI.

Analysis of rTPJ and ldlPFC responses while controlling the behavioral effects. For this reanalysis of the BOLD signals, we controlled the behavioral effect to be comparable between the other-bonus and self-bonus values by selecting a subgroup of subjects and using a subset of the trials. For this purpose, we down-sampled both the pool of self-bonus trials and the pool of subjects to analyze the BOLD signals of self-bonus trials in which the behavioral effect was reasonably comparable to that of other-bonus trials.

First, we took the ratio of behavior weights of the other-bonus relative to that of the self-bonus (i.e., w_o/w_s), and noted that the range of this ratio for all the fMRI subjects ($n = 36$) was ~ 0.5 , specifically, from 0.012 to 0.965, mean = 0.426, SD = 0.241. Therefore, we decided to exclude the trials with a self-bonus magnitude equal to 20, which allowed us to control the maximum of the self-bonus and other-bonus as 10 and 20, respectively, and to have their ratio be 0.5. Second, we sorted the subjects in order of their ratio and excluded subjects one by one from the low end of the ratio and monitored the mean of the ratio of the remaining subjects. We stopped excluding subjects when the mean of the ratio first became close to 0.5 (mean = 0.509, SD = 0.205, from 0.224 to 0.965). The number of remaining subjects was 28. Therefore, with the subsampling of the self-bonus trials and the remaining 28 subjects, we performed BOLD signal reanalysis for trials in which the behavioral effect was comparable between the other- and self-bonus.

In the reanalysis, our GLM was in the same form as the original GLM in our main analysis except that the regressors for the self-bonus value were modified to treat the trials having a self-bonus of only 10 or 0 points. Analyzing the corresponding effect sizes in ROIs (generated by the activations in Table 3), we found that activations in the rTPJ remain significant ($t_{(27)} = 3.088$, $p = 0.005$) for the other-bonus value, but not for the self-bonus value ($t_{(27)} = -0.686$, $p = 0.498$), with a significant difference between the values ($t_{(27)} = 3.139$, $p = 0.004$). We also found activation in the left dlPFC remained significant for both the self-bonus and other-bonus values ($t_{(27)} = 3.587$, $p = 0.001$, $t_{(27)} = 3.047$, $p = 0.005$, respectively) and with no significant difference ($t_{(27)} = 0.147$, $p = 0.884$).

Analysis of the effect size covariate with the individuals' behavioral variability. For the results shown in Figure 3F, we first used the leave-one-out, cross-validated ROI of the right anterior insula (rAI) and extracted the BOLD signal by averaging the signal from all voxels in the ROI from the left-out subject. The effect size was estimated by regressing this BOLD signal in the original GLM. Second, to investigate the group-level variability of the effect sizes in relation to behavior weights (e.g., w_O^i), we calculated Pearson's correlation coefficient between the corresponding variables and tested the corresponding statistical significance by *t* test.

Psychophysiological interaction (PPI) analysis. We conducted two types of PPI analyses: ROI PPI and voxelwise PPI (Friston et al., 1997; O'Reilly et al., 2012). In the former, the target brain region was determined based on the activations from the earlier GLM analyses (Table 4). We chose to use both PPIs for clarity because they complement each other. ROI PPI probed the significantly activated regions by using the interaction term to specifically target the relationship of the two activated brain regions obtained from our GLM results with a hypothesized psychological seed. Voxelwise PPI probed the activated regions by using the interaction term for each voxel, allowing us to examine the activation including the regions surrounding the target brain area. We could thus examine the overlap of the significant activations from the voxelwise PPI analysis with the activations from the original GLM.

The PPI regressors were constructed as follows. First, we determined each ROI with respect to a variable of interest based on the group-level activation maps generated from the GLM analysis. We then extracted the averaged BOLD time courses within a given ROI for each subject based on their preprocessed BOLD signals. These extracted signals were used as the signals of the physiological seed. We estimated the signals by deconvolving the time course signals with a canonical hemodynamic response function. Together with a given psychological seed of the variable of interest, we generated the interaction term (i.e., the variable of main interest). We first normalized each of the physiological and psychological terms to [0, 1] and then multiplied them together, further orthogonalizing the product to each of the two first-order terms. Third, the three terms (the interaction and the two first-order terms) were mean-corrected and convolved with a canonical hemodynamic response function. In both PPI analyses, to guard against possible confounding effects, we included not only the first-order terms of the interaction term as usual, but also other regressors, such as the signed effective other-bonus value.

For voxelwise PPI, the maps were analyzed as random effects by *t* test. The significant activations were determined and reported as $p < 0.05$, using small volume correction (SVC; the target ROI of the corresponding PPI defined the SVC ROI) and using FWE correction for multiple comparisons by permutation test. For ROI PPI, we determined the target ROI based on the significant brain activations found by the original GLM analysis and then conducted the PPI analysis. The mean and SEM of the effect size in the target ROI was computed across subjects and the statistical significance was tested using a one-sample *t* test against a null hypothesis of zero mean. The second-level covariates in the original GLM were also included in the both types of PPI analysis.

Dynamic causal modeling analysis. Dynamic causal modeling (DCM; by SPM12, <http://www.fil.ion.ucl.ac.uk/spm/>, RRID:SCR_007037) (Friston et al., 2003) was used to analyze the coupling directions and structures over the rTPJ, ldlPFC, rAI, and medial prefrontal cortex (mPFC) signals in relation the three-stage processing revealed by our PPI results. For each subject, we extracted average activation time courses from an 8-mm-radius sphere around the peak of each ROI. By way of example, using cases of Figures 4C and 5, we constructed eight models including a main model based on the three-stage processing and seven models with possible revised directions of connectivity (Fig. 5A). In this DCM analysis (and also all other DCM analyses), we set offered other-bonus value for the rTPJ and ldlPFC signals as driving inputs because we were interested in how the other-bonus value is propagated in social value conversion. Following our PPI results, we set as modulatory inputs the effectivity of other-bonus value for the rTPJ-rAI and ldlPFC-rAI connectivity and the signed effectivity of other-bonus value for the rAI-mPFC connectivity. Second, we performed random-effects Bayesian model comparison to identify which of the eight models best fit our

dataset (Stephan et al., 2009). We calculated the exceedance probability for each model relative to the seven other models (Fig. 5B). Third, a model family comparison analysis (also used in Fig. 8) was also conducted similarly using a random-effects Bayesian model selection procedure (Stephan et al., 2009). Further, we also used Bayesian model averaging to examine the estimate of the connection parameters (Penny et al., 2010). The parameters over all models in the model space were averaged using weights corresponding to the posterior model probabilities.

Comparison between prosocial and individualistic subjects. To explore possible between-group differences, we added the SVO classification variable as another covariate for the brain activations. The effect sizes and statistical tests of the covariate and constant terms are reported for the rAI responses in the prosocial and individualistic groups (Fig. 6B).

To analyze coupling, we examined each of the PPI effects in each SVO group by both voxelwise and ROI analyses (Figs. 7A,B, respectively). Repeated-measures ANOVA was used to identify the between-group difference in SVO, with the seed regions as a within-group factor. Further, to control for the effects of the self-bonus value, we performed a control PPI using [rAI \times signed effective self-bonus value] on the mPFC responses (separately for each group) and then explored the effect sizes of the ROI PPI results in the rAI and ldlPFC cases for both other-bonus and self-bonus values on a 2D map (Fig. 7C; details in the following section). DCM analyses were also conducted to examine each SVO group separately (Figs. 7D,E, 8C–E).

Circular statistical analysis for SVO PPI results. To examine the nature or the relative difference between the other-bonus and self-bonus values between SVO groups in Figure 7C, we defined a vector for each subject, starting from the self-bonus and ending at other-bonus point. We then obtained the mean vector for each group. Using the projection from each vector onto the group mean vector, we assessed the major component (the value of the projection) of each subject's vector and examined the group effect across the subjects in each group, using the *t* test. The difference between the self-bonus and other-bonus values was significant in the individualistic subjects ($t_{(11)} = 2.254, p = 0.046, p = 0.027$ by bootstrap test), but not in the prosocial subjects ($t_{(20)} = 0.913, p = 0.372, p = 0.352$ by bootstrap test) and the difference in differences between the two groups was significant ($t_{(31)} = 2.145, p = 0.039, p = 0.017$ by bootstrap test). In addition, we traversed the vectors within each group to the origin and examined their vectorial changes; the distribution over the vector directions was examined by the Rayleigh test within each group and compared between the two groups by the Watson-Williams two-sample test. The vector difference in the effect sizes was significantly nonuniformly distributed across the individualistic subjects ($z = 3.054, p = 0.044$ by Rayleigh test), but not across the prosocial subjects ($z = 0.079, p = 0.926$). Focusing on angular changes, we also observed that the difference between the two groups was significant ($F = 4.255, p = 0.026$ by Watson-Williams two-sample test).

Analysis of a larger sample of individualistic subjects. The number of the individualistic subjects was relatively small ($n = 12$) in the main experiment, so we conducted two additional analyses and obtained essentially the same results with a larger sample size, both when adding four new individualistic subjects ($n = 16$, additional sample I, first approach) and when further including the seven originally excluded subjects ($n = 23$, additional sample II, second approach).

In the first approach, we added four new individualistic subjects for analysis. We scanned seven new, individualistic subjects (four female; all right-handed; age, 20–22 years; range, 21.0 ± 1.00 years) by the same experimental procedure as in the main experiment except that, before the experiment, they were prescreened using the SVO questionnaire to ensure that they were individualistic. Two of them were excluded because we found that the signal-to-noise ratio (SNR) of their fMRI data was significantly lower than the SNR of the data reported in the main experiment (by one-sample *t* test, $p < 0.001$). Among the remaining five subjects, one subject's choice behavior was found to be insensitive to the others' reward and therefore this subject was excluded, as was done in the main experiment. Therefore, the remaining four individualistic subjects (three female; age, 20–22 years; 21.3 ± 0.96 years) were added to the

Table 5. Results in the samples with larger number of individualistic subjects ($n = 16$, and $n = 23$) replicate our main results ($n = 12$)

| | Main sample ($n = 12$) | Additional sample I ($n = 16$) | Additional sample II ($n = 23$) |
|--|--------------------------------|----------------------------------|-----------------------------------|
| Figure 6B | | | |
| Effect of covariate | $t_{(10)} = 0.086, p = 0.934$ | $t_{(14)} = -0.340, p = 0.738$ | $t_{(21)} = 0.435, p = 0.668$ |
| Effect of constant | $t_{(10)} = -2.623, p = 0.026$ | $t_{(14)} = -2.594, p = 0.021$ | $t_{(21)} = -2.207, p = 0.039$ |
| Difference with prosocial subjects | $t_{(31)} = 2.471, p = 0.019$ | $t_{(35)} = 2.470, p = 0.044$ | $t_{(42)} = 2.574, p = 0.017$ |
| Figure 7B | | | |
| [rAI \times signed effectivity of the other-bonus value] | $t_{(11)} = 2.325, p = 0.040$ | $t_{(15)} = 2.602, p = 0.020$ | $t_{(22)} = 3.811, p < 0.001$ |
| [dlPFC \times other-bonus value] | $t_{(11)} = -0.725, p = 0.484$ | $t_{(15)} = -0.111, p = 0.913$ | $t_{(22)} = 0.146, p = 0.885$ |
| 2 \times 2 repeated-measure ANOVA | $F = 6.278, p = 0.018$ | $F = 5.670, p = 0.023$ | $F = 8.248, p = 0.007$ |

With the samples with relatively larger number of individualistic subjects ($n = 16$, adding 4 new individualistic subjects, and $n = 23$, further including the 7 originally excluded subjects; for details, see Materials and Methods), we confirmed the findings from the original main samples ($n = 12$).

original 12 individualistic subjects and we then performed our analyses (additional sample I in Table 5).

In the second approach, we reanalyzed the data by adding some of the originally excluded subjects to the 16 subjects in the first approach (resulting in $n = 23$, additional sample II). We originally excluded seven subjects in the main experiment and one subject in the additional experiment because their choice behavior was not significantly positively influenced by the other-bonus, despite the intention of this study and our experimental design (our preexperimental procedure). More precisely, among the eight subjects, we found one subject's choice behavior was significantly negatively influenced by the other-bonus, so that subject was excluded in the following. All of the remaining seven subjects were individualistic according to the SVO classification and their behavior was not found to be influenced by the other-bonus in our task. Nevertheless, their social value conversion process may still occur and thus may possibly be detected in BOLD signals. We therefore added the seven subjects to the 16 individualistic subjects in the first approach and conducted our analyses (additional sample II in Table 5).

Bootstrap test. We used a bootstrap test to further support our statistical examinations, for instance, to guard against potential confounds such as outlier samples, in the analysis of both behavior and brain signals because the bootstrap procedure is distribution independent and thus expected to be less susceptible to different distribution shapes (DiCiccio and Efron, 1996). We primarily used bootstrap tests to estimate the distributions of group means (Efron and Tibshirani, 1993; Adèr et al., 2008). We used the following procedure: resample the original data (with replacement), repeat the resampling to obtain a set of bootstrapped samples of the same size as the original sample, repeat this resampling procedure 10,000 times to obtain 10,000 bootstrapped sets, and then compute their means to obtain the bootstrapped estimated distribution of the means. We used the estimated distribution to perform bootstrap hypothesis testing, finding the probability that the mass differs from zero (in a two-tailed manner). This approach is similarly applied to the case examining the difference between two original samples. For the correlation analysis, similar estimations were applied to the distribution of correlation coefficients.

Results

Our experiment comprised a main task and a control task (Fig. 1; for details, see Materials and Methods). The main task consisted of standard trials interleaved with other-bonus trials, whereas the control task consisted of standard and self-bonus trials (Fig. 1B). The subject repeatedly chose between two options in every trial to indicate his or her preference. In the standard trial, each of two options was associated with the subject's own probabilistic outcomes (Fig. 1B). An additional number was displayed below only one of the two options to indicate a bonus to others in other-bonus trials and a bonus to the self in self-bonus trials (see Materials and Methods). These bonuses would always be given when the option was chosen regardless of whether the probabilistic standard reward was given. The bonus magnitude was controlled to be moderate relative to the standard reward, so that concerns about reward allocation balance were minimized in this task.

These settings allowed us to measure the conversion with respect to proximal self-regarding decision making to avoid interpretative complications introduced from higher-order processes such as the balancing of reward allocation between the self and others.

Behavior

We used choice behavior in the standard trial as the reference condition and assessed modification by other-bonus and self-bonus to define their bonus values. Choice behavior was significantly modified by both the other-bonus and self-bonus, but to a lesser extent by the other-bonus (Fig. 2). Compared with standard trials that had a similarly sized value difference between the two options, the subject chose the option with a bonus more often than the other option in both other-bonus and self-bonus trials and this tendency became stronger as the magnitude of the bonus became larger (Fig. 2A). Given the same face amount, however, the choices were less strongly modified by the other-bonus than by the self-bonus. These observations were confirmed by an increasing trend in the extent of behavioral change at the point of indifference to bonus magnitude (Fig. 2B).

We quantified these observations by modeling the choice behaviors based on the differences in value between standard, self-bonus, and other-bonus reward, respectively (Eq. 1, $DV = \Delta V_S + w_S \Delta S + w_O \Delta O$ in Materials and Methods; Table 1 for model selection and comparisons with other models). The term for the standard reward was modeled as the standard value difference between the two options, with its probabilistic nature incorporated as risk dependency. Both the self- and other-bonus terms were inserted to the corresponding option with corresponding weights, w_S and w_O , respectively, where the subscripts "S" and "O" indicate self and other, respectively. By fitting to the individuals' behavior, we found that the estimated weights of both bonuses were significantly larger than zero (Fig. 2C; statistics are given in the legend). The weight of the other-bonus was significantly smaller than the weight of the self-bonus (Fig. 2C). In addition, we confirmed that the choice behavior was not affected by the mere presence of a bonus and that it was not influenced by accumulated payment of the others' reward (see "Model fitting and selection" section). In the following BOLD analyses, the weights estimated from the behavioral regression were used to define various forms derived from the self- and other-bonus values.

Neural activation for converting others' reward to decision signals

We then analyzed BOLD signals using a whole-brain GLM approach to examine the value conversion processes in three stages for signals underlying the decision (i.e., decision value, DV) and for signals relating to two forms of the other-bonus and self-bonus value (their offered and effective values).

First, we examined brain activations according to the DV. We found significant activation by the DV in the mPFC and right middle temporal gyrus (rMTG) (Fig. 3A, Table 3; $[x, y, z] = [0, 59, 4]$ and $[60, -4, 2]$, respectively (coordinates are in Talairach space and $p < 0.05$, FWE corrected; for details, see Materials and Methods). We confirmed that the mPFC activation was enhanced with increasing DV for all three trial types (Fig. 3B). These two findings on the mPFC responses (corresponding to Fig. 3A, B) were also confirmed in a variant of the GLM, in which we changed the regressor of interest from the DV to the choice probability (center of activation, $[0, 56, 1]$ in Talairach space, cluster size = 238, and the increased mPFC activations with increasing choice probability for all three trial types, by Page's test, standard trial, $p = 0.001$; self-bonus trial, $p = 0.009$, other-bonus trial, $p = 0.024$). These results indicate that the mPFC responses contained the brain signals for making behavioral choices in our task, which is generally consistent with previous findings on activation in the mPFC and also in the ventromedial PFC (vmPFC) (Rushworth and Behrens, 2008). We refer to this as mPFC activation rather than vmPFC activation to err on the side of caution, partly because the peak's z -coordinates was 4, higher than the AC-PC line; however, we note that the activation might well be regarded as being in the vmPFC activation because our activated were extended to the vmPFC region (Fig. 3A). Indeed, the z -coordinate of peak signals close to our activation have often been called vmPFC responses (Phelps et al., 2004; Tobler et al., 2007; De Martino et al., 2013; Liljeholm et al., 2015).

For the offered values (e.g., w_{O} for other-bonus), we found that the rTPJ uniquely encoded the other-bonus value, whereas the ldlPFC was commonly activated by both the other- and self-bonus value. First of all, the BOLD responses were significantly correlated with the other-bonus value in the rTPJ and ldlPFC ($[57, -61, 34]$ and $[-39, 44, 10]$, respectively; Fig. 3C) and some other brain regions (Table 3), including the dorsomedial prefrontal cortex (dmPFC, $[9, 17, 49]$). The BOLD responses that significantly correlated with the self-bonus value were found in the ldlPFC ($[-27, 50, 16]$), right dorsolateral prefrontal cortex (rdlPFC) ($[39, 56, 16]$), dmPFC ($[9, 26, 37]$), and several other brain regions (Table 3). No brain region showed significant activation correlated to the covariates of either the offered other- or self-bonus value.

Using a contrast analysis of the other-bonus with the self-bonus (other > self), we found that only the rTPJ activation was unique to the other-bonus value ($[58, 58, 30]$; Table 3). No brain region was activated by the contrast of the self-bonus with the other-bonus (self > other). However, because there was a notable difference in the above findings in the ldlPFC and rdlPFC activations, we further conducted contrast analysis in their respective

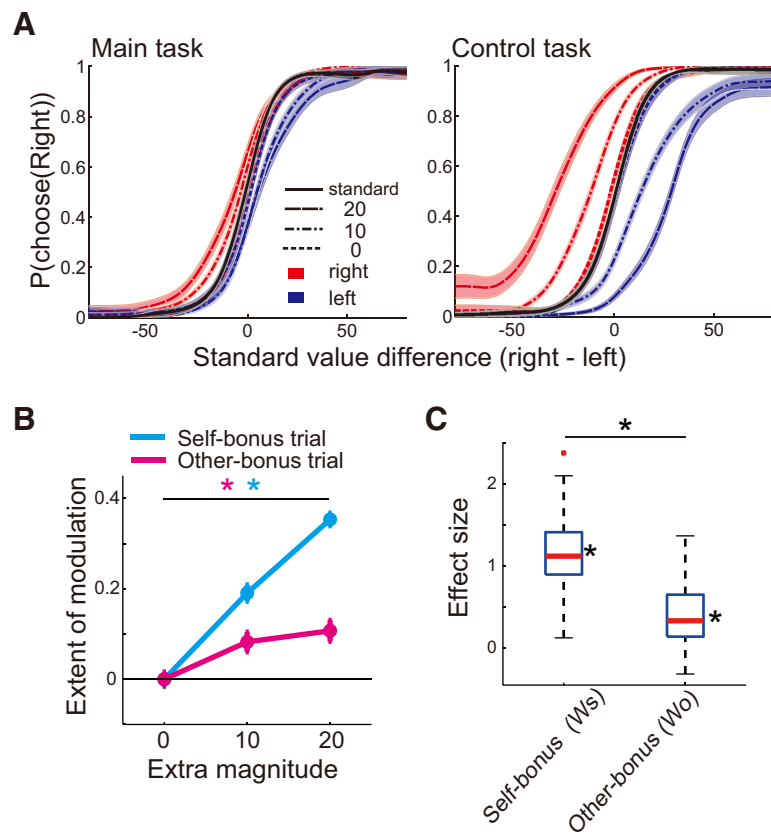


Figure 2. Behavioral results. **A**, The probability of choosing the option on the right side as a function of the standard value difference (ΔV_s , right minus left) in the main and control tasks (left and right panels, respectively); the groupwise mean (the shaded region indicates the SE) is shown separately in each type of trial for different reward magnitudes (different line styles) and when the bonus was attached to the right or left option (different colors). **B**, Extent of behavioral change (mean and SE) in the other-bonus and self-bonus trials (magenta and cyan lines, respectively) versus reward magnitudes. The extent of change is derived by taking the difference in choice probability between the corresponding sigmoidal choice curve and that of the trial with zero additional reward at the indifference point ($\Delta V_s = 0$). *Significant increasing trend by Page's test: for, self-bonus trials, $p < 0.001$ and for other-bonus trials, $p < 0.001$. **C**, Behavioral weights from model fitting. * $p < 0.001$, significantly larger than zero or significantly different by group-level t tests; self-bonus: $w_s = 1.225 \pm 0.471$, $t_{(42)} = 17.041$, $p < 0.001$; other-bonus: $w_o = 0.388 \pm 0.367$, $t_{(42)} = 6.939$, $p < 0.001$; difference, $t_{(42)} = 10.824$, $p < 0.001$, paired t test. For the box-plots, red lines in the boxes indicate the medians; the box limits indicate the top and bottom quartiles; the length of each whisker indicates 1.5 times the interquartile range; and the red dot indicates an outlier.

ROIs and confirmed that the ldlPFC region had significant activation by only the other-bonus value, whereas the rdlPFC region had significant activation by both the other-bonus and the self-bonus value. Using ROIs derived from the activation shown in Figure 3D (Table 3), we found the effect sizes for offered other-bonus value were significantly larger than zero ($t_{(35)} = 6.272$, $p < 0.001$) in the ldlPFC but not in the rdlPFC ($t_{(35)} = 1.254$, $p = 0.218$), with a significant difference between the two sides ($t_{(35)} = 3.068$, $p = 0.004$). By contrast, the effect sizes for offered self-bonus value were significantly larger than zero for both sides of the dlPFC (left: $t_{(35)} = 3.798$, $p < 0.001$; right: $t_{(35)} = 3.390$, $p = 0.002$), with no significant difference between them ($t_{(35)} = 0.067$, $p = 0.947$). Moreover, using a conjunction analysis (Table 3; Nichols et al., 2005), we found that the ldlPFC activation was common to the offered other- and self-bonus values ($[-44, 43, 21]$), and also activations of several other brain regions including the dmPFC ($[9, 17, 49]$), precuneus ($[9, -70, 37]$), and inferior parietal lobule (IPL; right, $[36, -64, 40]$; left, $[-42, -49, 40]$).

We then addressed the issue of whether these findings might be due to the difference in the size of the behavioral effect rather than the difference between rewards to others versus the self,

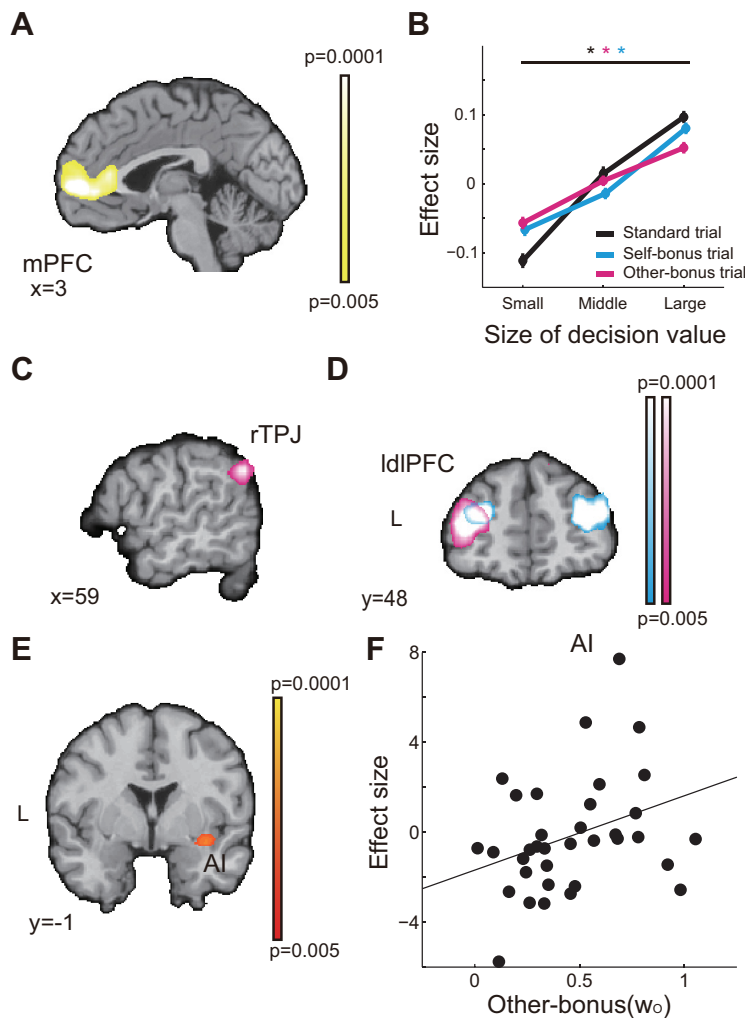


Figure 3. Signals for decision value, other-bonus value, and self-bonus value. **A**, Activation by DV (yellow) in the mPFC, which survived at the $p < 0.05$ level, FWE (whole brain)-corrected by a nonparametric method. For display, individual voxels in the activation map were first subjected to thresholding at $p < 0.005$ and then blurred at the resolution of anatomic imaging. The activation figures hereinafter were made in the same way unless stated otherwise. **B**, BOLD signal change (mean and SE) in the mPFC (extracted based on the ROI in **A**) in standard, self-bonus, and other-bonus trials as a function of the size of the DV, normalized within each trial type. *Significant increasing trend by Page’s test: standard trial, $p = 0.002$; self-bonus trial, $p < 0.001$, and other-bonus trial, $p = 0.008$. **C**, rTPJ activation by offered other-bonus value (magenta). **D**, Activations by both other-bonus value (magenta) and self-bonus values (cyan) in the ldlPFC. **E**, rAI response to signed effective other-bonus value signals is significantly positively correlated with the behavioral weight for the other-bonus (w_o^i), which survived at the $p < 0.05$ level, FWE (whole brain)-corrected. **F**, Groupwise scatterplot between the weight of other-bonus and the rAI response for signed effective other-bonus value signals in the rAI ROI, based on leave-one-out cross-validated ROI analysis. The result showed a significant positive correlation, $r = 0.384$, $p = 0.021$, $n = 36$; $p < 0.017$ by bootstrap test. Each dot corresponds to a subject. See Table 3 for a summary of GLM activations.

because the self-bonus has a larger effect than the other-bonus for a given face amount. We addressed this issue for the rTPJ and ldlPFC activations by reanalyzing the BOLD signals with the subjects and trials selected to control for the size of the behavioral effect (see “Analysis of rTPJ and ldlPFC responses while controlling the behavioral effects” section). Even after controlling for the behavioral effect size, we confirmed that the activation in the rTPJ was unique to the offered other-bonus value and the ldlPFC was common to the offered other- and self-bonus values.

For the second stage, we defined a variable called signed effective value that links the offered value in the first stage to the final choice in the third stage, representing the effective impact of the offer on the choice. We note that even the same offer could have different impacts on the choice, dependent on the standard value

difference, due to the nonlinearity of the logit function connecting a final decision value to choice probability (for details, see “Signed effective value” section). This effect could be assessed using the (choice) effectivity, which is defined as the derivative of the choice probability with respect to the standard value difference (Eq 2). The signed effective value was then formulated as the product of the signed value and effectivity (Eq. 3). The sign of the signed effective value indicates the choice. We found that rAI activation was related to the signed effective other-bonus via the weight of the other-bonus as a covariate in the second-level analysis. Indeed, for both the self-bonus and other-bonus, we did not find any significant activation by the signed effective value alone. However, we observed a significant activation of the rAI for only the other bonus via the behavioral weight w_o as a covariate (Fig. 3E: whole-brain voxelwise activation, $[39, -4, -11]$; Fig. 3F: ROI analysis with leave-one-out cross-validation; Table 3). Therefore, we found that the rAI activations of the signed effective value covary with the behavioral weights across subjects, for instance, suggesting that a subject who places greater weight on others’ value in making choices (i.e., with the larger behavioral weight) would have a stronger rAI activation by the signed effective value. To confirm that the rAI activation was due to the other-bonus only, we also verified that there was no significant correlation between the effect sizes of behavioral weight (w_o^i) and the signed effective self-bonus value in the ROI of the rAI ($r = -0.081$, $p = 0.638$). These results suggest that the activation of the rAI by the signed effective other-bonus value reflects or correlates with individual differences in converting the others’ reward into decisions.

In summary, we found brain signals for three stages of social value conversion: the offered other-bonus value especially in the rTPJ and the ldlPFC, the signed effective other-bonus value covarying with the weight for the other-bonus in the rAI, and the DV in the mPFC.

Coupling of rTPJ and ldlPFC to rAI responses and rAI to mPFC responses

We next conducted PPI analyses to determine whether brain activation patterns followed our hypothesized processing of the three stages. We probed for corresponding interactions from the offered value to the signed effective value, and from the signed effective value to the DV. For a control, we probed for an interaction from the offered value directly to the DV. We report activation when both of the following two types of PPIs indicated significant activation (unless explicitly stated otherwise): ROI PPI, in which the target region is based on the activations from

our GLM results, and voxelwise PPI, which visualizes the activated region and provides additional confirmation of the ROI PPI (in the following, all results by voxelwise PPI are at the $p < 0.05$ level, FWE-corrected unless stated otherwise). Care must be taken regarding several points in the PPI analysis. Importantly, we should note that PPI analysis assesses correlation rather than directionality or causality. Therefore, to complement our PPI analyses, we also conducted DCM analysis as a second, independent method of analysis.

First, in the interaction of the offered value to the signed effective value, we found that both the rTPJ and ldlPFC responses had an impact on rAI responses, modulated by the (unsigned) effectivity of the other-bonus value. We chose the unsigned effectivity as a psychological seed because it quantifies possible impacts of a same offer on different decisions and thus helps in tracing social value conversion; this is unsigned because both the rTPJ and ldlPFC responses were unsigned (not aligned on choices), and effectivity is used because it is readily available given the absolute standard value difference between two options without a final decision (DV). The impact of an offer may not be appreciable when the standard value difference is large (i.e., when the unsigned effectivity is small) because the decision does not change much in the first place, whereas the impact may be appreciable when the standard value difference is small. Therefore, we performed a PPI analysis of [rTPJ \times effectivity of the other-bonus value] and [ldlPFC \times effectivity of the other-bonus value] and found that the rAI activations were significantly correlated with both of the interaction terms (ROI PPI: rTPJ seed, $t_{(35)} = 2.320$, $p = 0.026$, $p = 0.005$ by bootstrap test, [30, -10, -8] in voxelwise PPI; ROI PPI: ldlPFC seed, $t_{(35)} = 2.171$, $p = 0.037$, $p = 0.013$ by bootstrap test, [48, 5, -11] in voxelwise PPI; Fig. 4A, Table 4). We also conducted the same PPI analysis with other brain regions that had the activations by the offered value, namely, the dmPFC, IPL and precuneus, and found that none of them had a significant impact on the rAI responses.

Second, for the signed effective value to the DV, we found that rAI responses had an impact on mPFC responses, modulated by the signed effectivity of the other-bonus value. We chose this as a psychological seed because the rAI responses were already signed. Therefore, using the interaction term [rAI \times signed effectivity of the other-bonus value], we found the mPFC activations to be significant (ROI PPI: $t_{(35)} = 2.627$, $p = 0.013$, $p = 0.012$ by bootstrap test, [1, 57, 4] in voxelwise PPI; Fig. 4B, Table 4).

Third, we examined as a control whether the responses in the ldlPFC, rTPJ, and other regions that had activations by the offered

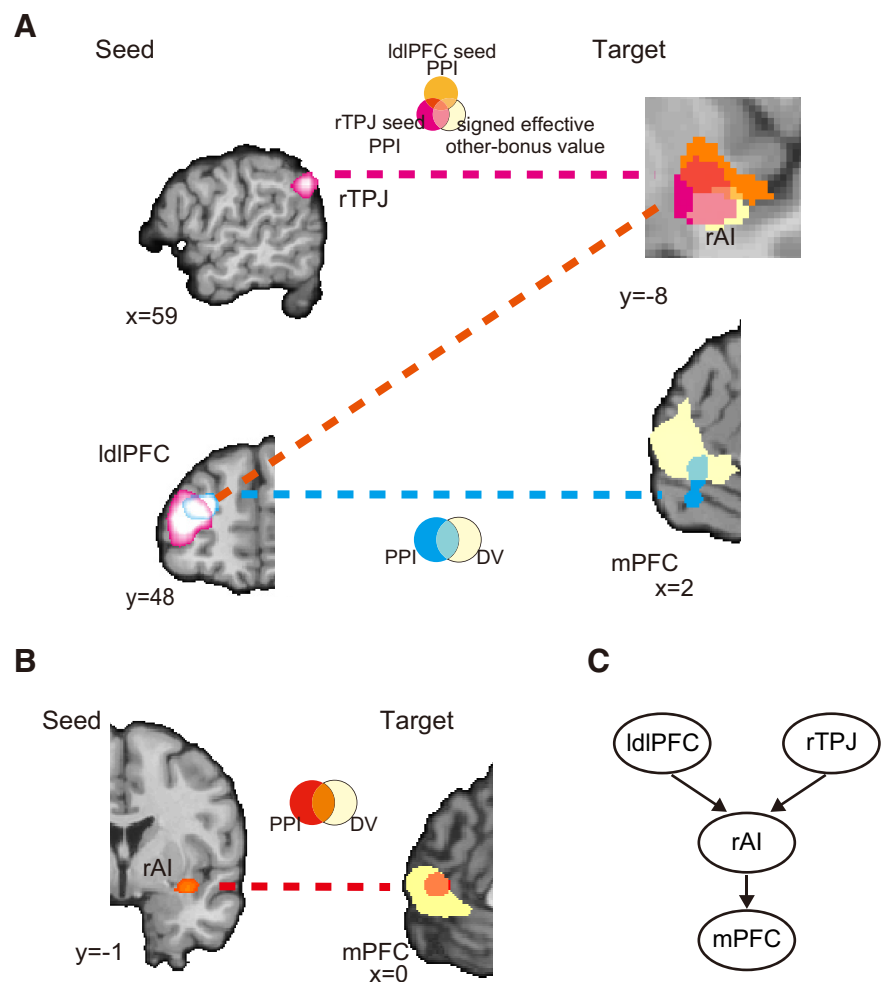


Figure 4. PPI and DCM results for all subjects. **A**, PPI results for the rTPJ and ldlPFC on the rAI and mPFC. Coupling with the rTPJ and ldlPFC responses (physiological seeds, derived from activation by the offered other-bonus value, Fig. 3C and D) as a function of the effectivity of the other-bonus value in the rAI responses (magenta, rTPJ seed; orange, ldlPFC seed) overlapped with the rAI activations (yellow, from Fig. 3E). Coupling with ldlPFC responses (physiological seed, derived from activation by the offered self-bonus value) as a function of offered self-bonus in the mPFC responses (cyan) overlapped with activation by the DV (yellow, from Fig. 3A). For display purposes, the activation identified by voxelwise PPI is shown as a whole cluster (uncorrected $p < 0.005$), including not only the activated voxels within the target ROI but also those in the surrounding region, and the significance of the activation is tested using SVC for only the activated voxels within the target ROI (all voxelwise PPI activations are at the $p < 0.05$ level, FWE [SVC]-corrected). The same methods for the SVC test and for display were used for Figure 7A. **B**, PPI results for the rAI to the mPFC. Coupling with rAI responses as a function of the signed effectivity of the other-bonus value on mPFC responses. Activated clusters identified by voxelwise PPI (red), indicated by the target in the panel, overlapped with the activation by the DV (yellow, derived from Fig. 3A). See Table 4 for a summary of PPI activations. **C**, DCM analysis supported three-stage processing; the structure shown in the figure was chosen as the best model by comparing exceedance probabilities among other models that had reversed connections over the brain regions (also see Fig. 5). Our DCM analysis included driving and modulatory inputs (derived from our findings for regional activations and connectivity, respectively), but for visibility, they were not shown in this and following figures (see Materials and Methods).

value (dmPFC, IPL, and precuneus) directly affected the mPFC responses. We used the interaction term derived from the respective response time series (physiological seed) and the offered other-bonus value (psychological seed). For these (seed region \times offered other-bonus value), we found no significant impact of the ldlPFC, rTPJ, or any of the other regions' responses on the mPFC responses for the other-bonus value (Table 4). Indeed, when we investigated whether the responses of all those regions directly affected the responses in the rMTG, which was the other region that had the significant activation by the DV, we found no significant impact from any of those seed regions' activations.

These results support our hypothesis of three stages of social value conversion from the offered value to the effective value and

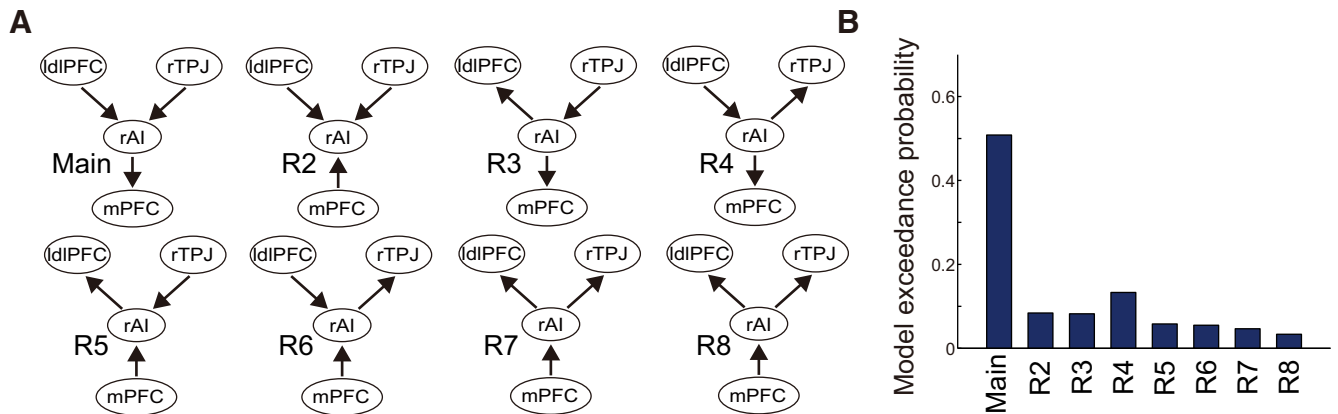


Figure 5. DCM model with feedforward connections compared with models having reversed connections. **A**, Main, the model with feedforward connections representing the three-stage processing for social value conversion shown in Figure 4C; R2–R8, alternative models that have possible reversed (feedback) connections between ROIs. **B**, Exceedance probabilities obtained by using random-effect Bayesian model selection indicated that the Main model is better than all of the models with the reversed connections (called the reversed family) at explaining the data.

then to the DV; in neural terms, from the rTPJ and ldlPFC responses to the rAI responses and then to the mPFC responses. For further verification, we used DCM analysis to examine the directionality of the interactions among the brain areas of the three computational stages. The main model with feedforward connections (Fig. 4C) was selected as the best model compared with other models involving opposite connections (called the reversed family; Fig. 5), thus supporting the notion that the three-stage processing is the dominant interaction among these areas in this behavioral task.

Finally, to clarify the relationships of the above findings with the processing of the self-bonus value, we examined additional PPIs. Because we did not find any significant activations by the signed effective value of the self-bonus, we examined the impact of the offered value on the DV; in other words, we investigated whether the responses to the offered value in the regions that had the significant activations common to both the other- and self-bonus values (namely, the ldlPFC, dmPFC, IPL, and precuneus) directly altered the mPFC responses. We found significant impacts by the ldlPFC and dmPFC responses, but by no other regions. Using the interaction term (ldlPFC \times offered self-bonus value), we found the mPFC activation to be significant (ROI PPI: $t_{(35)} = 2.063$, $p = 0.047$, $p = 0.039$ by bootstrap test, [9, 47, -5] in voxelwise PPI at the $p < 0.05$ level, FWE-corrected; Fig. 4A, Table 4). Using [dmPFC \times offered self-bonus value], we found the mPFC activation to be significant (ROI PPI: $t_{(35)} = 2.617$, $p = 0.013$ and also $p = 0.005$ by bootstrap test, [9, 62, 4] in voxelwise PPI at the $p < 0.05$ level, FWE-corrected; Table 4). These results indicate that the offered self-bonus value encoded in the respective ldlPFC and dmPFC responses had a direct impact on the mPFC responses.

rAI activation differs between prosocial and individualistic subjects

We tested the hypothesis that variability in the value conversion process may underlie different sociobehavioral isotypes. To do this, we investigated whether relative differences in this process underpin one example of sociobehavioral isotype, namely, SVO, which is a social preference about resource allocation between the self and others. Differences in SVO give rise to different behavioral phenotypes, including prosocial and individualistic (selfish). Prosocial people prefer to maximize the resource gain for both the self and others, whereas individualistic people prefer to maximize gains for themselves (Van Lange et al., 1997; Van

Lange, 1999). Through a postexperiment SVO questionnaire, we found that 21 subjects were prosocial, 12 were individualistic, and three were unclassified, largely consistent with the results of previous studies (Van Lange et al., 1997; Haruno and Frith, 2010). Although their responses to the SVO questionnaire clearly differed, the prosocial and individualistic groups behaved similarly within our task. The behavioral weight for the other-bonus was not significantly different between prosocial and individualistic subjects ($t_{(31)} = 0.870$, $p = 0.391$; Fig. 6A, details in the legend). Two additional tests also supported the finding (Table 2). It is important to note that the aims are quite different between our experimental task and the SVO. Our task was designed to examine social value conversion in proximal self-oriented decision making (in which the other-bonus was given in addition to the standard reward to the self in the same option), whereas the SVO aims to capture individual differences in preference for balancing reward allocation between the self and others.

To explore possible between-group differences of brain signals in social value conversion, we first examined the activations related to the other-bonus values in the rTPJ, ldlPFC, and rAI (Fig. 3) by entering the SVO classification variable (a binary indicating either prosocial or individualistic) as another covariate for the ROI GLM second-level analysis of the corresponding other-bonus value variable, in addition to the original covariate (weight for the other-bonus). We found a notable difference in only the rAI, showing significant effects of both the covariates (weight for the other-bonus: $t_{(33)} = 3.517$, $p = 0.001$; SVO classification: $t_{(33)} = 2.142$, $p = 0.040$ mPFC, but no covariate effect on the ROIs activated either by offered other-bonus value in the rTPJ and ldlPFC, or by DV in the mPFC.

For further examinations, we reanalyzed the rAI activation separately for each SVO group by using the signed effective value of the other-bonus together with the original covariate (weight for the other-bonus). We found that the rAI activation worked in an opposing manner between the two groups (Fig. 6B). In our original all-subject analysis, the rAI was significantly activated by the signed effective other-bonus value through the weight for the other-bonus (a covariate in the second-level analysis), but not by signed effective other-bonus value as a constant in the second-level analysis. In a comparison between the groups, we found that the significant effect of the weight for the other-bonus as a covariate remained for the prosocial subjects ($t_{(19)} = 2.350$, $p = 0.030$) but not for the individualistic subjects ($t_{(10)} = 0.086$, $p = 0.934$), with a marginal difference between the groups ($t_{(31)} = 1.570$, $p =$

0.063). By contrast, the constant had a significant negative effect in the individualistic subjects ($t_{(10)} = -2.623, p = 0.026$), whereas the effect tended to be positive but not significant in the prosocial subjects ($t_{(19)} = 1.493, p = 0.152$); there was a significant difference between the groups ($t_{(31)} = 2.471, p = 0.019$). These results suggested an opposing influence of the rAI activation on decision making between the groups, such that the increased rAI activation would promote or suppress the choice of the option with the other-bonus.

In the above and following analyses, one might wonder about the validity of the results for the individualistic subjects because the number of the subjects was relatively small ($n = 12$). To address this concern, we conducted two additional analyses (Table 5; for details, see Materials and Methods) and obtained essentially the same results with larger samples, when adding four new individualistic subjects ($n = 16$) or when further including the seven originally excluded subjects ($n = 23$).

ldlPFC and rAI coupling to the mPFC in prosocial and individualistic subjects

We next examined differences between the two groups in the coupling between these activations, by applying all the PPI analyses described above to each group separately. We found that impacts of the rAI and ldlPFC on the mPFC responses differed between the groups (but the impact of the TPJ responses did not). The mPFC responses were significantly modulated by the rAI response only in the individualistic subjects ([rAI \times signed effectiveness of the other-bonus value]), but by the ldlPFC response only in the prosocial subjects ([ldlPFC \times offered other-bonus value]), supported by voxelwise PPI (Fig. 7A, Table 4) and by ROI PPI analysis (Fig. 7B). To further examine the between-group difference, we performed a 2×2 repeated-measures ANOVA in the ROI PPI analysis (two levels of groups [prosocial, individualistic] by 2 levels of seed regions [rAI, ldlPFC]). We found a significant interaction ($F = 6.278, p = 0.018$; Fig. 7B) but no significant main effects ($F = 0.010, p = 0.923$ for prosocial vs individualistic, $F = 1.544, p = 0.223$ for rAI vs ldlPFC). These results suggest a stronger impact of the rAI and ldlPFC responses on the mPFC responses for the individualistic and prosocial groups, respectively.

This characteristic remained evident even after controlling for the responses related to the self-bonus value (Fig. 7C; Materials and Methods). First, the impacts of the rAI and ldlPFC responses on the mPFC response for the self-bonus value did not significantly differ between the two groups. As shown in Figure 7C, the prosocial subjects responded similarly in cases with the other- and self-bonus values, whereas the responses of individualistic subjects in the other-bonus case were quite different from those in the remaining three cases (i.e., the individualistic subjects with the self-bonus value, and the prosocial subjects with the self-bonus value and the other-bonus value). Using ROI PPI [ldlPFC \times self-bonus value], we found that the activation was significant in the prosocial subjects ($t_{(20)} = 2.841, p = 0.010$ and

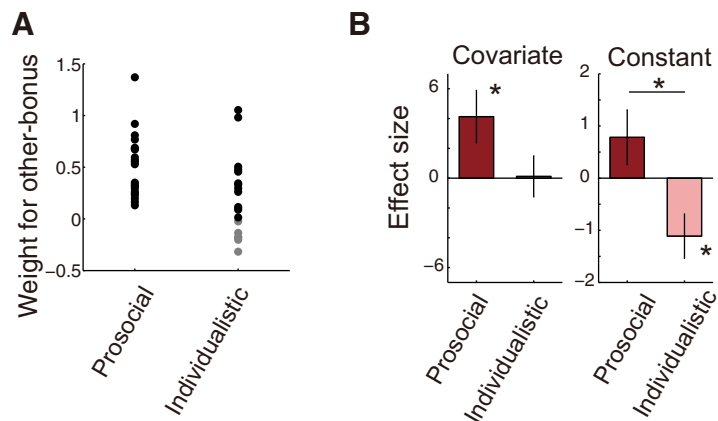


Figure 6. rAI results for the prosocial and individualistic groups. **A**, Behavioral weight for the other-bonus distribution in each of the individualistic and prosocial subjects. Each dot corresponds to a subject. The weight for the other-bonus was not significantly different between prosocial and individualistic subjects ($t_{(31)} = 0.870, p = 0.391$). However, for clarification, stars indicate the seven subjects whose choice behavior was insensitive to the other-bonus in our task (see “Model selection” section). The difference became significant ($t_{(38)} = 2.758, p = 0.009$) when including the gray dots, that is, when including the seven subjects. This is reasonable because all seven subjects were also insensitive to others’ reward in the SVO and were therefore classified as individualistic, thus lowering the average weight for the other-bonus of the individualistic subjects in this case. **B**, Reanalysis of rAI activation by signed effective other-bonus value (as indicated in Fig. 3E) separately for the prosocial and individualistic subjects (brown and pink); therefore, as a second-level analysis, multiple regression was performed on the rAI β sizes using the regressors of the constant (corresponding to the signed effective other-bonus value) and behavioral weights (covariates) in each SVO group; left, effect size of the covariate term; right panel, effect size of the constant term (error bar: SE); * $p < 0.05$, significantly larger than zero for each of covariate and constant effect and also significantly different between the two groups by group-level t tests.

$p = 0.005$ by bootstrap test) and marginally significant in the individualistic subjects ($t_{(11)} = 1.876, p = 0.086$, but significant by bootstrap test $p = 0.039$), and this was also confirmed as significant by voxelwise PPI (Table 4). These observations were further supported by two types of circular statistical analysis of the directional change in the 2D map at both the group and individual levels (for details, see “Circular statistical analysis for SVO PPI” section).

Our DCM analysis further supported these results. First, we examined by DCM analysis whether the subjects in each SVO group had differences in connectivity from the ldlPFC to the mPFC (Fig. 7D). Comparing exceedance probabilities in Bayesian model selection (Stephan et al., 2009), we found that the three-stage processing together with the connection from the ldlPFC to the mPFC was dominant in the prosocial subjects, whereas the three-stage processing alone was dominant in the individualistic subjects (Fig. 7E). Next, using Bayesian model averaging, we examined the connectivity strength in each SVO group (for connectivity from the ldlPFC to the mPFC, and from the rAI to the mPFC). First, we found that the ldlPFC to mPFC connectivity was significantly positive in only the prosocial subjects ($t_{(20)} = 2.596, p = 0.017, p < 0.001$ by bootstrap test; in the individualistic subjects, $t_{(11)} = 1.104, p = 0.293, p = 0.240$ by bootstrap test). Second, we found that the rAI to mPFC connectivity was marginally significant in the individualistic subjects by t test ($t_{(11)} = 2.034, p = 0.067$), but significant by bootstrap ($p < 0.001$), which we take to suggest overall consistency with our PPI results. Interestingly, this connectivity was also significant in the prosocial subjects ($t_{(20)} = 2.630, p = 0.016, p < 0.001$ by bootstrap test), which notably differs from the PPI results above (Fig. 7B; see Discussion).

We did not find significant between-group differences in the rTPJ’s impact on the mPFC, or in the rTPJ’s and ldlPFC’s impacts on the rAI. First, the PPI effects of rTPJ ([rTPJ \times offered other-bonus value]) on the mPFC was not significant within each group

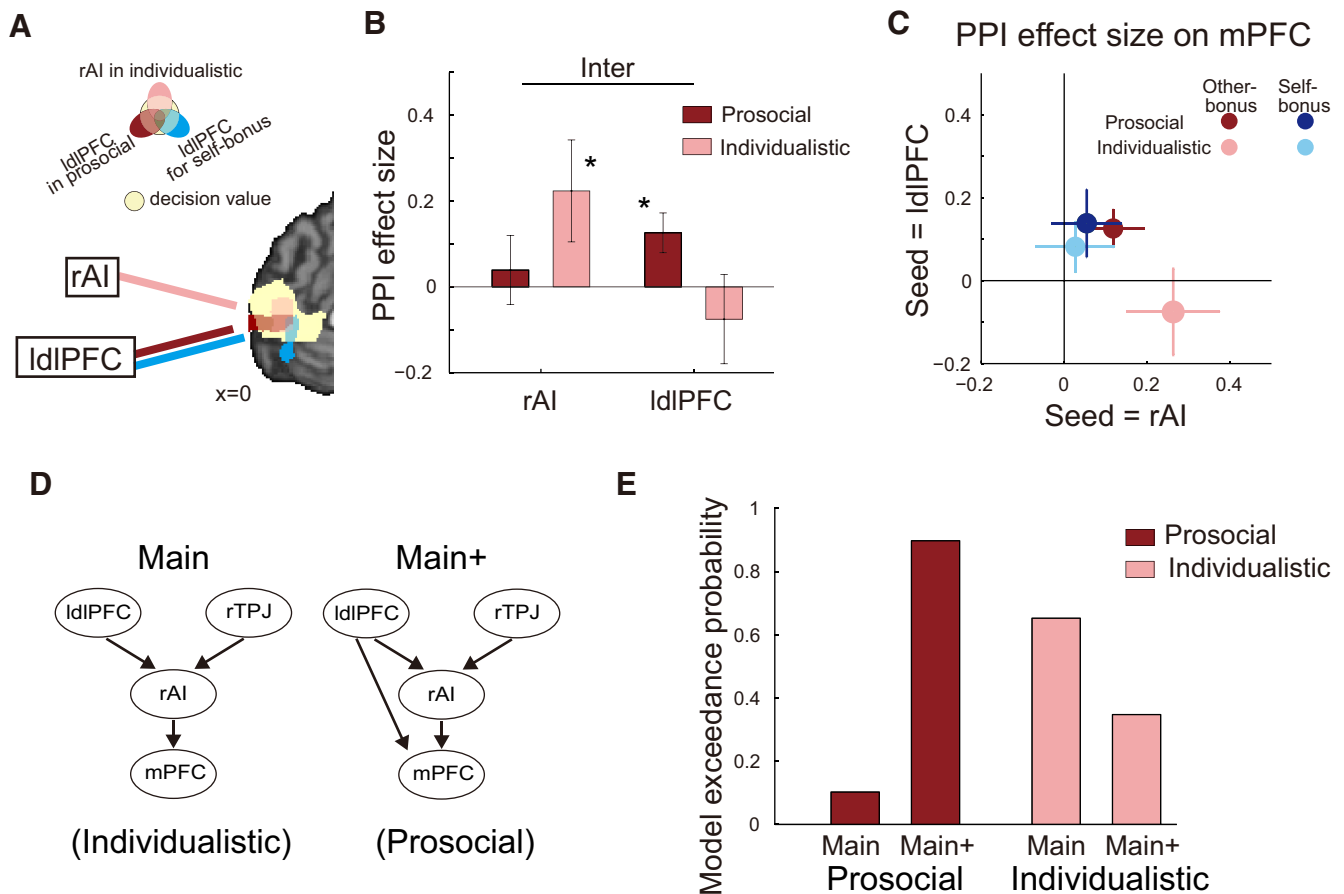


Figure 7. PPI and DCM results for the prosocial and individualistic groups. **A**, PPI activations in the mPFC (by voxelwise analysis) overlapped with activation for the DV (yellow, from Fig. 3A): Pink, [rAI × signed effectivity of other-bonus value] for the individualistic group. Brown, [ldIPFC × offered other-bonus value] for the prosocial group. Cyan, [ldIPFC × offered self-bonus value] for both groups (similar to Fig. 4A, but with three unclassified subjects excluded). All voxelwise PPI activations are at the $p < 0.05$ level, FWE (SVC)-corrected. **B**, ROI PPI effect sizes in the mPFC indicated by mean and SE. *Significant effect of [rAI × signed effectivity of other-bonus value] for the individualistic group (pink, $t_{(11)} = 2.325, p = 0.040$ by bootstrap test, $p = 0.008$, but n.s. for the prosocial group, brown, $t_{(20)} = 1.554, p = 0.136$, bootstrap $p = 0.096$), and [ldIPFC × offered other-bonus value] for the prosocial group ($t_{(20)} = 2.722, p = 0.013$, by bootstrap $p = 0.004$, but n.s. for the individualistic group, $t_{(11)} = -0.725, p = 0.484$, by bootstrap $p = 0.449$). “inter,” Significant interaction effect ($F = 6.278, p = 0.018$) by 2×2 repeated-measures ANOVA for the SVO groups with the two brain regions (rAI and ldIPFC). **C**, Comparison of the impact between the two SVO groups, relative to that for the self-bonus. Effect sizes from PPI results for mPFC responses are plotted on a 2D map (mean and SE, by ROI PPI), separately for prosocial and individualistic subjects for the other-bonus and self-bonus values (reddish and greenish, respectively), with the rAI and ldIPFC cases shown on the horizontal and vertical axes. **D**, Based on the main model of our DCM analysis (left), a new DCM model was constructed with connectivity from the ldIPFC to the mPFC (right, indicated by Main+). Words in parentheses at the bottom are shown only for the ease of understanding in correspondence to the results in **E**. **E**, Exceedance probabilities obtained by using random-effect Bayesian model selection to indicate which of the two models fit better each SVO group; the model with dlPFC–mPFC connectivity (Main+) was better for the prosocial group, whereas the original main model (Main) was better for the individualistic group.

or between groups (individualistic: $t_{(11)} = -0.489, p = 0.634$, bootstrap $p = 0.600$; prosocial: $t_{(20)} = 1.217, p = 0.238$, bootstrap $p = 0.215$; difference: $t_{(31)} = 1.333, p = 0.192$, bootstrap $p = 0.278$). Second, for the PPI on rAI responses by the rTPJ ([rTPJ × choice effectivity of other-bonus value]) and by the ldIPFC ([ldIPFC × choice effectivity of other-bonus value]), we found no significant between-group difference by the rTPJ ($t_{(31)} = 0.606, p = 0.549$, bootstrap $p = 0.649$, independent t test) or by the ldIPFC ($t_{(31)} = 0.995, p = 0.328$, bootstrap $p = 0.422$), and no significant main effects or interaction by 2×2 repeated-measures ANOVA (prosocial vs individualistic, $F = 0.717, p = 0.403$; rTPJ vs ldIPFC, $F = 0.842, p = 0.366$, interaction $F = 0.608, p = 0.442$). We note that the effect of the rTPJ on the rAI was significant in the prosocial ($t_{(20)} = 3.163, p = 0.005$, bootstrap $p = 0.002$) but not in the individualistic subjects ($t_{(11)} = 1.446, p = 0.176$, bootstrap $p = 0.078$), but this finding might be attributable to the relatively small number of the original individualistic subjects ($n = 12$) rather than a significant difference between the two groups, because the effect became significant with

a larger sample size ($n = 23, t_{(22)} = 1.842, p = 0.079$, bootstrap $p = 0.024$).

These results suggest that the prosocial subjects have a greater relative reliance on the ldIPFC activation by the other-bonus value affecting the mPFC response similarly to the self-bonus value (the process of which is common between the two groups), whereas the individualistic subjects have a greater relative reliance on the rAI activation by the other-bonus value affecting the mPFC response.

Before proceeding to the Discussion, we have a cautionary remark on roles played by feedback processing. This study focused on feedforward processing of social value conversion through the areas identified above, however, these should not be considered as indicating a strict feedforward process for the conversion. Indeed, we consider it possible to involve feedback interactions among these areas. For this caution, we expanded the main DCM model (Figs. 4C, 5A) to a family of models having additional feedback connections (Fig. 8) and then found that the Main model family with feedback connections was relatively bet-

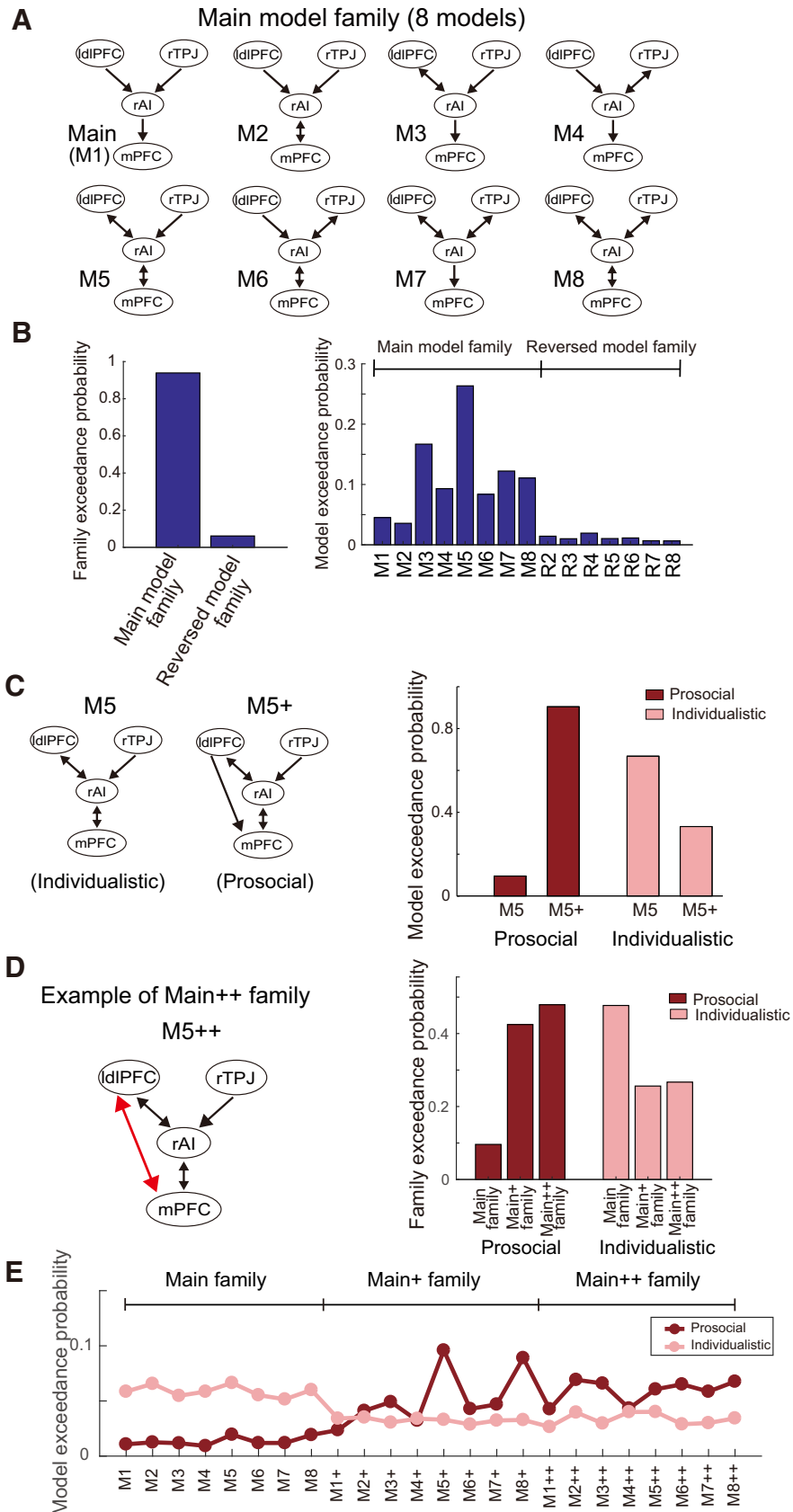


Figure 8. DCM models and results for difference between SV0 groups. **A**, Main model family (eight models): M1, the main model representing the three-stage processing for social value conversion as shown in Figure 4C, and other 7 models (M2–M8) with possible bidirectional connections. **B**, Across all subjects, comparison of exceedance probabilities between the main model family vs the reversed model family (left) and comparison of all models of the two families (right). The reversed model family was composed of 7 models (R2–R8) with reversed connections (see Fig. 5A). The dataset is much better explained by the main model family and by model M5 in particular. **C**, Left, As in the analysis in Figure 7D and E, a new DCM model (Figure legend continues.)

ter at explaining the data for both all subjects (Fig. 8B) and each SVO group (Fig. 8C,D). First, for all subjects, we expanded the main DCM model (Fig. 4C) to a family of models having different single bidirectional connections (the Main model family in Fig. 8A). We found that the Main family still had a better fit than the reversed family (Fig. 8B, left) and, within the Main family, models with bidirectional connections between the rAI and the mPFC (e.g., M5, Fig. 8B, right) were relatively better than the other models at explaining the data. Second, we merged this best model (M5) with the best model of each SVO group (Figs. 7D, 8C, left) and found that the respective models (Main and Main + family, e.g., M5 and M5+) better accounted for the individualistic and prosocial subjects, respectively (Fig. 8C, right). We also performed the Bayesian model averaging test for this case and the case below Figure 8D, and the results were similar to the original analysis. Further verification involved another family, called the Main++ family, that had a bidirectional dlPFC-mPFC connection (e.g., M5++ in Fig. 8D, left). In prosocial subjects, the Main+ and Main++ families were better, whereas in the individualistic subjects, the Main model family was better (Fig. 8D,E). In sum, these results suggest that such feedback processing would likely be involved in social conversion of the feedforward processing identified in this study, and remains a topic for further investigation.

Discussion

This study aimed to characterize a fundamental brain computation in human social cognition: How are social signals activated by rewards to others converted into self-oriented valuation and decision making? Our experimental paradigm enabled a detailed examination of social value conversion during proximal self-oriented decision making, and the isolation of brain signals associated with the computational stages from the offered valuation to the choice-effective valuation and finally to the decision valuation. Further, we indexed differences in value conversion between SVO behavioral phenotypes. Collectively, our results reveal a neural circuit for converting the social value of the others' reward into a self-oriented decision signal and demonstrate that individual variability in the conversion likely underlies discrete sociobehavioral phenotypes.

Value conversion signals for three computational stages were identified: rTPJ and dlPFC activations for the offered value, rAI activation for the signed effective value, and mPFC activation for the integrated DV. These value conversion signals can help in interpreting the neural activation patterns observed in various social contexts in relation to a broader understanding of their functions (Suzuki et al., 2012; Stanley and Adolphs, 2013; Tusche et al., 2016). For example, rTPJ activations are generally consid-

ered critical for inferring others' mental states, as in theory of mind and perspective taking (Saxe, 2010; Lee and Seo, 2016). Our findings show that the rTPJ activation signals the offered value of others even without apparent involvement of theory of mind (Hare et al., 2010). Likewise, dlPFC activations are generally linked to higher-order cognitive operations such as cognitive and self-control (Miller and Cohen, 2001; Rudolf and Hare, 2014; Schmidt et al., 2018), thus mediating the balance between self-impulses and other-regarding considerations in social behavior (Hare et al., 2009; Baumgartner et al., 2011). Our findings on dlPFC activation provide another view, showing the nature of the shared encoding or equivalent treatment of the offered other- and self-bonus values. The rAI region responded to the effective impact of the other-bonus offer on the choice (assessed by the signed effective value) via individual variability, thus mediating the self-relevance of a possible others' outcome in the individual's choices. Although AI activations are typically thought to reflect the involvement of emotion and subjective feelings in decisions (Zaki and Ochsner, 2012; Hein et al., 2016), the observed rAI activation signaled relevance to self-decisions in a setting without strong emotional involvement, which may instead suggest the self-relevance underpinning emotional responses (Rilling and Sanfey, 2011; Gospic et al., 2014). The mPFC responses to the DV (which might well be vmPFC responses) are concordant with the region's role as a general mechanism in signaling decisions (Hare et al., 2008; Ruff and Fehr, 2014).

Our findings on the coupling of these areas' responses from PPI and DCM analyses support the three-stage processing from the offer to the signed effective value, with the dlPFC and rTPJ responses modifying the rAI responses, and from the signed effective value to the decision, with the rAI responses modifying the mPFC responses. These findings are corroborated by previous findings regarding these regions' functional and anatomical connections (Petrides and Pandya, 1999; Ongür and Price, 2000; Mars et al., 2013; Schmidt et al., 2018). Therefore, these findings highlight a feedforward processing of social value conversion through these areas. However, these results do not imply that the conversion is strictly feedforward only, as explored at the end of the Results. Future studies are required to delineate the roles of feedforward and feedback processes in the conversion.

The value conversion process we identified should provide a common denominator for more complex forms of social conversion in decision making. Let us consider social preferences as an example. The conversion treats an option's utility as a simple sum of self-regarding value and other-bonus value, whereas this second term is replaced in many social preference models by terms concerning the balance of reward allocation (Fehr and Schmidt, 1999; Fehr and Camerer, 2007; Gu et al., 2015). If reflecting this balance is correct, then a relevant social factor may modulate the other-bonus value signals found in this study, for instance, shifting each signal relative to the preference or reward allocation balance. Such modulation might occur in parallel across all the stages and would also engage other value-related brain regions such as the ventral striatum, amygdala, and dmPFC, including the anterior cingulate (Nicolle et al., 2012; Báez-Mendoza et al., 2013; Marsh et al., 2014; Bilek et al., 2015; Apps et al., 2016; Wittmann et al., 2016; Hill et al., 2017). In the first stage, such modulation in the rTPJ signals would reflect offered outcomes for others under the allocation balance (Morishima et al., 2012; Kameda et al., 2016). We found that the dlPFC rather than the rdIPFC was involved in offered other-bonus valuations, whereas previous studies have placed greater emphasis on the rdIPFC in social behavior (Nihonsugi et al., 2015; but also see Steinbeis et

←

(Figure legend continued.) based on the model M5 was constructed with connectivity from the dlPFC to the mPFC (M5+). Right, Exceedance probabilities obtained by using random-effect Bayesian model selection to indicate which of the two models better fit data of each SVO group; M5+ was better for the prosocial group, whereas M5 was better for the individualistic group. D, Left, One example (M5++) of the Main++ family. Right, Familywise comparison of fit for each SVO group among the three model families (Main, Main+, and Main++). The Main+ and Main++ model families were each composed of eight models corresponding to those of the main model family, in which feedforward and bidirectional connections were added between the dlPFC and mPFC, respectively. The results indicate that Main+ and Main++ family were better than the Main family for the prosocial group, whereas the Main family was better for individualistic group. E, Exceedance probabilities obtained by using random-effect Bayesian model selection of all the models in the three families in each SVO group.

al., 2012; Crockett et al., 2017). Modulated dlPFC signals might selectively engage the rdlPFC (Buckholtz et al., 2015), thus indexing the allocation balance, for example, with respect to social norms (Baumgartner et al., 2011; Ruff et al., 2013). In the second stage, the modulation of the AI under the allocation balance would likely be related to social dispositions such as fairness, inequity aversion, and guilt (Zaki and Mitchell, 2011; Gluth and Fontanesi, 2016). In future studies, it will be important to investigate integrative modulation processes in social value conversion while taking into account other key considerations. For example, recent studies have suggested that rTPJ signals provided a top-down signal of commitment to altruistic behavior (Park et al., 2017) or resolve moral conflict in response to material rewards (Obeso et al., 2018). How these factors are integrated into and modulate social value conversion, leading to a final decision variable, remains a topic for future study.

As a common denominator process, social value conversion would constitute a primitive process for a range of social behaviors and their conditional variability. For example, our findings suggest that differences in social value conversion may underlie different SVO behavioral phenotypes. For the difference in the rAI activation between prosocial and individualistic subjects, our results suggest that the rAI responses in the prosocial subjects encode individual variability of the self-relevance of others' reward for the decision, whereas in the individualistic subjects, the elevated rAI response might lead to choices that disregard the others' reward (Sanfey et al., 2003; Chang et al., 2011; Dawes et al., 2012). From the findings on the coupling from the dlPFC to the mPFC responses, we consider that the prosocial subjects convert offered other-bonus and self-bonus values into decisions in a similar way, with a relative reliance on dlPFC–mPFC coupling. In contrast, the individualistic subjects might convert other-bonus values via an effective influence on self-oriented decisions, with a relative reliance on the rAI–mPFC coupling (Hare et al., 2010). Indeed, our findings from the DCM indicated that the rAI–mPFC coupling was strong in the individualistic subjects and also considerable in the prosocial subjects. In sum, these findings indicate the importance of value conversion in understanding the variable relationships of brain social functions and motivated behavior.

Because of the simplicity of our task, we could isolate the social value conversion process, but this finding also has limitations. First, in complex social environments involving social inferences and interactions, higher-order processes in the brain would modulate the signals of the conversion that we studied (Coricelli and Nagel, 2009; Yoshida et al., 2010; Crockett et al., 2014; Fareri et al., 2015; Garvert et al., 2015). Second, we did not address relationships involved in how social value is generated, for instance, various motives that are psychologically or economically distinguishable, such as pure altruism, reputation, inequity aversion, direct/indirect reciprocity, and warm-glow giving (Harbaugh et al., 2007; Tricomi et al., 2010; Rilling and Sanfey, 2011). Also, our study considered charity donation as merely an addition to self-gain and thus did not directly address the social conversion in more complex forms of charity donation behavior, which would entail a reduction of self-gain and thus probe the nature of altruism (Moll et al., 2006; Kuss et al., 2013; Kuss et al., 2015). The interaction of generative motives with the primitive conversion, and the conversion process under balancing of other-gain with reduction of self-gain, remain for future studies. In relation to these, a promising research direction is to examine the possibility that the primitive conversion might differ over a

wider spectrum of populations such as antisocial extremes (Marsh, 2019).

In summary, our results help to explain an important but heretofore unaddressed human social computation: how social value migrates into the self-oriented decision-making process via social value conversion. We further demonstrate the behavioral relevance of value conversion in different sociobehavioral phenotypes. These primitive computations could be used as building blocks to enable value conversion for a wide variety of behavioral outputs in human social behaviors, preferences, and interactions based on the social cognition of other-regarding information involving layers of social valuation with increasing complexity. Combining our methodological approach with a systematic exploration of social cognitive behaviors holds great promise for addressing these questions quantitatively.

References

- Adèr HJ, Mellenbergh GJ, Hand DJ (2008) *Advising on research methods: a consultant's companion*: Huizen, The Netherlands: Johannes van Kessel Publishing.
- Apps MA, Rushworth MF, Chang SW (2016) The anterior cingulate gyrus and social cognition: tracking the motivation of others. *Neuron* 90:692–707.
- Báez-Mendoza R, Harris CJ, Schultz W (2013) Activity of striatal neurons reflects social action and own reward. *Proc Natl Acad Sci U S A* 110:16634–16639.
- Baumgartner T, Knoch D, Hotz P, Eisenegger C, Fehr E (2011) Dorsolateral and ventromedial prefrontal cortex orchestrate normative choice. *Nat Neurosci* 14:1468–1474.
- Behrens TE, Hunt LT, Woolrich MW, Rushworth MF (2008) Associative learning of social value. *Nature* 456:245–249.
- Behrens TE, Hunt LT, Rushworth MF (2009) The computation of social behavior. *Science* 324:1160–1164.
- Bilek E, Ruf M, Schäfer A, Akdeniz C, Calhoun VD, Schmahl C, Demanuele C, Tost H, Kirsch P, Meyer-Lindenberg A (2015) Information flow between interacting human brains: identification, validation, and relationship to social expertise. *Proc Natl Acad Sci U S A* 112:5207–5212.
- Bogaert S, Boone C, Declercq C (2008) Social value orientation and cooperation in social dilemmas: a review and conceptual model. *Br J Soc Psychol* 47:453–480.
- Bornstein AM, Daw ND (2012) Dissociating hippocampal and striatal contributions to sequential prediction learning. *Eur J Neurosci* 35:1011–1023.
- Bornstein AM, Daw ND (2013) Cortical and hippocampal correlates of de-liberation during model-based decisions for rewards in humans. *PLoS Comput Biol* 9:e1003387.
- Buckholtz JW, Martin JW, Treadway MT, Jan K, Zald DH, Jones O, Marois R (2015) From blame to punishment: disrupting prefrontal cortex activity reveals norm enforcement mechanisms. *Neuron* 87:1369–1380.
- Chang LJ, Smith A, Dufwenberg M, Sanfey AG (2011) Triangulating the neural, psychological, and economic bases of guilt aversion. *Neuron* 70:560–572.
- Coricelli G, Nagel R (2009) Neural correlates of depth of strategic reasoning in medial prefrontal cortex. *Proc Natl Acad Sci U S A* 106:9163–9168.
- Crockett MJ, Kurth-Nelson Z, Siegel JZ, Dayan P, Dolan RJ (2014) Harm to others outweighs harm to self in moral decision making. *Proc Natl Acad Sci U S A* 111:17320–17325.
- Crockett MJ, Siegel JZ, Kurth-Nelson Z, Dayan P, Dolan RJ (2017) Moral transgressions corrupt neural representations of value. *Nat Neurosci* 20:879–885.
- Daw ND, Doya K (2006) The computational neurobiology of learning and reward. *Curr Opin Neurobiol* 16:199–204.
- Dawes CT, Loewen PJ, Schreiber D, Simmons AN, Flagan T, McElreath R, Bokemper SE, Fowler JH, Paulus MP (2012) Neural basis of egalitarian behavior. *Proc Natl Acad Sci U S A* 109:6479–6483.
- Dayan P, Nakahara H (2018) Models and methods for reinforcement learning. In: *Stevens' handbook of experimental psychology*, Vol V (Wagmakers EJ, Wixted J, eds), pp 507–546. New York: Wiley.
- De Martino B, O'Doherty JP, Ray D, Bossaerts P, Camerer C (2013) In the

- mind of the market: theory of mind biases value computation during financial bubbles. *Neuron* 79:1222–1231.
- DiCiccio TJ, Efron B (1996) Bootstrap confidence intervals. *Statistical Science* 11:189–212.
- Efron B, Tibshirani RJ (1993) *An introduction to the bootstrap*. New York, London: Chapman and Hall.
- Fareri DS, Chang LJ, Delgado MR (2015) Computational substrates of social value in interpersonal collaboration. *J Neurosci* 35:8170–8180.
- Fehr E, Camerer CF (2007) Social neuroeconomics: the neural circuitry of social preferences. *Trends Cogn Sci* 11:419–427.
- Fehr E, Krajbich I (2014) Social preferences and the brain. In: *Neuroeconomics*, Ed 2. (Glimcher P, Fehr E, eds), pp 193–218. Amsterdam: Elsevier.
- Fehr E, Schmidt KM (1999) A theory of fairness, competition, and cooperation. *The Quarterly Journal of Economics* 114:817–868.
- Friston KJ, Buechel C, Fink GR, Morris J, Rolls E, Dolan RJ (1997) Psychophysiological and modulatory interactions in neuroimaging. *Neuroimage* 6:218–229.
- Friston KJ, Harrison L, Penny W (2003) Dynamic causal modelling. *Neuroimage* 19:1273–1302.
- Frith CD, Frith U (2012) Mechanisms of social cognition. *Annu Rev Psychol* 63:287–313.
- Garvert MM, Moutoussis M, Kurth-Nelson Z, Behrens TE, Dolan RJ (2015) Learning-induced plasticity in medial prefrontal cortex predicts preference malleability. *Neuron* 85:418–428.
- Gläscher J (2009) Visualization of group inference data in functional neuroimaging. *Neuroinformatics* 7:73–82.
- Glimcher P, Fehr E (2014) *Neuroeconomics: decision making and the brain*. In: *Neuroeconomics*, Ed 2 (Fehr PWG, ed). San Diego: Academic.
- Gluth S, Fontanesi L (2016) Wiring the altruistic brain. *Science* 351:1028–1029.
- Gospic K, Sundberg M, Maeder J, Fransson P, Petrovic P, Isacsson G, Karlström A, Ingvar M (2014) Altruism costs—the cheap signal from amygdala. *Soc Cogn Affect Neurosci* 9:1325–1332.
- Gu X, Wang X, Hula A, Wang S, Xu S, Lohrenz TM, Knight RT, Gao Z, Dayan P, Montague PR (2015) Necessary, yet dissociable contributions of the insular and ventromedial prefrontal cortices to norm adaptation: computational and lesion evidence in humans. *J Neurosci* 35:467–473.
- Harbaugh WT, Mayr U, Burghart DR (2007) Neural responses to taxation and voluntary giving reveal motives for charitable donations. *Science* 316:1622–1625.
- Hare TA, O’Doherty J, Camerer CF, Schultz W, Rangel A (2008) Dissociating the role of the orbitofrontal cortex and the striatum in the computation of goal values and prediction errors. *J Neurosci* 28:5623–5630.
- Hare TA, Camerer CF, Rangel A (2009) Self-control in decision making involves modulation of the vmPFC valuation system. *Science* 324:646–648.
- 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.
- Haruno M, Frith CD (2010) Activity in the amygdala elicited by unfair divisions predicts social value orientation. *Nat Neurosci* 13:160–161.
- Hein G, Morishima Y, Leiberg S, Sul S, Fehr E (2016) The brain’s functional network architecture reveals human motives. *Science* 351:1074–1078.
- Hill CA, Suzuki S, Polania R, Moisa M, O’Doherty JP, Ruff CC (2017) A causal account of the brain network computations underlying strategic social behavior. *Nat Neurosci* 20:1142–1149.
- Hsu M, Anen C, Quartz SR (2008) The right and the good: distributive justice and neural encoding of equity and efficiency. *Science* 320:1092–1095.
- Hu X, Le TH, Parrish T, Erhard P (1995) Retrospective estimation and correction of physiological fluctuation in functional MRI. *Magnetic resonance in medicine* 34:201–212.
- Hutcherson CA, Bushong B, Rangel A (2015) A neurocomputational model of altruistic choice and its implications. *Neuron* 87:451–462.
- Kameda T, Inukai K, Higuchi S, Ogawa A, Kim H, Matsuda T, Sakagami M (2016) Rawlsian maximin rule operates as a common cognitive anchor in distributive justice and risky decisions. *Proc Natl Acad Sci U S A* 113:11817–11822.
- Kriegeskorte N, Simmons WK, Bellgowan PS, Baker CI (2009) Circular analysis in systems neuroscience: the dangers of double dipping. *Nat Neurosci* 12:535–540.
- Kuss K, Falk A, Trautner P, Elger CE, Weber B, Fließbach K (2013) A reward prediction error for charitable donations reveals outcome orientation of donors. *Soc Cogn Affect Neurosci* 8:216–223.
- Kuss K, Falk A, Trautner P, Montag C, Weber B, Fließbach K (2015) Neuronal correlates of social decision making are influenced by social value orientation—an fMRI study. *Front Behav Neurosci* 9:40.
- Lee D, Seo H (2016) Neural basis of strategic decision making. *Trends Neurosci* 39:40–48.
- Liljeholm M, Dunne S, O’Doherty JP (2015) Differentiating neural systems mediating the acquisition vs expression of goal-directed and habitual behavioral control. *Eur J Neurosci* 41:1358–1371.
- Mars RB, Sallet J, Neuberta XF, Rushworth MF (2013) Connectivity profiles reveal the relationship between brain areas for social cognition in human and monkey temporoparietal cortex. *Proc Natl Acad Sci U S A* 110:10806–10811.
- Marsh AA (2019) The caring continuum: evolved hormonal and proximal mechanisms explain prosocial and antisocial extremes. *Annu Rev Psychol* 70:347–371.
- Marsh AA, Stoycos SA, Brethel-Haurwitz KM, Robinson P, VanMeter JW, Cardinale EM (2014) Neural and cognitive characteristics of extraordinary altruists. *Proc Natl Acad Sci U S A* 111:15036–15041.
- Miller EK, Cohen JD (2001) An integrative theory of prefrontal cortex function. *Annu Rev Neurosci* 24:167–202.
- Moll J, Krueger F, Zahn R, Pardini M, de Oliveira-Souza R, Grafman J (2006) Human fronto-mesolimbic networks guide decisions about charitable donation. *Proc Natl Acad Sci U S A* 103:15623–15628.
- Montague PR, King-Casas B, Cohen JD (2006) Imaging valuation models in human choice. *Annu Rev Neurosci* 29:417–448.
- Morishima Y, Schunk D, Bruhin A, Ruff CC, Fehr E (2012) Linking brain structure and activation in temporoparietal junction to explain the neurobiology of human altruism. *Neuron* 75:73–79.
- Nichols T, Brett M, Andersson J, Wager T, Poline JB (2005) Valid conjunction inference with the minimum statistic. *Neuroimage* 25:653–660.
- Nicolle A, Klein-Flügge MC, Hunt LT, Vlaev I, Dolan RJ, Behrens TE (2012) An agent independent axis for executed and modeled choice in medial prefrontal cortex. *Neuron* 75:1114–1121.
- Nihonsugi T, Ihara A, Haruno M (2015) Selective increase of intention-based economic decisions by noninvasive brain stimulation to the dorsolateral prefrontal cortex. *J Neurosci* 35:3412–3419.
- Obeso I, Moisa M, Ruff CC, Dreher JC (2018) A causal role for right temporo-parietal junction in signaling moral conflict. *Elife* 7:e40671.
- O’Doherty JP, Hampton A, Kim H (2007) Model-based fMRI and its application to reward learning and decision making. *Ann N Y Acad Sci* 1104:35–53.
- Ongür D, Price JL (2000) The organization of networks within the orbital and medial prefrontal cortex of rats, monkeys and humans. *Cereb Cortex* 10:206–219.
- O’Reilly JX, Woolrich MW, Behrens TE, Smith SM, Johansen-Berg H (2012) Tools of the trade: psychophysiological interactions and functional connectivity. *Social cognitive and affective neuroscience* 7:604–609.
- Page EB (1963) Page ordered hypotheses for multiple treatments: a significance test for linear ranks. *J Am Stat Assoc* 58:216–230.
- Park SQ, Kahnt T, Dogan A, Strang S, Fehr E, Tobler PN (2017) A neural link between generosity and happiness. *Nat Commun* 8:15964.
- Penny WD, Stephan KE, Daunizeau J, Rosa MJ, Friston KJ, Schofield TM, Leff AP (2010) Comparing families of dynamic causal models. *PLoS Comput Biol* 6:e1000709.
- Petrides M, Pandya DN (1999) Dorsolateral prefrontal cortex: comparative cytoarchitectonic analysis in the human and the macaque brain and corticocortical connection patterns. *Eur J Neurosci* 11:1011–1036.
- Phelps EA, Delgado MR, Nearing KI, LeDoux JE (2004) Extinction learning in humans: role of the amygdala and vmPFC. *Neuron* 43:897–905.
- Platt ML, Huettel SA (2008) Risky business: the neuroeconomics of decision making under uncertainty. *Nat Neurosci* 11:398–403.
- Prelec D (1998) The probability weighting function. *Econometrica* 66:497–527.
- Rangel A, Camerer C, Montague PR (2008) A framework for studying the neurobiology of value-based decision making. *Nat Rev Neurosci* 9:545–556.
- Rilling JK, Sanfey AG (2011) The neuroscience of social decision making. *Annu Rev Psychol* 62:23–48.
- Rudolf S, Hare TA (2014) Interactions between dorsolateral and ventrome-

- dial prefrontal cortex underlie context-dependent stimulus valuation in goal-directed choice. *J Neurosci* 34:15988–15996.
- Ruff CC, Fehr E (2014) The neurobiology of rewards and values in social decision making. *Nat Rev Neurosci* 15:549–562.
- Ruff CC, Ugazio G, Fehr E (2013) Changing social norm compliance with noninvasive brain stimulation. *Science* 342:482–484.
- Rushworth MF, Behrens TE (2008) Choice, uncertainty and value in prefrontal and cingulate cortex. *Nat Neurosci* 11:389–397.
- Rushworth MF, Noonan MP, Boorman ED, Walton ME, Behrens TE (2011) Frontal cortex and reward-guided learning and decision making. *Neuron* 70:1054–1069.
- Sanfey AG, Rilling JK, Aronson JA, Nystrom LE, Cohen JD (2003) The neural basis of economic decision making in the ultimatum game. *Science* 300:1755–1758.
- Saxe R (2010) The right temporo-parietal junction: a specific brain region for thinking about thoughts. In: *Handbook of theory of mind* (Leslie A, German T, eds), pp 1–35. Philadelphia, PA: Psychology Press, Taylor & Francis Group.
- Schmidt L, Tusche A, Manoharan N, Hutcherson C, Hare T, Plassmann H (2018) Neuroanatomy of the vmPFC and dlPFC predicts individual differences in cognitive regulation during dietary self-control across regulation strategies. *J Neurosci* 38:5799–5806.
- Schultz W, Dayan P, Montague PR (1997) A neural substrate of prediction and reward. *Science* 275:1593–1599.
- Stanley DA, Adolphs R (2013) Toward a neural basis for social behavior. *Neuron* 80:816–826.
- Steinbeis N, Bernhardt BC, Singer T (2012) Impulse control and underlying functions of the left DLPFC mediate age-related and age-independent individual differences in strategic social behavior. *Neuron* 73:1040–1051.
- Stephan KE, Penny WD, Daunizeau J, Moran RJ, Friston KJ (2009) Bayesian model selection for group studies. *Neuroimage* 46:1004–1017.
- Strombach T, Weber B, Hangebrauk Z, Kenning P, Karipidis II, Tobler PN, Kalenscher T (2015) Social discounting involves modulation of neural value signals by temporoparietal junction. *Proc Natl Acad Sci U S A* 112:1619–1624.
- Suzuki S, Harasawa N, Ueno K, Gardner JL, Ichinohe N, Haruno M, Cheng K, Nakahara H (2012) Learning to simulate others' decisions. *Neuron* 74:1125–1137.
- Talairach J, Tournoux P (1988) *Co-planar stereotaxic atlas of the human brain*. Stuttgart, New York: Georg Thieme Verlag.
- Tobler PN, O'Doherty JP, Dolan RJ, Schultz W (2007) Reward value coding distinct from risk attitude-related uncertainty coding in human reward systems. *J Neurophysiol* 97:1621–1632.
- Tricomi E, Rangel A, Camerer CF, O'Doherty JP (2010) Neural evidence for inequality-averse social preferences. *Nature* 463:1089–1091.
- Tusche A, Böckler A, Kanske P, Trautwein FM, Singer T (2016) Decoding the charitable brain: empathy, perspective taking, and attention shifts differentially predict altruistic giving. *J Neurosci* 36:4719–4732.
- Van Lange PAM (1999) The pursuit of joint outcomes and equality in outcomes: an integrative model of social value orientation. *J Pers Soc Psychol* 77:337–349.
- Van Lange PA, Otten W, De Bruin EM, Joireman JA (1997) Development of prosocial, individualistic, and competitive orientations: theory and preliminary evidence. *J Pers Soc Psychol* 73:733–746.
- Winkler AM, Ridgway GR, Webster MA, Smith SM, Nichols TE (2014) Permutation inference for the general linear model. *Neuroimage* 92:381–397.
- Wittmann MK, Kolling N, Faber NS, Scholl J, Nelissen N, Rushworth MF (2016) Self-other merge in the frontal cortex during cooperation and competition. *Neuron* 91:482–493.
- Yoshida W, Seymour B, Friston KJ, Dolan RJ (2010) Neural mechanisms of belief inference during cooperative games. *J Neurosci* 30:10744–10751.
- Zaki J, Mitchell JP (2011) Equitable decision making is associated with neural markers of intrinsic value. *Proc Natl Acad Sci U S A* 108:19761–19766.
- Zaki J, Ochsner KN (2012) The neuroscience of empathy: progress, pitfalls and promise. *Nat Neurosci* 15:675–680.