

Spatial gradient in value representation along the medial prefrontal cortex reflects individual differences in prosociality

Sunhae Sul^{a,1}, Philippe N. Tobler^b, Grit Hein^b, Susanne Leiberg^b, Daehyun Jung^c, Ernst Fehr^b, and Hackjin Kim^{a,2}

^aDepartment of Psychology, Korea University, Seoul 136-701, Republic of Korea; ^bDepartment of Economics, University of Zurich, CH-8006 Zurich, Switzerland; and ^cDepartment of Brain and Cognitive Engineering, Korea University, Seoul 136-701, Republic of Korea

Edited by Susan T. Fiske, Princeton University, Princeton, NJ, and approved May 5, 2015 (received for review December 13, 2014)

Despite the importance of valuing another person's welfare for prosocial behavior, currently we have only a limited understanding of how these values are represented in the brain and, more importantly, how they give rise to individual variability in prosociality. In the present study, participants underwent functional magnetic resonance imaging while performing a prosocial learning task in which they could choose to benefit themselves and/or another person. Choice behavior indicated that participants valued the welfare of another person, although less so than they valued their own welfare. Neural data revealed a spatial gradient in activity within the medial prefrontal cortex (MPFC), such that ventral parts predominantly represented self-regarding values and dorsal parts predominantly represented other-regarding values. Importantly, compared with selfish individuals, prosocial individuals showed a more gradual transition from self-regarding to other-regarding value signals in the MPFC and stronger MPFC–striatum coupling when they made choices for another person rather than for themselves. The present study provides evidence of neural markers reflecting individual differences in human prosociality.

medial prefrontal cortex | striatum | anterior insula | reinforcement learning | computational model

Ranging from small acts of kindness in daily life to self-sacrificing altruism under life-threatening situations, we often observe large individual differences in how humans value another person's welfare. This differential valuation process seems to be the key to understanding various human prosocial behaviors, which are fundamental to the sustainability of human society (1). The underlying neural mechanisms and their relationship to individual differences in prosociality remain unclear, however.

Perhaps the most powerful way of assessing how an outcome is valued is to use an instrumental learning paradigm that examines whether the occurrence of a response increases when it is followed by that outcome (2). The mechanisms underlying this type of learning have been described more formally with a computational model, known as the advantage learning model (3–5), which has been used successfully to reveal the neuroanatomical substrates of subjective valuation (3, 4, 6). Previous research has further refined the neurobiological model of reinforcement learning by emphasizing the specific roles played by the medial frontal cortex and the striatum; the medial frontal cortex computes the value of the chosen action, whereas the striatum processes reward prediction errors during reinforcement learning (4, 6–10).

Unlike our current understanding of the valuation process for self-regarding choices (3, 6–12), it is much less clear whether learning also can be driven by other-regarding values, and whether this other-regarding valuation relies on the same mechanisms of reinforcement learning as those used for self. Moreover, despite the rapidly accumulating research on reward processing in social domains (13–19), the question remains of how neural representation of self-regarding vs. other-regarding values is related to individual differences in altruistic behavior.

In this work, we designed a novel version of an instrumental learning task (i.e., a prosocial learning task) to assess behavioral and

neural processes associated with self- and other-regarding valuation in a comparable, principled way. In the prosocial learning task, participants chose between two alternatives to achieve a higher probability of benefiting either themselves and/or another person by reducing the duration of exposure to unpleasantly loud noise. Thirty pairs of healthy right-handed female college students participated in the study. The scanned participant of each pair performed the prosocial learning task (Fig. 1). In each trial of the task, participants were presented with two options and had to choose one of them. In different conditions, the two options were represented by specific fractal images and associated with points only for the participant in the scanner (SELF condition), for both participants (BOTH condition), or only for the participant outside the scanner (OTHER condition). One of the two options always had a higher probability of yielding points than the other (70% vs. 30%). By trial and error, the participants in the scanner would learn about these probabilities and subsequently choose the option that they preferred. Participants were told that they would be exposed to unpleasant noise for 5 min after the task, and that the points earned in the task would be used to reduce the duration of the noise for themselves and/or the paired participant outside the scanner (see *SI Appendix* for details).

We predicted that if participants valued others' welfare, then the other-regarding outcome (i.e., points earned to reduce the

Significance

How do selfish and prosocial brains function differently with regard to valuing the welfare of others? The present study addresses this question by combining neuroimaging, computational modeling, and an instrumental conditioning paradigm. Contrary to the conventional notion of the dorsal medial prefrontal cortex (MPFC) implicated in mentalization, we found that it was selfish individuals who showed greater spatial segregation between ventral and dorsal MPFC, which encoded self- and other-regarding values, respectively. Prosocial individuals, on the other hand, were characterized by overlapping self-other representation in the ventral MPFC and by stronger functional coupling between MPFC and striatum while representing and updating the value of other-regarding choices. These findings provide rigorous scientific evidence of neural markers reflecting individual differences in human prosociality.

Author contributions: S.S., P.N.T., G.H., S.L., E.F., and H.K. designed research; S.S. and D.J. performed research; S.S. and H.K. analyzed data; and S.S., P.N.T., and H.K. wrote the paper.

The authors declare no conflict of interest.

This article is a PNAS Direct Submission.

Freely available through the PNAS open access option.

¹Present address: Department of Psychological and Brain Sciences, Dartmouth College, Hanover, NH 03755.

²To whom correspondence should be addressed. Email: hackjinkim@korea.ac.kr.

This article contains supporting information online at www.pnas.org/lookup/suppl/doi:10.1073/pnas.1423895112/-DCSupplemental.

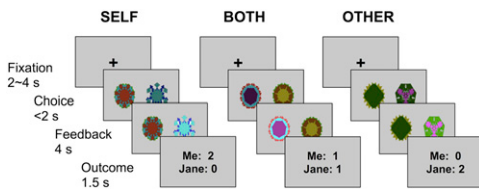


Fig. 1. Prosocial learning task and within-subject experimental conditions. Only rewarded trials are shown.

duration of aversive noise for the other) would increase their performance above chance level, such that they would earn more points for the other participant than if they chose randomly. In line with the idea that avoidance of punishment is reinforcing and has been shown to activate brain regions similar to those involved in reward learning (3), the points earned in the task, which could be used later, just like money, were presumed to have appetitive motivational value. Therefore, regarding neural representation of the chosen value, we expected a spatial segregation within the medial prefrontal cortex (MPFC) in computing self- and other-regarding values, consistent with previous studies showing that the ventral and dorsal parts of the MPFC are involved in self- and other-regarding processes, respectively (19–25). More importantly, we hypothesized that the degree of spatial segregation would provide a neural index of the individual propensity to help others. Given that positive subjective valuation of others' welfare can lead to prosocial decisions (17, 26–29), we expected to find that decreased segregation would be associated with greater prosociality. In addition, we examined whether and how corticostriatal communications contribute to individual differences in representing and updating self- and other-regarding values.

Behavioral Results

We tested whether participants valued another person's welfare at all, and found that they did. In particular, the proportions of choosing the high reward probability (HRP) option in the OTHER condition were significantly higher than chance level [0.5; $t(25) = 2.68, P < 0.05$] (Fig. 2A; trial-by-trial learning curves in *SI Appendix, Fig. S1*). Although this finding clearly indicates that participants did learn to help others even when they had nothing to gain, there were considerable individual differences in their propensity to help. Some individuals showed equal preference for the HRP option in the SELF and the OTHER conditions, whereas other individuals showed such a preference only in the SELF condition. Owing to this individual variability, preference for the HRP option was weaker on average in the OTHER condition compared with the SELF or BOTH condition ($P < 0.05$ for both); main effect of condition: $F(2, 50) = 5.97, P < 0.01$; pairwise comparisons for SELF vs. BOTH: not significant (Fig. 2A).

Functional Magnetic Resonance Imaging Results

Spatial Gradient Within MPFC for Self- vs. Other-Regarding Value Computation. We hypothesized that the ventral and dorsal subregions of the MPFC would be involved in computing self- and other-regarding values, respectively. A minimal requirement to support this hypothesis is that the MPFC as a whole would be associated with choice values in all conditions. Thus, we conducted a parametric modulation analysis using subject-specific value parameters estimated by the advantage learning model (*Materials and Methods*) and found that the MPFC ($x = 4, y = 52, z = 8$ mm, $Z = 3.73$) was engaged in computing the value of the chosen option at the time of stimulus presentation across all three conditions (*SI Appendix, Table S1 and Fig. S3A*). We next ran the same parametric modulation analysis separately for each

condition. These analyses revealed that the chosen value-related MPFC activation clusters in the SELF and OTHER conditions were located somewhat more ventrally and more dorsally, respectively, than the cluster in the BOTH condition (*SI Appendix, Table S1 and Fig. S3B*), consistent with our prediction of spatial specificity.

For a more quantitative examination of this spatial segregation within the MPFC, we defined five regions of interest (ROIs) along the ventral–dorsal midline axis. Specifically, we obtained a sagittal view of the statistical parametric map from the aforementioned parametric modulation analysis with a lenient threshold of $P < 0.05$ uncorrected and then selected five equally spaced coordinates spanning the ventral-to-dorsal extent of the MPFC (Fig. 3A), similar to a previous approach (30). The parameter estimates of the neural activation associated with chosen value at the time of stimulus presentation were extracted from the ROIs for each individual. A 2 (condition: SELF, OTHER) \times 5 (ROI locations) repeated-measures ANOVA showed a clear spatial distinction between the ventromedial prefrontal cortex (VMPFC) and the dorsomedial prefrontal cortex (DMPFC) in computing values for self-regarding vs. other-regarding choices, respectively [interaction of condition with ROI locations: $F(4, 100) = 4.49, P < 0.005$] (Fig. 3B). The value signal for SELF was stronger in more ventral ROIs, whereas the value signal for OTHER was stronger in more dorsal ROIs within the MPFC. Value signals for the BOTH condition were dominant in intermediate ROIs, and the interaction effect remained significant when we included the BOTH condition in the analysis [$F(8, 200) = 2.28, P < 0.05$]. In line with the gradient hypothesis, the difference in the linear effect of ROI location between conditions was also significant [$F(1, 25) = 10.13, P < 0.005$]. More specifically, the strength of the value signal for self-regarding choices decreased linearly from the VMPFC to the DMPFC [$F(1, 25) = 4.81, P < 0.05$], whereas the value signal for other-regarding choices showed the opposite trend [$F(1, 25) = 3.10, P = 0.09$].

Self–Other Distinction Within the MPFC Reflects a Propensity to Help Others in the Prosocial Learning Task.

We expected that the spatial pattern of value representations within the MPFC would reliably track individual variability in prosocial propensity, as measured by choice behavior in the prosocial learning task. To capture individual differences in a categorical manner, we formed two groups, namely, prosocial and selfish groups, based on the parameters estimated by the advantage learning model (*Materials and Methods*). Fig. 2B illustrates the behavioral characteristics of the prosocial and selfish individuals in terms of their propensity to choose HRP options across conditions (learning curves in *SI Appendix, Fig. S1*). Group membership (prosocial and selfish) interacted with condition (SELF, BOTH, and OTHER) with respect to the proportions of high reward probability choice [$F(2, 46) = 11.78, P < 0.001$]. Compared with the selfish group, the prosocial group had a smaller difference between the SELF and OTHER conditions. (*SI Appendix, Figs. S2 and S4* presents validation of the categorization.)

To better characterize the effects of prosociality on value representation, we conducted a mixed ANOVA with reward type and ROI location as within-subject factors and group

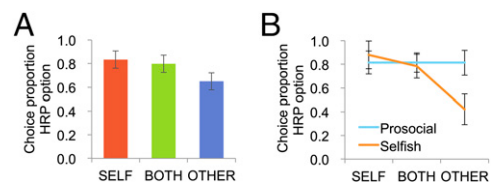


Fig. 2. Behavioral results. (A) Proportions of choosing HRP options in SELF, BOTH, and OTHER conditions. (B) Comparison of the prosocial and selfish groups. Error bars indicate 95% CI.

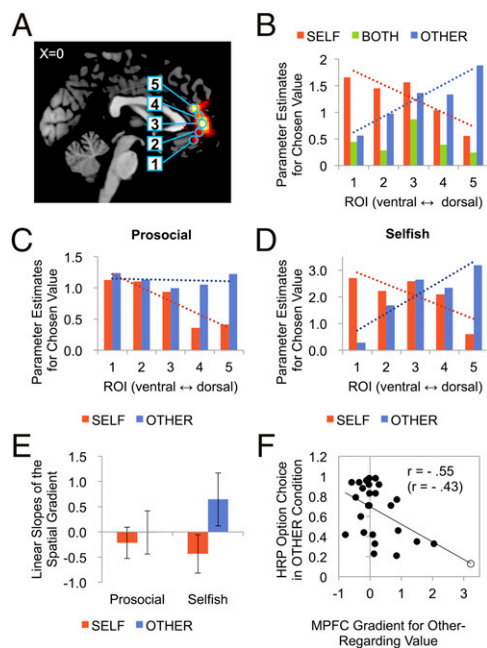


Fig. 3. Regions representing the value of the chosen option. (A) Definition of ROIs. (B) Spatial gradient for self- and other-regarding value computation within the MPFC. Dotted lines indicate linear fits of the spatial gradient for SELF (red) and OTHER (blue) conditions. (C and D) Spatial gradient for self- and other-regarding value computation within the MPFC, depicted separately for prosocial ($n = 15$) and selfish groups ($n = 10$) as defined by the advantage learning model (*SI Appendix*). Dotted lines indicate linear fits of the spatial gradient for SELF (red) and OTHER (blue) conditions. (E) Linear slopes of the spatial gradient within the MPFC in SELF and OTHER conditions among prosocial and selfish participants. Negative slope indicates greater value signal in VMPFC than DMPFC, and vice versa. (F) Participants with greater gradient showed a less clear preference for the HRP option in the OTHER condition. The result remained significant after excluding the subject with the strongest gradient (open circle; the correlation coefficient without this data point is reported in brackets).

membership as a between-subjects factor. We found a significant three-way interaction [$F(4, 92) = 3.10, P < 0.05$], such that the spatial distinction of self-regarding and other-regarding value representations was stronger in the selfish group than in the prosocial group (Fig. 3 C and D). Analyses performed separately for the prosocial and selfish groups further supported this finding: the spatial separation of self- and other-regarding value signals was prominent only among selfish individuals, [$F(4, 36) = 5.37, P < 0.005$], and not among prosocial individuals [$F(4, 56) < 1, P$ not significant]. We quantified the degree of spatial separation within the MPFC by fitting linear functions to the self- and other-regarding value signals along the ventral–dorsal axis for each individual. A between-groups comparison of the linear slopes fitted to the spatial gradient revealed that slopes were steeper in selfish individuals compared with prosocial individuals [$F(1, 23) = 6.72, P < 0.05$] (Fig. 3E), consistent with a greater separation between self- and other-regarding values in selfish individuals compared with prosocial individuals. The spatial gradient revealed that the difference between selfish and prosocial groups arose mainly in the OTHER condition [$F(1, 23) = 4.015, P = 0.057$]; that is, the other-regarding value signal was stronger in the DMPFC than in the VMPFC only in selfish individuals. Such a group difference was not observed in the SELF condition [$F(1, 23) < 1, P$ not significant], where the self-regarding value signal was stronger in the VMPFC than in the DMPFC for both groups. Control analyses showed that the extent of MPFC activation was not merely associated with performance

level (*SI Appendix* provides additional analyses addressing alternative explanations).

Furthermore, spatial gradient tracked the choices in the prosocial learning task. The slopes of the spatial gradient for self- and other-regarding values correlated negatively with the average proportion of choosing HRP options in the OTHER condition ($r = -0.55, P < 0.01$), and the correlation remained significant after excluding the most extreme value, which could be considered an outlier ($r = -0.43, P < 0.05$) (Fig. 3F). That is, participants with greater other-regarding value signals in the VMPFC than in the DMPFC were more likely to choose to help their partners in the OTHER condition, whereas those with the opposite gradient were more likely to behave selfishly. The spatial gradient for self-regarding values was not associated with the choices in the SELF condition ($r = -0.075, ns$).

Functional Connectivity of the MPFC During Other-Regarding vs. Self-Regarding Choices.

It has been well established by previous studies that communication between the medial frontal regions computing the chosen values and the striatum processing the reward prediction error (RPE) plays an essential role in updating and maintaining value-related information during reinforcement learning (7, 8). We confirmed that activity in the striatum, including the nucleus accumbens and parts of the caudate and putamen, was correlated with the RPE at the time of the outcome presentation phase, irrespective of an individual's propensity to behave prosocially (*SI Appendix*). We then performed psychophysiological interaction (PPI) analyses to test whether and how the individual differences observed in the present study are reflected in the pattern of functional connectivity between the MPFC subregions and the striatum during other-regarding vs. self-regarding choices. We selected the ventral (VMPFC; $x = 0, y = 56, z = 2$; peak voxel computing chosen value for SELF condition), middle (MMPFC; $x = 4, y = 52, z = 8$; peak voxel computing chosen value for both SELF and OTHER conditions), and dorsal (DMPFC; $x = 2, y = 44, z = 12$; peak voxel computing chosen value for OTHER condition) parts of the MPFC as seed regions. Then we performed three separate PPI analyses to identify the regions showing differential functional coupling with the three seed regions in the OTHER vs. SELF condition at the option presentation phase. Finally, we performed two-sample t tests of the difference between selfish and prosocial groups. As expected, we found a significant group difference in functional connectivity between the striatum and the VMPFC ($x = 16, y = 20, z = 0, Z = 3.54$), MMPFC ($x = 6, y = 14, z = 2, Z = 3.36$), and DMPFC ($x = 14, y = 18, z = 4, Z = 4.14$) (Fig. 4 A and B). To better understand these group differences, we performed post hoc ROI analyses and found that the difference between the OTHER and SELF condition was robust among prosocial individuals, such that the connectivity was stronger in the OTHER condition than in the SELF condition for all three MPFC seeds [all $F(1, 23) > 9.98$, all $P < 0.01$] (Fig. 4 C–E). In contrast, selfish individuals tended to show the opposite patterns, with stronger connectivity in the SELF condition than in the OTHER condition, although the differences between the two conditions were not statistically significant. It is also worth noting that the part of the striatum communicating with the DMPFC covered a large area extending from the nucleus accumbens to the dorsal caudate and putamen. In contrast, the regions communicating with the MMPFC and VMPFC were more restricted, and the peak voxels were located more ventrally, than those connected with the DMPFC (Fig. 4 A and B and *SI Appendix*, Fig. S5 and Table S3 for regions other than the striatum). *SI Appendix* provides for additional analyses addressing alternative explanations.

Mode of Decisions for Self and Other. The differences between prosocial and selfish individuals in representing and updating

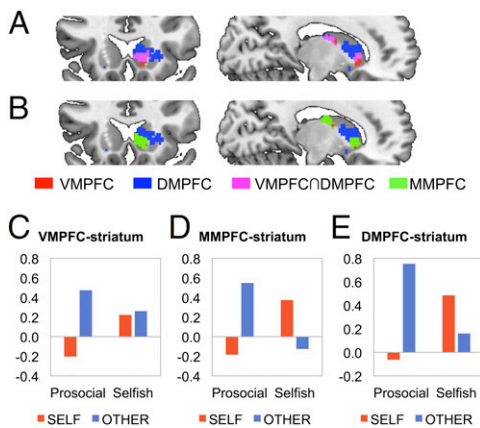


Fig. 4. Comparison of prosocial and selfish individuals for functional connectivity (PPI) with the VMPFC, MMPFC, and DMPFC as seed regions during OTHER vs. SELF conditions. (A) Result of two-sample *t* tests for the PPIs with the VMPFC and DMPFC as seed regions masked with the striatum ROI. The different seed regions are color-coded (red, VMPFC; blue, DMPFC; magenta, VMPFC ∩ DMPFC; for illustration purposes; $P < 0.005$, uncorrected). (B) Result of two-sample *t* tests for all three PPIs masked with the striatum ROI. Voxels connected with the VMPFC, MMPFC, and DMPFC as seed regions are color-coded with red, green, and blue, respectively. Note that the regions connected with VMPFC and MPFC largely overlap; for illustration purposes, $P < 0.005$, uncorrected. (C–E) Average connectivity strength between the MPFC subregions and the striatum in SELF and OTHER conditions, showing a close match with the differential prosociality of the prosocial and selfish groups.

self-regarding and other-regarding values lead us to the question of whether the two groups engage in different modes of decision in SELF and OTHER conditions. Our response time (RT) data suggested the possibility that additional cognitive processes might be required for selfish individuals to make other-regarding decisions, because selfish individuals were significantly slower in the OTHER than the SELF condition [$F(2, 18) = 3.84, P < 0.05$] (*SI Appendix, Fig. S8*). Prosocial individuals did not show such a difference [$F(2, 28) = 0.49, P$ not significant]. In line with this behavioral finding, regions known to be involved in cognitive control, such as the right anterior insula (AI) extending to the inferior frontal gyrus (IFG), showed greater activation when selfish participants made choices for the other participant rather than for themselves, whereas prosocial individuals showed no significant difference across conditions (AI/IFG; $x = 36, y = 26, z = -4, Z = 3.65$; two-sample *t* test for the group difference in the contrast maps of OTHER vs. SELF condition at the time of option presentation; *SI Appendix, Fig. S9A*; see Table S4 for whole-brain result). This result remained the same even when trial-to-trial RTs were included in the analysis (*SI Appendix, Fig. S10*), ruling out the possibility that the increased AI/IFG activation during the OTHER condition among selfish participants may merely reflect differences in RT. Interestingly, AI/IFG activation was correlated with the average RTs, such that participants with slower RTs in the OTHER condition compared with the SELF condition showed greater AI/IFG activation in the OTHER condition ($r = 0.50, P < 0.05$), a difference that remained significant even after controlling for the effect of the need for cognition (31). In sum, these findings suggest that the AI/IFG may be involved in decision mode switching, which is then indirectly related to additional information processing. The RTs and neural responses in the SELF condition were not related to performance.

To further examine how AI/IFG activation influences the corticostriatal communication underlying the process of updating and representing other-regarding vs. self-regarding values, we tested the correlation between the AI/IFG activation (i.e., the beta estimates of the OTHER vs. SELF contrasts at the option

presentation phase) and the MPFC–striatum connectivity (i.e., the beta estimates of the PPIs with the MPFC subregions as seeds during OTHER vs. SELF trials) across participants. Interestingly, the greater AI/IFG activation in the OTHER condition than the SELF condition, the weaker DMPFC–striatum ($r = -0.48, P = 0.01$), MMPFC–striatum ($r = -0.47, P = 0.01$), and VMPFC–striatum ($r = -0.37, P = 0.11$) coupling in the OTHER condition vs. the SELF condition (*SI Appendix, Fig. S9B*).

Discussion

The present study investigated the neural mechanisms of valuing and representing another’s welfare and their relation to an individual’s propensity for prosocial behavior. Combined with a computational approach, our prosocial learning task provides a novel behavioral measure to quantify individual differences in prosociality and allowed us to explore the question that we raised in the beginning: What makes some people more prosocial than others, and how does our brain enable us to value the welfare of others?

Our finding that the spatial specificity for self-regarding vs. other-regarding value representation within the MPFC was robust only among selfish individuals and was attenuated among prosocial individuals supports the idea that prosociality requires a shared value representation for self and other (1, 17, 26). Interestingly, a closer examination of the spatial gradient revealed that the difference between prosocial and selfish groups was especially prominent in the VMPFC. This might seem puzzling, considering that the VMPFC has been strongly implicated in the processes of self-relevant information (32–34), and the DMPFC has been implicated in theory of the mind and mentalizing (35–37). However, there is converging evidence for functional specialization within the MPFC. For example, the VMPFC has been suggested as a domain-general valuation system that processes significant and motivating information, such as reward (38, 39), and the DMPFC has been suggested to be part of the attentional system that is predominantly involved in cognitively demanding tasks, such as strategic social inference (38–40). This idea does not necessarily contradict the idea that self-relevance is a major factor distinguishing VMPFC and DMPFC (41) because self-relevant information is often most significant and motivating (19, 42). Our results suggest that the VMPFC is tightly associated with subjective valuation regardless of the choice’s beneficiary (30), whereas the DMPFC is involved in more general other-specific processing commonly required for other-regarding choices invariant across individuals (22–24, 40, 43, 44).

Another interesting finding is that the pattern of functional coupling between the MPFC subregions computing the values of choices and the striatum encoding RPEs was significantly correlated with an individual’s propensity to help others. In support of this finding, many previous studies have reported strong anatomical and functional links between MPFC subregions and the striatum, which play a key role in reinforcement learning (7, 8, 45). More specifically, recent theoretical work proposed a hierarchical model in which reinforcement learning in vertebrates occurs through multiple independent corticostriatal loops that interact with one another, allowing information to propagate mostly from ventral to dorsal corticostriatal loops (45). Although the spatial resolution of functional magnetic resonance imaging (fMRI) is far lower than that of animal neurophysiological studies (45), we found a similar spatial segregation between ventral and dorsal corticostriatal networks. Our PPI data suggest that prosocial individuals may be characterized by active propagation of subjective value signals between the ventral and dorsal loops, which could be crucial for maximizing their capacity to represent, update, and maintain the value of other-regarding choices. This in turn may enable the shared value representation for self and other within the MPFC.

Despite the stronger functional coupling between the striatum and the MPFC among prosocial individuals, it appears that selfish individuals required greater cognitive effort and control during the OTHER vs. SELF condition, where they had slower RTs and showed increased activity in the AI/IFG. Although we cannot completely rule out the possibility that the AI/IFG activation may reflect an aversive response to other-regarding choices among selfish individuals, the correlation between AI/IFG activation and the average RTs across individuals suggests that additional cognitive processes may be engaged during the OTHER condition. It is also noteworthy that AI/IFG has been strongly implicated in cognitive control and self-regulation (46–49). Given that an increase in AI/IFG activity weakened the MPFC–striatum coupling during choices for others, selfish individuals seem to use additional cognitively demanding processes that interrupt the process of prosocial valuation, which may involve signal propagation through the MPFC–striatum loops. The exact nature of these additional cognitive processes used by selfish individuals merits further investigation.

In summary, the present study reveals that spatial segregation within the MPFC in computing values for self-regarding vs. other-regarding choice is critically involved in determining individual variability in prosociality. Furthermore, weaker segregation of self- and other-regarding value signals in the MPFC and stronger MPFC–striatal coupling are associated with being prosocial rather than being selfish. Despite having yet to be tested with more direct measures of altruism, our findings provide important insight into human prosociality/altruism. First, the shared neural representation for self- and other-regarding values found among prosocial individuals supports the view that altruism requires value extension from self to others, a process in which another person’s welfare becomes valuable (1, 16, 17, 26). Second, the other-regarding valuation process subserved by the VMPFC as a part of the cortico-striatal network in prosocial individuals emphasizes the automatic and intuitive nature of prosocial motivation. This finding is in line with recent perspectives that prosociality and morality are rooted in intuition acquired, formed, and maintained through socialization (50–53). Third, our findings provide neural evidence of social norms internalized within an individual, which may have evolved to benefit groups by promoting prosocial behaviors (54). In conclusion, our present findings shed some light on the mystery of human altruism and support the notion that this mystery can be better understood by adopting rigorous scientific methods and theoretical frameworks in the ripening field of decision neuroscience.

Materials and Methods

Participants. Thirty pairs of healthy right-handed female college students participated in the experiment. The two participants in each pair were strangers to each other. One-half of the participants (one participant per pair; mean age, 21.9 y; range, 19–29 y) were randomly assigned to perform a prosocial learning task in the scanner. Four participants were excluded owing to excessive head movement or random responses in all of the conditions, leaving 26 participants included in the fMRI analysis. All participants were compensated with 30,000 KRW (~30 USD). The study protocol was approved by Korea University’s Institutional Review Board, and all participants provided written consent to participate before the start of the experiment.

Prosocial Learning Task. During the prosocial learning task, participants made choices between two fractal images, each of which was associated with different reward probabilities (30% vs. 70%). Each trial began with one of three pairs of fractal images, and each pair of images was associated with one of three different types of condition: SELF, BOTH, or OTHER (Fig. 1). Participants could earn two points for self and none for other in the SELF condition, one point for self and one point for other in the BOTH condition, and none for self and two points for other in the OTHER condition. The points earned from 10 randomly selected trials across all three conditions were to be used to reduce the duration of exposure to stressful noise (i.e., 10 s per point) for self and/or other. Each condition comprised 48 trials, resulting in a total of 144 trials in one functional run (~25 min). The conditions were presented in a pseudorandom order, and the conditions and

reward probabilities associated with different fractal images were counterbalanced across participants (SI Appendix).

Estimation of Chosen Value and RPE. The chosen value and RPE for each trial for each individual were estimated using the advantage learning model (3, 4) (SI Appendix).

Behavioral Measures of Individual Prosociality. To measure individual differences in prosociality within the task, we estimated an experienced magnitude of reward outcomes separately for each of the three conditions (SI Appendix). Applying the conceptual framework of social value orientation (55), we grouped participants into a prosocial group ($n = 15$), which valued other-regarding outcomes the same as or more than self-regarding outcomes, and a selfish group ($n = 10$), which valued self-regarding outcomes more than other-regarding outcomes. One subject who did not value either of the outcomes was excluded from the analyses examining group differences. To test the validity of our model-based categorization, we used participants’ self-reports in the social value orientation questionnaire (55) to group them into prosocial ($n = 14$) and prosel ($n = 10$) groups (SI Appendix). With this grouping, the behavioral and fMRI results remained the same (SI Appendix, Figs. S2 and S4).

fMRI Data Acquisition and Analysis. Brain images were acquired on a 3-T MRI scanner (MAGNETOM Tim Trio; Siemens Medical Solutions) at the Korea University Brain Imaging Center. T2*-weighted functional images were obtained through gradient echo planar imaging (EPI) with blood oxygenation level-dependent (BOLD) contrast [response time (TR) = 2,000 ms; echo time (TE) = 30 ms; flip angle = 90°; field of view (FOV) = 240 mm; 80 × 80 matrix; 36 axial slices; 3 × 3 × 3 mm in-plane resolution]. High-resolution T1-weighted structural images were collected as well (TR = 1,900 ms; TE = 2.52 ms; flip angle = 9°; 256 × 256 matrix; 1 × 1 × 1 mm in-plane resolution). The fMRI data were preprocessed and analyzed using SPM8 (Wellcome Department of Imaging Neuroscience, University College of London, London, UK). Images were realigned, normalized to the standard Montreal Neurological Institute EPI template, and spatially smoothed using a Gaussian kernel with an 8-mm full width at half maximum.

We created a first-level general linear model (GLM) with parametric modulators (SI Appendix). Trial-by-trial fluctuations of subject-specific chosen values and RPEs were estimated using the advantage learning model and entered into the first-level GLM model as parameters that modulated the hemodynamic responses at the time of option presentation and outcome presentation, respectively. Linear contrasts of regression coefficients for the parametric modulators of value and for RPE were computed and subjected to a random-effects group-level analysis using one-way ANOVA with condition (i.e., SELF, BOTH, and OTHER) as a repeated-measures factor. We quantified spatial gradients in self- and other-regarding value representations within the MPFC by extracting parameter estimates from five anatomical ROIs (spheres with a 4-mm radius) along the midline axis from the VMPFC to the DMPFC within the activation cluster correlating with the value parameters (30). For each ROI, we extracted the value parameter estimates for each condition, and then entered them into a repeated-measures one-way ANOVA. In addition, we fitted a linear slope to the OTHER vs. SELF contrasts across the five ROIs along the ventral-to-dorsal axis for each individual to estimate the degree of spatial gradient within the MPFC in terms of other-regarding vs. self-regarding valuation.

We assessed differential functional connectivity with the MPFC subregions during other-regarding compared with self-regarding choices at option presentation with a PPI analysis. The MMPFC seed was the peak voxel from the region ($x = 4, y = 52, z = 8$) found to be correlated with the value parameters from all three conditions in the parametric modulation analysis, and the VMPFC ($x = 0, y = 56, z = 2$) and DMPFC ($x = 2, y = 44, z = 12$) seeds were the peak voxels found to be correlated with the value parameters from SELF and OTHER conditions, respectively. The OTHER vs. SELF contrast at the time of option presentation was included as a psychological variable. Individual PPI maps were entered into a group-level two-sample t test comparing prosocial and selfish groups.

In addition, to examine whether prosocial and selfish individuals use different modes of decision for self vs. other, we contrasted neural responses to the presentation of options between SELF and OTHER conditions. The contrast maps of SELF vs. OTHER and OTHER vs. SELF at option presentation were entered into random-effects group-level two-sample t tests comparing prosocial and selfish groups.

All statistical thresholds were set to $P < 0.05$ corrected for multiple comparisons, using a cluster threshold determined at an uncorrected $P < 0.001$ by Monte Carlo simulations implemented in AlphaSim within AFNI software (afni.nimh.nih.gov/afni/) (56) for each search volume described

below. To assess value signals in the MPFC, we formed an a priori anatomical search volume that included the superior medial frontal cortex and anterior cingulate based on the AAL atlas (57) as implemented in the WFU_PickAtlas toolbox (www.ansir.wfubmc.edu) (58). In searching for RPE signals in the striatum and for the PPI analyses, we created an a priori search volume including the bilateral caudate and putamen (extending to nucleus accumbens), based on the AAL atlas. For the group comparison analyses of SELF vs. OTHER and OTHER vs. SELF contrasts at option presentation, the correction

was confined within the whole brain, because we had no specific hypothesis for this analysis.

ACKNOWLEDGMENTS. This work was supported by a research grant from Korea University (to S.S.), National Research Foundation of Korea grants from the Korean Government (NRF-2012-S1A3-A2033375, to S.S. and H.K.; 2006-2005110, to H.K.), and grants from the Swiss National Science Foundation (PP00P1_128574 and PP00P1_150739, to P.N.T.; CRSII3_141965, to E.F. and P.N.T.).

- Batson CD (2011) *Altruism in Humans* (Oxford Univ Press, New York).
- Thorndike EL (1911) *Animal Intelligence: Experimental Studies* (Macmillan, New York).
- Kim H, Shimojo S, O'Doherty JP (2006) Is avoiding an aversive outcome rewarding? Neural substrates of avoidance learning in the human brain. *PLoS Biol* 4(8):e233.
- O'Doherty J, et al. (2004) Dissociable roles of ventral and dorsal striatum in instrumental conditioning. *Science* 304(5669):452–454.
- Sutton RS, Barto AG (1998) *Introduction to Reinforcement Learning* (MIT Press, Cambridge, MA).
- 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(22):5623–5630.
- Haber SN, Knutson B (2010) The reward circuit: Linking primate anatomy and human imaging. *Neuropsychopharmacology* 35(1):4–26.
- Haruno M, Kawato M (2006) Heterarchical reinforcement-learning model for integration of multiple cortico-striatal loops: fMRI examination in stimulus-action-reward association learning. *Neural Netw* 19(8):1242–1254.
- Kim H, Adolphs R, O'Doherty JP, Shimojo S (2007) Temporal isolation of neural processes underlying face preference decisions. *Proc Natl Acad Sci USA* 104(46):18253–18258.
- Schultz W (1998) Predictive reward signal of dopamine neurons. *J Neurophysiol* 80(1):1–27.
- Lebreton M, Jorge S, Michel V, Thirion B, Pessiglione M (2009) An automatic valuation system in the human brain: Evidence from functional neuroimaging. *Neuron* 64(3):431–439.
- Plassmann H, O'Doherty J, Rangel A (2007) Orbitofrontal cortex encodes willingness to pay in everyday economic transactions. *J Neurosci* 27(37):9984–9988.
- Behrens TE, Hunt LT, Rushworth MF (2009) The computation of social behavior. *Science* 324(5931):1160–1164.
- Burke CJ, Tobler PN, Baddeley M, Schultz W (2010) Neural mechanisms of observational learning. *Proc Natl Acad Sci USA* 107(32):14431–14436.
- Christopoulos GI, King-Casas B (2015) With you or against you: Social orientation-dependent learning signals guide actions made for others. *Neuroimage* 104:326–335.
- Fehr E, Camerer CF (2007) Social neuroeconomics: The neural circuitry of social preferences. *Trends Cogn Sci* 11(10):419–427.
- Harbaugh WT, Mayr U, Burghart DR (2007) Neural responses to taxation and voluntary giving reveal motives for charitable donations. *Science* 316(5831):1622–1625.
- Izuma K, Saito DN, Sadato N (2008) Processing of social and monetary rewards in the human striatum. *Neuron* 58(2):284–294.
- Seid-Fatemi A, Tobler PN (2015) Efficient learning mechanisms hold in the social domain and are implemented in the medial prefrontal cortex. *Soc Cogn Affect Neurosci* 10(5):735–743.
- Chang SWC, Gariépy JF, Platt ML (2013) Neuronal reference frames for social decisions in primate frontal cortex. *Nat Neurosci* 16(2):243–250.
- Denny BT, Kober H, Wager TD, Ochsner KN (2012) A meta-analysis of functional neuroimaging studies of self- and other judgments reveals a spatial gradient for mentalizing in medial prefrontal cortex. *J Cogn Neurosci* 24(8):1742–1752.
- Mitchell JP, Macrae CN, Banaji MR (2006) Dissociable medial prefrontal contributions to judgments of similar and dissimilar others. *Neuron* 50(4):655–663.
- Moran JM, Lee SM, Gabrieli JDE (2011) Dissociable neural systems supporting knowledge about human character and appearance in ourselves and others. *J Cogn Neurosci* 23(9):2222–2230.
- Qin P, Northoff G (2011) How is our self related to midline regions and the default-mode network? *Neuroimage* 57(3):1221–1233.
- Saxe R, Moran JM, Scholz J, Gabrieli J (2006) Overlapping and non-overlapping brain regions for theory of mind and self-reflection in individual subjects. *Soc Cogn Affect Neurosci* 1(3):229–234.
- Hare TA, Camerer CF, Knoeplfle 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(2):583–590.
- Moll J, et al. (2006) Human fronto-mesolimbic networks guide decisions about charitable donation. *Proc Natl Acad Sci USA* 103(42):15623–15628.
- Shenhav A, Greene JD (2010) Moral judgments recruit domain-general valuation mechanisms to integrate representations of probability and magnitude. *Neuron* 67(4):667–677.
- Tricomi E, Rangel A, Camerer CF, O'Doherty JP (2010) Neural evidence for inequality-averse social preferences. *Nature* 463(7284):1089–1091.
- Nicolle A, et al. (2012) An agent independent axis for executed and modeled choice in medial prefrontal cortex. *Neuron* 75(6):1114–1121.
- Cacioppo JT, Petty RE, Kao CF (1984) The efficient assessment of need for cognition. *J Pers Assess* 48(3):306–307.
- Heatherton TF, et al. (2006) Medial prefrontal activity differentiates self from close others. *Soc Cogn Affect Neurosci* 1(1):18–25.
- Northoff G, et al. (2006) Self-referential processing in our brain—a meta-analysis of imaging studies on the self. *Neuroimage* 31(1):440–457.
- Sul S, Choi I, Kang P (2012) Cultural modulation of self-referential brain activity for personality traits and social identities. *Soc Neurosci* 7(3):280–291.
- Frith CD, Frith U (1999) Interacting minds—a biological basis. *Science* 286(5445):1692–1695.
- Hampton AN, Bossaerts P, O'Doherty JP (2008) Neural correlates of mentalizing-related computations during strategic interactions in humans. *Proc Natl Acad Sci USA* 105(18):6741–6746.
- Saxe R (2006) Uniquely human social cognition. *Curr Opin Neurobiol* 16(2):235–239.
- Bartra O, McGuire JT, Kable JW (2013) The valuation system: A coordinate-based meta-analysis of BOLD fMRI experiments examining neural correlates of subjective value. *Neuroimage* 76:412–427.
- Bzdok D, et al. (2013) Segregation of the human medial prefrontal cortex in social cognition. *Front Hum Neurosci* 7:232.
- Seo H, Cai X, Donahue CH, Lee D (2014) Neural correlates of strategic reasoning during competitive games. *Science* 346(6207):340–343.
- Wagner DD, Haxby JV, Heatherton TF (2012) The representation of self and person knowledge in the medial prefrontal cortex. *Wiley Interdiscip Rev Cogn Sci* 3(4):451–470.
- Verplanken B, Holland RW (2002) Motivated decision making: Effects of activation and self-centrality of values on choices and behavior. *J Pers Soc Psychol* 82(3):434–447.
- Jung D, Sul S, Kim H (2013) Dissociable neural processes underlying risky decisions for self versus other. *Front Neurosci* 7:15.
- Kang P, Lee J, Sul S, Kim H (2013) Dorsomedial prefrontal cortex activity predicts the accuracy in estimating others' preferences. *Front Hum Neurosci* 7:686.
- Yin HH, Knowlton BJ (2006) The role of the basal ganglia in habit formation. *Nat Rev Neurosci* 7(6):464–476.
- Aron AR, et al. (2007) Converging evidence for a fronto-basal-ganglia network for inhibitory control of action and cognition. *J Neurosci* 27(44):11860–11864.
- Corbetta M, Shulman GL (2002) Control of goal-directed and stimulus-driven attention in the brain. *Nat Rev Neurosci* 3(3):201–215.
- Heatherton TF (2011) Neuroscience of self and self-regulation. *Annu Rev Psychol* 62:363–390.
- Sridharan D, Levitin DJ, Menon V (2008) A critical role for the right fronto-insular cortex in switching between central-executive and default-mode networks. *Proc Natl Acad Sci USA* 105(34):12569–12574.
- Greene JD, Paxton JM (2009) Patterns of neural activity associated with honest and dishonest moral decisions. *Proc Natl Acad Sci USA* 106(30):12506–12511.
- Rand DG, Greene JD, Nowak MA (2012) Spontaneous giving and calculated greed. *Nature* 489(7416):427–430.
- Haidt J (2007) The new synthesis in moral psychology. *Science* 316(5827):998–1002.
- Sauer H (2012) Educated intuitions: Automaticity and rationality in moral judgement. *Philos Explor* 15:255–275.
- Sober E, Wilson DS (1998) *Unto Others: The Evolution and Psychology of Unselfish Behavior* (Harvard Univ Press, Cambridge, MA).
- Van Lange PAM, Otten W, De Bruin EM, Joireman JA (1997) Development of pro-social, individualistic, and competitive orientations: Theory and preliminary evidence. *J Pers Soc Psychol* 73(4):733–746.
- Cox RW (1996) AFNI: Software for analysis and visualization of functional magnetic resonance neuroimages. *Comput Biomed Res* 29(3):162–173.
- Tzourio-Mazoyer N, et al. (2002) Automated anatomical labeling of activations in SPM using a macroscopic anatomical parcellation of the MNI MRI single-subject brain. *Neuroimage* 15(1):273–289.
- Maldjian JA, Laurienti PJ, Kraft RA, Burdette JH (2003) An automated method for neuroanatomic and cytoarchitectonic atlas-based interrogation of fMRI data sets. *Neuroimage* 19(3):1233–1239.