



# Neurocomputational evidence that conflicting prosocial motives guide distributive justice

Yue Li<sup>a,b,1</sup>, Jie Hu<sup>a,c,1,2</sup> , Christian C. Ruff<sup>1</sup> , and Xiaolin Zhou<sup>a,b,d,e,2</sup>

Edited by René Marois, Vanderbilt University; received June 13, 2022; accepted October 18, 2022 by Editorial Board Member Michael S. Gazzaniga

In the history of humanity, most conflicts within and between societies have originated from perceived inequality in resource distribution. How humans achieve and maintain distributive justice has therefore been an intensely studied issue. However, most research on the corresponding psychological processes has focused on inequality aversion and has been largely agnostic of other motives that may either align or oppose this behavioral tendency. Here we provide behavioral, computational, and neuroimaging evidence that distribution decisions are guided by three distinct motives—inequality aversion, harm aversion, and rank reversal aversion—that interact with each other and can also deter individuals from pursuing equality. At the neural level, we show that these three motives are encoded by separate neural systems, compete for representation in various brain areas processing equality and harm signals, and are integrated in the striatum, which functions as a crucial hub for translating the motives to behavior. Our findings provide a comprehensive framework for understanding the cognitive and biological processes by which multiple prosocial motives are coordinated in the brain to guide redistribution behaviors. This framework enhances our understanding of the brain mechanisms underlying equality-related behavior, suggests possible neural origins of individual differences in social preferences, and provides a new pathway to understand the cognitive and neural basis of clinical disorders with impaired social functions.

striatum | frontostriatal circuitry | decision-making | distributive justice | prosocial motives

Most proposals for structuring human societies—from Aristotle's *Nicomachean Ethics* to *Marxism* and the *Declaration of Independence*—highlight that the pursuit of fairness and equality is a cornerstone of social justice and is essential for productive coexistence and collaboration (1). Fairness principles not only affect everyone's individual situation (e.g., work income) but also shape collective political ideology and social welfare (e.g., taxation and health-resource distribution policies) (2, 3). In line with this universal importance, people usually approach issues of distributive justice from the perspective of fairness norms (4), which are considered to be the most fundamental principle by which humans distribute resources (5, 6). This view is increasingly supported by evidence that people not only help disadvantaged parties to gain more equally distributed outcomes (7, 8) but also punish fairness norm violators (9–12).

However, fairness norms and inequality aversion alone cannot fully account for choices in situations requiring resource redistribution, which often reflect different motives (5). Imagine that two colleagues have made similar contributions to a project, but their employer gave one of them 1,000 dollars as bonus and the other only 100 dollars (A: \$1000 / B: \$100). Most people would feel frustrated by such an unequal distribution (9, 13) and would be willing to help the disadvantaged colleague (6, 14), albeit within certain limits. For example, most people would be happy to transfer 200 dollars from the advantaged to the disadvantaged (A: \$800 / B: \$300) but would be reluctant to transfer 700 dollars since this would reverse the initial rankings of each party (A: \$300 / B: \$800). This gives an example of the core motive conflicts in distributive justice, which in real life often lead to intense debates, e.g., on how to increase taxation on wealthy people while at the same time protecting everyone's interests and maintaining social order (15). This real-life example emphasizes the necessity to explore the boundaries of inequality aversion and to understand the natural limits of what people would do in the name of "fairness" (16, 17).

In situations like the above dilemma, and taxation debates in general, a primary aim is to reduce social inequality. However, this always involves trade-offs between inequality aversion and at least two other motives that support the status quo—harm aversion (2, 18) and rank reversal aversion (8). Specifically, moral decision studies suggest that people generally take into account the "do-no-harm" principle and tend to avoid helping one group at the expense of harming another group, even when the benefits outweigh the harm (2, 18). This entails that people are reluctant to redistribute wealth by transferring money

## Significance

Resource allocation in human societies usually triggers discussions about fairness, but satisfactory solutions to distribution problems also involve other prosocial motives that may prescribe different actions. Here, we address how the human brain mitigates such conflicts between multiple prosocial motives (fairness, harm aversion, and rank reversal aversion) during wealth distribution. Combining an experimental paradigm with fMRI and integrated neurocomputational modeling, we show that different prosocial motives are separately represented and integrated into choices by neural activity in striatum and its interactions with different brain regions. These findings extend unidimensional economic theories of third-party social preferences, characterize biological bases for individual and contextual differences in resource distribution behavior, and have economic and political implications for the design of taxation policies.

Author contributions: Y.L., J.H., and X.Z. designed research; Y.L. and J.H. performed research; Y.L., J.H., C.C.R., and X.Z. analyzed data; and J.H., C.C.R., and X.Z. wrote the paper.

The authors declare no competing interest.

This article is a PNAS Direct Submission. R.M. is a guest editor invited by the Editorial Board.

Copyright © 2022 the Author(s). Published by PNAS. This article is distributed under [Creative Commons Attribution-NonCommercial-NoDerivatives License 4.0 \(CC BY-NC-ND\)](https://creativecommons.org/licenses/by-nc-nd/4.0/).

<sup>1</sup>Y.L. and J.H. contributed equally to this work.

<sup>2</sup>To whom correspondence may be addressed. Email: [hujie0223@gmail.com](mailto:hujie0223@gmail.com) or [xz104@pku.edu.cn](mailto:xz104@pku.edu.cn).

This article contains supporting information online at <https://www.pnas.org/lookup/suppl/doi:10.1073/pnas.2209078119/-/DCSupplemental>.

Published November 29, 2022.

from the advantaged to the disadvantaged party (19, 20). Supporting this tendency, people are averse to overturn stable hierarchies in a society even though such preexisting hierarchies may conflict with their inequality aversion (21, 22). During wealth redistribution, it is widely observed that people are anchored to the initially unequal distribution and support such inequality to avoid reversal of preexisting income rankings (8). Thus, while harm aversion and rank reversal aversion can be seen as prosocial motives (in that they promote social welfare), they can work against inequality aversion and deter people from pursuing equality.

To establish the boundaries of these different motives, we have to uncouple them and examine how each of them contributes to redistribution behaviors in situations where they are in conflict. However, previous studies often employed paradigms specialized to study each motive in isolation, potentially biasing participants to act in line with just one of them. For instance, since in most of the previous paradigms, participants either played as victims of unfair distributions (6, 14, 23) or played as irrelevant third-party to punish intentional norm violations (7, 24), motives to maximize one's own interests or to punish norm violators may have amplified observed inequality aversion in these situations. Moreover, due to the limitations of previous paradigms and econometric models (25, 26), it is difficult to differentiate harm aversion and rank reversal aversion from inequality aversion and to clarify how humans weigh between these motives to make redistribution decisions. The trade-off between these motives may challenge the basic assumption of many econometric social preference models that distribution behaviors depend on ultimate outcomes rather than the changes between the ultimate and initial outcomes (25, 26).

In the current study, we aim to develop an integrated approach to examine how inequality aversion, harm aversion, and rank reversal aversion interact with each other to guide wealth redistribution choices. Specifically, we present a paradigm and a modeling approach that allows us to establish the boundaries and relative strengths of each motive and to elucidate the neural mechanisms underlying their effects on redistribution. We employ functional magnetic resonance imaging (fMRI) to clarify how information relevant for the different motives is represented and integrated in the human brain when people make redistribution decisions. One hypothesis is that equality-related information may be represented in the reward system [e.g., striatum and VMPFC, (6, 27, 28)] and that individuals' preferences related to equality seeking can be predicted by this activity, as well as the connectivity strengths between these regions and other systems (e.g., prefrontal regions) (7, 14, 29). With respect to harm aversion and rank reversal aversion, the literature suggests that social cognition (e.g., temporal parietal junction (TPJ)) and executive control systems (e.g., prefrontal regions) may underlie expression of these motives, since these structures have been found to be associated with greater preferences to minimize others' loss or pain (30–32). Thus, TPJ and prefrontal cortex may be sensitive to information concerning harm to others, which may be expressed as harm aversion and rank reversal aversion.

After identifying the systems involved in representing the information relevant for each motive, we examined how these motives are weighed and coordinated in the brain to guide redistribution decisions. To this end, we focus on how neural systems representing the different signals interact with each other to affect decisions in line with the latent motives. This allows us to differentiate between two potential scenarios regarding the motive-weighing process. On the one hand, while similar neural responses to equality signals have been observed in the striatum across different contexts, the connectivity of striatum with other brain regions has varied (6, 14). Therefore, one possible scenario is that equality signals are represented invariantly in the human brain, but conveyed differentially

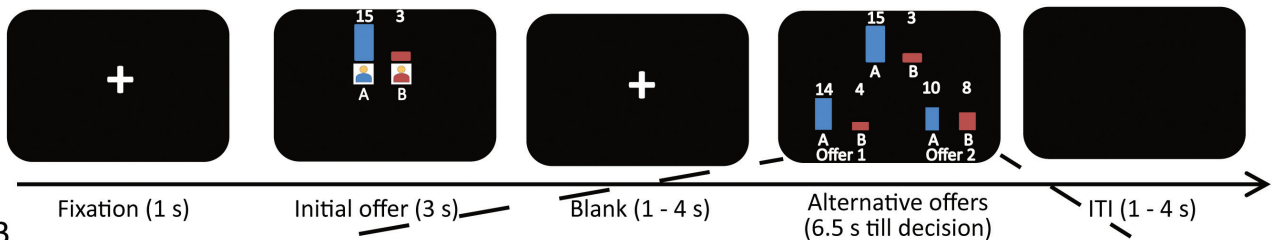
to other systems during conflicts with other motives (Scenario 1: Conflict gating of equality signals). On the other hand, previous studies have suggested that neural sensitivity to equality signals can depend on how strongly individuals weigh equality and that equality signals may in fact only be expressed when individuals' decisions are actually guided by equality (33). Therefore, neural equality representations may vary in their strength when other motives conflict with inequality aversion (Scenario 2: Conflict modulation of equality signals).

To address these questions, we developed a redistribution game that allowed us to measure individuals' inequality aversion, harm aversion, and rank reversal aversion during wealth redistribution. In the redistribution game, the participant played as a third-party to redistribute wealth between two anonymous strangers. They were first presented with a monetary distribution offer between two strangers (e.g., initial offer: Person A: ¥15, Person B: ¥3) and were told that these initial endowments were allocated randomly by a computer. They could choose between two alternative offers to reach a more equal distribution. Critically, we included two conditions: In the No Rank-reversal condition, the two alternative offers were both more equal than the initial offer but maintained the payoff ranking across the initially advantaged and disadvantaged person (e.g., Offer 1: Person A: ¥14, Person B: ¥4; Offer 2: Person A: ¥10, Person B: ¥8). In the Rank-reversal condition, by contrast, participants were presented with the same initial offer and the same more unequal alternative offer (e.g., Offer 1: Person A: ¥14, Person B: ¥4), but with a different alternative offer (e.g., Offer 2: Person A: ¥8, Person B: ¥10) that had the same inequality level as the alternative in the No Rank-reversal condition but that reversed the initially relative rankings (Fig. 1 *A* and *B*). If redistribution decisions are only driven by inequality aversion, people will choose the more equal offer more often regardless of whether or not the more equal offer will reverse the initially relative rankings. But if harm aversion and rank reversal aversion are at play, people will choose the more equal offer less often in the Rank-reversal condition than No Rank-reversal condition. This allows us to capture harm aversion (via participants' decision weights on how much money is taken away from the advantaged party) and rank reversal aversion (by a binary weight on choices that would reverse the initial rankings). We set up the offer matrix carefully so that the different motives were uncorrelated across trials, and our paradigm and model could capture the effects of each motive (for details, see *SI Appendix, SI Materials and Methods*). Combining this paradigm with computational modeling and fMRI approaches, we could thus distinguish the effects of the different motives on redistribution behaviors and the underlying neural mechanisms.

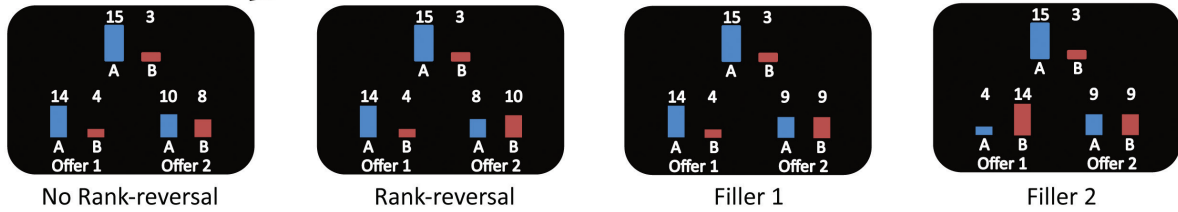
## Results

**Model-Free Results.** To examine whether individuals' choices reflected only inequality aversion or also the other motives, we first performed a generalized mixed-effect regression of redistribution choices on inequality differences between the two offers (relevant for inequality aversion), amount taken from the initially more advantaged player in the alternative offer (relevant for harm aversion), and a binary variable indicating whether the alternative offer would reverse the rank between players (relevant for rank reversal aversion) (see *SI Appendix, SI Materials and Methods* for model details). We estimated the effect of each predictor on choices of the more equal offer with the odds ratio estimate (ORE); ORE values between 0 and 1 (larger than 1) indicate a decrease (increase) in choice likelihood for a unit change in the predictor. This showed that participants' redistribution decisions indeed depended on these three factors (*SI Appendix, Table S1*): Indicative of inequality aversion, participants chose the more equal offer more frequently when it more strongly

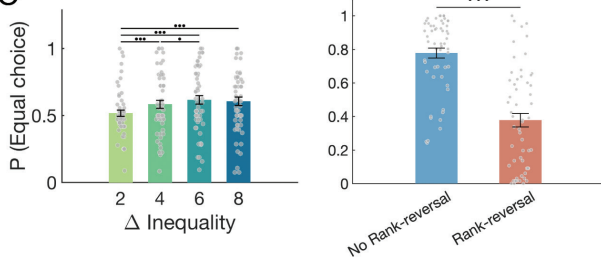
## A Experimental procedure



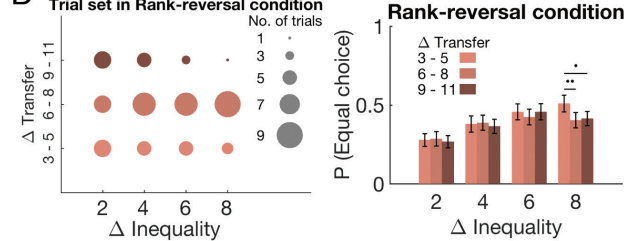
## B



## C



## D



**Fig. 1.** Experimental design and behavioral results. (A) Experimental task. In each trial, participants were first presented with an unequal monetary distribution between two anonymous strangers. Then, they could redistribute money by choosing between two alternative offers. (B) Experimental design. In all the trials, both alternative offers were more equally distributed than the initial offer. In the No Rank-reversal condition (Leftmost panel), the initially advantaged party will still gain more than the initially disadvantaged party from both alternative offers. In the Rank-reversal condition (Middle left panel), the initially advantaged party will still gain more from the relatively more unequal alternative offer, but will gain less from the relatively more equal alternative offer than the initially disadvantaged party. The inequality levels for both alternative offers were matched across the No Rank-reversal and Rank-reversal condition. We also included two types of filler trials in which the more equal alternative offers were always equally divided between the two parties, but the more unequal alternative offers would keep the initially relative rankings in one type of filler trials (middle Right panel) and reverse it in the other type of filler trials (Rightmost panel). (C) Main effects of inequality difference (Left) and condition (Right) on probability of more equal choice.  $P$  (Equal choice) increased when the more equal offer reduced the inequality level more strongly (Left), and  $P$  (Equal choice) decreased in the Rank-reversal versus No Rank-reversal condition (Right). Each dot represents one participant, and error bars represent SEMs.  $***, P < 0.001$ ;  $** , P < 0.01$ ;  $* , P < 0.05$ . (D) To differentiate the effect of inequality and the amount of transferred money (i.e., harm to the advantaged party), we orthogonalized the differences in inequality and the transferred money between the two alternative offers in the Rank-reversal condition. Left, x-axis represents the difference in inequality level ( $\Delta$  Inequality), and y-axis represents the difference in transferred money ( $\Delta$  Transfer) between the two offers. The size of the circle is proportional to the number of trials in each type of  $\Delta$  Inequality- $\Delta$  Transfer combination. Right, behavioral results in the Rank-reversal condition, probability of more equal choice [ $P$  (Equal choice)] is depicted as a function of  $\Delta$  Inequality and  $\Delta$  Transfer. For this post-hoc analysis, trials were divided for each level of  $\Delta$  Inequality into low (i.e., 3–5), middle (i.e., 6–8), and high levels (i.e., 9–11) of  $\Delta$  Transfer.

reduced the inequality level (effect of  $\Delta$  Inequality with ORE = 1.58, 95% CI [1.37–1.83],  $P < 0.001$ , Fig. 1C, Left and *SI Appendix, Table S2*) and when the initial inequality was greater (effect of  $\Delta$  Initial endowment with ORE = 1.12, 95% CI [1.01–1.24],  $P = 0.04$ , *SI Appendix, Fig. S2A and Table S3*). However, individuals' probability to choose the more equal offer was lower in the Rank-reversal condition than in the No Rank-reversal condition (ORE = 0.37, 95% CI [0.33–0.42],  $P_{\text{No Rank-reversal}}(\text{Equal}) = 0.78 \pm 0.03$  (MEAN  $\pm$  SE),  $P_{\text{Rank-reversal}}(\text{Equal}) = 0.38 \pm 0.04$ ,  $t(56) = 8.88$ ,  $P < 0.001$ , Fig. 1C, Right), demonstrating that rank reversal aversion influences choices independently from inequality considerations (which were matched across the two conditions). Importantly, participants also chose the more equal offer less frequently when it entailed larger transfers of money from the advantaged to the disadvantaged party (effect of  $\Delta$  Transfer with ORE = 0.46, 95% CI [0.43–0.50],  $P < 0.001$ , *SI Appendix, Fig. S2B and Table S4*), showing that harm aversion also affected choices on top of rank reversal aversion. This was also evident in a two-way  $\Delta$  Inequality \*  $\Delta$  Transfer interaction (ORE = 0.69, 95% CI [0.50–0.96],  $P = 0.03$ ), and a three-way  $\Delta$  Inequality \*  $\Delta$  Transfer \* condition interaction (ORE = 1.44, 95% CI [1.16–1.79],  $P < 0.001$ ).

To visualize and examine the patterns of the effects in the big regression model, we divided all trials based on condition,

$\Delta$  Inequality and  $\Delta$  Transfer, and inspected how individuals' choices varied as functions of these variables. Since we had orthogonalized the differences in initial endowment and in transfer/inequality between the two alternative offers, the effects reported here are not confounded by the effect of initial endowment (please see Fig. 1D, Left and *SI Appendix, Fig. S1*). These post-hoc tests confirmed that harm aversion had a stronger effect on redistribution for higher levels of inequality difference (i.e.,  $\Delta$  Inequality = 8,  $t_{\Delta \text{Transfer: low vs middle, } \Delta \text{Inequality} = 8} = 2.71$ ,  $p_{\Delta \text{Transfer: low vs middle, } \Delta \text{Inequality} = 8} = .009$ ;  $t_{\Delta \text{Transfer: low vs high, } \Delta \text{Inequality} = 8} = 2.36$ ,  $p_{\Delta \text{Transfer: low vs high, } \Delta \text{Inequality} = 8} = .022$ , Fig. 1D, Right, *SI Appendix, Table S5*). In addition, we observed a significant  $\Delta$  Initial endowment \*  $\Delta$  Transfer interaction (ORE = 1.19, 95% CI [1.10–1.28],  $P < 0.001$ ), suggesting that when the inequality level of the initial payoff is low, harm aversion leads people to choose the more equal offer less often when this requires greater monetary transfers and thus more strongly harms the initially advantaged person (*SI Appendix, Fig. S2C and Table S6*).

**Model-Based Results.** To better understand the effects of different motives on redistribution behaviors, we developed, fitted, and compared four families of computational models identifying how people weigh different motives to make redistributive decisions. We



**Table 1. Quality of model fits for computational models of redistribution decision-making**

Model	Equation	Parameters	LL	BIC	BF	Cross-validated prediction accuracy (Mean ± SE)
M1	$\Delta U = \alpha \Delta F$	$\alpha, \lambda$	-2,567	5,151	$6.33 \times 10^{60}$	0.478 ± 0.025
M2	$\Delta U = \alpha \Delta F - \delta$	$\alpha, \delta, \lambda$	-2,505	5,035	$4.09 \times 10^{35}$	0.536 ± 0.021
M3a	$\Delta U = \alpha \Delta F - \beta \Delta T - \delta$	$\alpha, \beta, \delta, \lambda$	-2,433	4,899	$1.20 \times 10^6$	0.744 ± 0.025
M3b	$\Delta U = \alpha \Delta F - \beta \Delta T$	$\alpha, \beta, \lambda$	-2,458	4,940	$9.62 \times 10^{14}$	0.726 ± 0.025
M3c	$\Delta U = \alpha \Delta F - \beta(\Delta T + \delta)$	$\alpha, \beta, \delta, \lambda$	-2,454	4,942	$2.61 \times 10^{15}$	0.738 ± 0.024
M4a	$\Delta U = \alpha \Delta F - \beta H - \delta$	$\alpha, \beta, \delta, \lambda$	-2,419	4,871	1	0.749 ± 0.025
M4b	$\Delta U = \alpha \Delta F - \beta H$	$\alpha, \beta, \lambda$	-2,457	4,938	$3.54 \times 10^{14}$	0.726 ± 0.025
M4c	$\Delta U = \alpha \Delta F - \beta(H + \delta)$	$\alpha, \beta, \delta, \lambda$	-2,434	4,900	$1.98 \times 10^6$	0.734 ± 0.024

$I_A(I_B)$ , payoff of the more unequal alternative offer for initially advantaged (disadvantaged) party;  $E_A(E_B)$ , payoff of the more equal alternative offer for initially advantaged (disadvantaged) party;  $\Delta F = |I_A - I_B| - |E_A - E_B|$ , difference in inequality between the two alternative offers;  $\Delta T = I_A - E_A$ , difference in transfer amount between the two alternative offers;  $H = E_B - E_A$ , harm for the initially advantaged party; all models have inverse temperature parameter  $\lambda$ ;  $\alpha$ , inequality aversion parameter;  $\beta$ , harm aversion parameter;  $\delta$ , rank reversal aversion parameter; LL, sum of log-likelihood over all participants; BIC, Bayesian information criterion over all participants; BF, Bayes factor. Models were estimated across all participants for model comparison.

focused these analyses on the Rank-reversal condition, which in contrast to the No Rank-reversal condition allowed us to differentiate inequality aversion from harm aversion and rank reversal aversion. We describe the principles and rationales of the four model families (M1–M4) in the following section and then report the results of the corresponding analyses. For detailed expositions of all the models and technical details of model selection and estimation, please see *SI Appendix, SI Materials and Methods* and Table 1.

**Model Construction.** The control model M1 only considered inequality aversion, whereas M2–M4 considered combinations of inequality aversion and the other motives.

The simplest model M1 followed the classical inequality aversion model proposed by Fehr and Schmidt (1999) in which people assign values to the outcomes of all parties but devalue the inequality they experience for any kinds of distribution. Model M2 quantified the additional effect of rank reversal aversion for the more equal offer, on top of inequality concerns.

Capturing the effects of harm aversion, on top of inequality concerns and rank reversal aversion, requires more complex model assumptions, which we embedded in different models that assumed different strategies of devaluing harms. In model family M3 (M3a–M3c), we assumed that people would devalue the utility of the alternative offer by the amount of money transferred from the initially advantaged party to the disadvantaged party. Therefore, in M3, in addition to the difference in inequality level and rank reversal, participants also weighted the difference in the amount of money transferred across the two parties between the two offers.

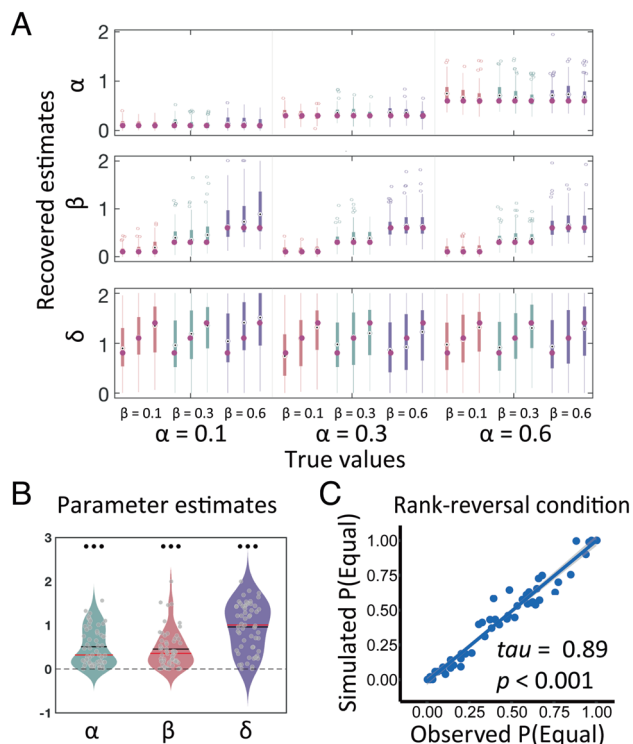
In model family M4 (M4a–M4c), we considered that people are not averse to transfer money away from the advantaged party as long as this transfer decreases the initial inequality level; but that they are averse to transferring more money than necessary to achieve a more equal payoff, which is the case in the Rank-reversal condition compared to the matched No Rank-reversal alternative. As shown in Fig. 1B, Offer 2 in the No Rank-reversal condition (A: ¥10, B: ¥8) transferred ¥5, and Offer 2 in the Rank-reversal condition (A: ¥8, B: ¥10) transferred ¥7 away from Person A's initial endowment. Therefore, the two types of Offer 2 achieved the same equality level (i.e., absolute payoff difference between parties), but the one in the Rank-reversal condition transferred “extra” money relative to the one in the No Rank-reversal condition (i.e., ¥7 – ¥5 = ¥2 in the above example). We thus considered this “extra” transferred money as unnecessary loss or harm for the initially advantaged party. Note that, with this assumption, it is not necessary for participants to memorize and compare the two counterpart equal offers between

the two conditions. Instead, they only needed to compare the more equal offer with a counterfactual offer in which the payoffs are flipped between the two parties. Therefore, this “extra” transferred money equals the payoff difference in the more equal offer. We referred to this amount of “extra” money as harm signal in following analyses. In M4, in addition to inequality and rank reversal, participants also weighted the harm signal for their choice.

For the model families M3 and M4, models within the same family calculated harm in the same way but assumed different types of devaluations of harm and rank reversal. Critically, M3a and M4a considered all the three components of inequality, harm, and rank reversal, M3b and M4b did not consider rank reversal, and M3c and M4c assumed that the harm aversion parameter captured effects of both the magnitudes of harm and rank reversal. For detailed expositions of the above models, see *SI Appendix, SI Materials and Methods* and Table 1.

**Model Comparison and Recovery.** Model comparison analyses first revealed that model M4a, which included the three components of inequality, harm, and rank reversal, outperformed all other models: It had the lowest BIC value (4,871 vs. 4,899 for the next model) and Bayes factors (BFs) relative to all alternative models that were higher than 100 (indicating very strong evidence favoring this model, *SI Appendix, SI Materials and Methods*, Table 1) (34). Detailed model comparison results, including model equations, free parameters, correspondence between parameters and cognitive components, log-likelihood, Bayesian information criterion (BIC), Bayes factor, and cross-validation prediction accuracy, are summarized in Table 1.

To ensure that the winning model could identify inequality aversion ( $\alpha$ ), harm aversion ( $\beta$ ), and rank reversal aversion ( $\delta$ ) in the Rank-reversal condition, we performed parameter recovery analysis. This showed that the three parameters in M4a could be recovered reliably and independently of each other (Fig. 2A), indicating that our paradigm and model could uncouple the effect of each motive on redistribution behaviors. Simulation analysis showed that the probability of more equal choice varied with all the three parameters (i.e.,  $\alpha$ ,  $\beta$ , and  $\delta$ , *SI Appendix, Fig. S3*), further confirming that different motives substantially affect individuals' decisions. Since the BIC scores for M3a and M4a were the closest, we performed model recovery to test how well data generated by model M3a and M4a could be recovered by each model. This revealed that choice data generated by M4a were more accurately recovered by M4a (prediction accuracy:  $0.92 \pm 0.01$ , MEAN ± SE) than by M3a ( $0.90 \pm 0.02$ ,  $t(53) = 3.00$ ,  $P = 0.004$ ), whereas



**Fig. 2.** Computational modeling results. (A) Parameter recovery results. Box plots show that the three parameters of the model M4a can be recovered reliably and independently of each other, indicating that our paradigm and model can clearly uncouple the effects of different motives on individuals' redistribution behaviors. We generated 27 datasets using all combinations of three plausible values for each parameter ( $\alpha$ : 0.1, 0.3, 0.6;  $\beta$ : 0.1, 0.3, 0.6; and  $\delta$ : 0.8, 1.1, 1.4). The boxes represent the distributions of the recovered parameters from 150 simulation sets of each combination of parameters. Each column corresponds to one type combination. Purple dots show the true values of the parameters. The recovered parameter distributions only vary with the true value of the parameter itself and not with the other parameters. (B) Violin plots show the distributions of the parameters of the winning model corresponding to different motives: inequality aversion ( $\alpha$ ), harm aversion ( $\beta$ ), and rank reversal aversion ( $\delta$ ). Black lines indicate the means, and red lines indicate the medians of the parameters. Each gray dot represents one participant. (C) Model simulation results. The scatter plot shows a strong correlation between observed probability of more equal choice and model simulated probability of more equal choice based on model M4a in the Rank-reversal condition.

there was no reliable difference in the accuracy with which choice data generated by M3a were recovered by M4a ( $0.84 \pm 0.02$ ) and M3a ( $0.82 \pm 0.02$ ,  $t(53) = 1.66$ ,  $P = 0.104$ ). Thus, the winning model M4a was indeed able to predict and capture unique aspects of the data compared to the closest alternative model.

**Model Parameters.** In line with the model-free analyses, model-based analyses confirmed that participants' redistribution behaviors in the Rank-reversal condition were driven by inequality aversion, harm aversion, and rank reversal aversion: Participants weighed the inequality difference between the two alternative offers ( $\alpha = 0.51 \pm 0.06$ ,  $t(56) = 8.90$ ,  $P < 0.001$ , Cohen's  $d = 1.18$ ), devalued the more equal offer by the extra harm for the initially advantaged party ( $\beta = 0.45 \pm 0.06$ ,  $t(56) = 7.83$ ,  $P < 0.001$ , Cohen's  $d = 1.04$ ), and valued rank reversal negatively ( $\delta = 0.96 \pm 0.07$ ,  $t(56) = 13.23$ ,  $P < 0.001$ , Cohen's  $d = 1.75$ , Fig. 2B). In line with expectations, greater inequality aversion ( $\alpha$ ) was associated with higher probability of more equal choice ( $\tau = 0.74$ ,  $P < 0.001$ , SI Appendix, Fig. S4, Left). By contrast, greater harm aversion ( $\beta$ ,  $\tau = -0.27$ ,  $P = 0.004$ ) and greater rank reversal aversion ( $\delta$ ,  $\tau = -0.63$ ,  $P < 0.001$ ) were associated with higher probability of more unequal choice (SI Appendix, Fig. S4, Middle and Right panels). Moreover, model simulation analyses showed that the choice probabilities predicted by the winning model

indeed captured the observed choice probabilities well ( $\tau = 0.89$ ,  $P < .001$ , Fig. 2C). Interestingly, inequality aversion ( $\alpha$ ) and rank reversal aversion ( $\delta$ ) were negatively correlated with each other ( $\tau = -0.62$ ,  $P < .001$ , SI Appendix, Fig. S5, Right). Given the posterior predictive checks and parameter recovery results, this correlation is very unlikely due to poor model performance and much more likely to indicate that more (less) inequality-averse participants indeed care less (more) about rank reversal.

Together, the model-based results suggest that people consider all three motives (inequality aversion, harm aversion, and rank reversal aversion) during wealth redistribution. Moreover, the specific form of the winning model M4a entails that people mainly consider offers harmful if these entail taking more money than would be necessary to reach a given equality level.

**Neuroimaging Results.** As our behavioral and modeling analyses suggested that participants jointly consider inequality aversion, harm aversion, and rank reversal aversion to make redistribution decisions, we investigated how these motives may be coordinated at the level of brain mechanisms. First, we clarified how each of these motives (e.g., equality and harm signals) is represented in the brain. To do so, we defined equality and harm signals based on the winning model (M4a) and inspected how these signals correlate with brain activity, either separately (general linear model 1, GLM1) or integrated into a common choice utility signals (GLM2). For these analyses, we focused on striatum and ventromedial prefrontal cortex (VMPFC), since these regions have been repeatedly suggested to be involved in equality and choice utility processing (6, 35). We also conducted exploratory analyses across the whole brain to identify other areas correlating with these signals. Second, we examined how these motives may interact to guide behavior, by investigating how the corresponding brain activity is functionally coupled, and how this relates to how strongly the motive is evident in the behavioral effects (psychophysiological interaction analyses PPIs, GLMs 3 and 4). That is, we tested whether neural responses to equality signals interacted with other regions related to harm processing or rank reversal, in a manner that correlates with the observed behavioral effects. These analyses were conducted at the whole-brain level, to identify any area that may show such functional interactions. Inference for all whole-brain analyses employed SnPM and used a cluster-level threshold of  $P < 0.05$  family-wise error (FWE) corrected for the whole brain, whereas region of interest (ROI) analyses were performed at a voxel-level  $P < 0.05$  FWE corrected for the ROI volume (see SI Appendix, SI Materials and Methods for details).

#### Striatum Represents Equality and Drives More Equal Choice.

We first examined how signals associated with inequality aversion and harm aversion were represented in the brain, by constructing a GLM 1 containing parametric regressors corresponding to equality in both conditions and harm ( $H$ ) in the Rank-reversal condition (see SI Appendix, SI Materials and Methods for details). We defined equality signals as  $-\Delta F = |E_A - E_B| - |I_A - I_B|$  so that higher equality values corresponded to smaller differences in inequality between the two alternative offers. The rationale for this definition was that people may perceive equality as something that is positively motivating and therefore assign increasingly larger values to more equal distributions. By contrast, when other motives conflict with equity-pursuing motives, responses to equality signals may be modulated, and motives to avoid harm may take over to guide decisions.

Our ROI analyses confirmed that activity in the striatum was related to equality. Specifically, activity in bilateral caudate/putamen (left peak MNI coordinates:  $[-18, 11, 1]$ , voxel-wise  $p(\text{FWE}) = 0.048$ ,  $t\text{-value} = 3.64$ ,  $k = 111$ ; right peak MNI

coordinates: [15, 20, -5], voxel-wise  $p(\text{FWE}) = 0.064$ ,  $t$ -value = 3.55,  $k = 76$ ) varied parametrically with equality ( $-\Delta F$ ) in the No Rank-reversal condition (Fig. 3A), but not in the Rank-reversal condition. A comparison between conditions confirmed a more positive striatal parametric effect of equality in the No Rank-reversal than Rank-reversal condition (peak MNI coordinates: [6, 14, -5], voxel-wise  $p(\text{FWE}) = 0.032$ ,  $t$ -value = 4.01,  $k = 45$ , Fig. 3B and *SI Appendix, Fig. S7* for a visualization of this effect). Note that this effect was also confirmed in the subsequent whole-brain analysis (*SI Appendix, Table S7*). The absence of striatum responses to equality in the Rank-reversal condition may be due to interactions between inequality aversion and the other motives that are stronger in this condition, a possibility that we tested explicitly in analyses described later.

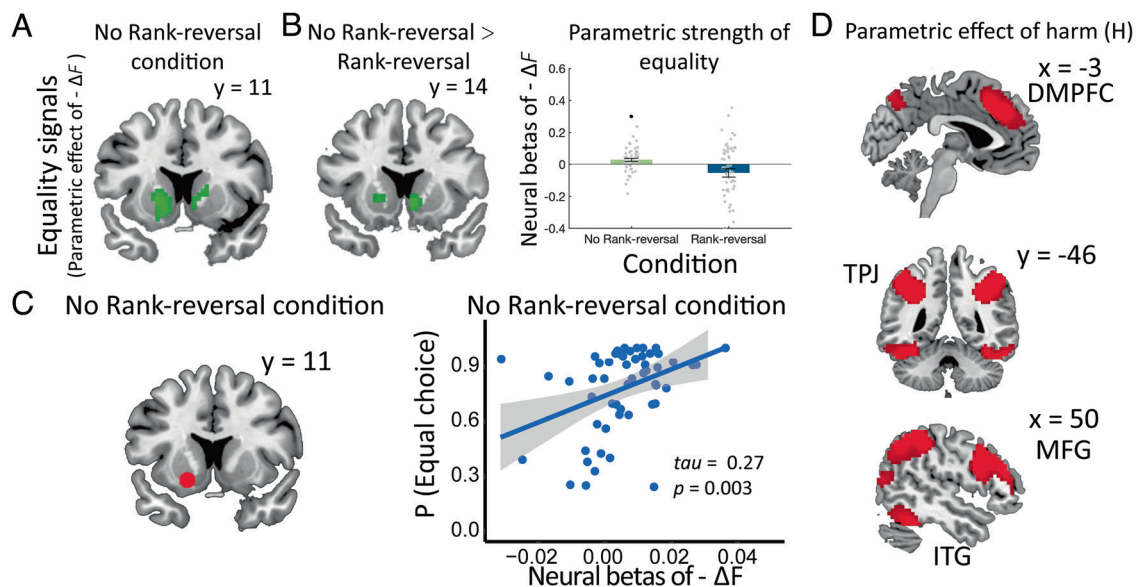
Our second ROI analysis showed that VMPFC was not involved in equality processing. However, consistent with prior studies (35, 36), this area (MNI peak coordinates: [3, 56, -14],  $t$ -value = 2.76, voxel-wise  $p(\text{FWE-SVC}) = 0.049$ ,  $k = 30$ , within VMPFC ROI with 8 mm radius centered on the peak MNI coordinates [0, 52, -8] involved in monetary incentive processing in ref. 35) was involved in representing the model-predicted value of the chosen option. This finding provides neural validation of our computational behavioral model.

Given that striatum was involved in signaling equality in the No Rank-reversal condition, we examined whether activity in this area can bias behavior in line with inequality aversion. A post-hoc correlation analysis showed that greater sensitivity to equality signals (i.e., more positive parametric estimates of  $-\Delta F$ ) in putamen (MNI peak coordinates: [-18, 11, -2], max  $t$ -value = 2.65, voxel-wise  $p(\text{FWE-SVC}) = 0.043$ ,  $k = 6$ , ROI center MNI coordinates [-12, 10, -6]) was indeed associated with a significantly higher probability of more equal choice in the No Rank-reversal condition (Kendall's  $\tau = 0.27$ ,  $P = 0.003$ , robust regression:  $b = 7.66$ ,  $P = 0.002$ , Fig. 3C) but not in the Rank-reversal condition (*SI Appendix, Fig. S8*). Whole-brain analyses revealed that no other region correlated with individuals' choices in either condition.

Taken together, these findings show that, in situations where inequality aversion is the main motive guiding behavior, the striatum plays a critical role in processing equality and biasing redistribution behaviors in line with these concerns.

**Cortical Regions Involved in Signaling Harm.** In the Rank-reversal condition, whole-brain analyses showed that activity in several brain areas correlated with the harm signals related to the more equal offer. These areas comprised dorsomedial prefrontal cortex/anterior cingulate cortex (DMPFC/ACC), inferior frontal gyrus (IFG), middle frontal gyrus (MFG), TPJ, and inferior temporal gyrus (ITG) (Fig. 3D and *SI Appendix, Table S7*). Thus, these areas could either represent the strength of the harm aversion motive, or they could be involved in processing/resolving the conflict between concerns about inequality and harm. The latter interpretation may be in line with previous findings that DMPFC/ACC, IFG, and MFG are often activated during cognitive control, conflict resolution, or behavioral adaptation (37, 38); and that TPJ is involved in mentalizing and perspective taking (39, 40). However, none of the neural effects in these areas were associated with the strength of behavioral harm aversion or inequality aversion, or the probability of more equal choice in the Rank-reversal condition. This motivated us to further examine whether and how the strength of the different motives was represented by interactions between the different neural systems representing harm and equality.

**DMPFC, as a Region Signaling Harm, Dampens Neural Sensitivity to Equality in Striatum.** We had observed weaker inequality aversion and dampened striatal sensitivity to equality in the Rank-reversal condition. These findings suggest that behaviorally relevant neural equality signals may not be represented invariably across different contexts, but may be modulated in situations where they conflict with harm signals. If this "conflict modulation" scenario held true, we should be able to observe that the reduction in striatal equality in the Rank-reversal condition relates to the strength of neural representations in harm-processing regions.



**Fig. 3.** Neural representations of equality and harm. (A) Activity in the striatum was associated with equality signals ( $-\Delta F$ ) in the No Rank-reversal condition. (B) More positive parametric strength of equality signal in the striatum in the No Rank-reversal than Rank-reversal condition (Left panel). For visualization, neural estimates of the significant cluster were extracted from both conditions (Right panel). Each dot represents one participant, and error bars indicate the SEMs.  $P < 0.05$ . (C) Scatter plot shows a correlation between the parametric strength of equality signal in the striatum (peak MNI coordinates [-18, 11, -2]) and individuals' probability of more equal choice in No Rank-reversal condition, suggesting that people whose striatum is more sensitive to equality have stronger preferences for more equal distribution. (D) Parametric effects of harm to the advantaged party in the Rank-reversal condition. Activity in DMPFC, TPJ, MFG, and ITG increased with the extent of harm to the advantaged party, suggesting processing of harm signals in these brain regions. Significant clusters are thresholded at voxel-wise  $P < 0.001$  uncorrected and cluster-wise FWE corrected  $P < 0.05$ . Correlation result in (C) is thresholded at voxel-wise  $P < 0.05$  FWE, small volume correction.



To test this hypothesis, we performed PPI analyses examining how interregional functional connectivity varies with inequality levels (GLMs 3 and 4; for ease of visualization  $-\Delta F$  was split into two bins (high  $-\Delta F$  vs. low  $-\Delta F$ ), but note that all effects are also present for a parametric regressor of  $-\Delta F$ ; for details, see *SI Appendix, SI Materials and Methods*). As the seed region for these analyses, we used an unbiased striatum region that was fully independent of the equality results described above (i.e., based on the peak coordinates in the Neurosynth “Striatum” activation map, Fig. 4A and *SI Appendix, SI Materials and Methods*). The PPI analyses were set up to identify brain regions that change their functional coupling with the striatum in line with how strongly equality concerns are relevant for the current choice. Evidence for this was assessed via the interaction term in the model, which quantifies for each voxel how much the correlation of the BOLD signal with that in the striatum changes as a function of the equality context (i.e., the equality concern triggered by the payoffs on the present trial), while simultaneously controlling for any main effects of (i.e., simple correlations with) the striatum time course and the equality context (41). These analyses revealed that dorsomedial prefrontal cortex (DMPFC, MNI peak coordinates: [0, 47, 40],  $k = 634$ ,  $t$ -value = 4.89, cluster-wise  $p$  (FWE) = 0.002) was functionally connected with striatum more strongly for high equality contexts (high  $-\Delta F$ ) in the Rank-reversal condition (Fig. 4C, *Left*; note that this effect was also present in control PPI analysis containing parametric inequality regressors; see *SI Results*). Importantly, the DMPFC region identified here largely overlapped with the DMPFC region involved in signaling harm to others (Fig. 4C, *Left*). A post-hoc comparison confirmed that this equality effect on DMPFC-Striatum connectivity was stronger in the Rank-reversal than No Rank-reversal condition (peak MNI coordinates: [3, 50, 34],  $t$ -value = 3.59, voxel-wise  $p$  (FWE-SVC) = 0.004,  $k = 63$ , ROI center MNI coordinates [0, 47, 40], Fig. 4C right, Rank reversal absence vs. presence).

To assess whether the pattern of DMPFC-Striatum connectivity may reflect functional influences on the striatum that change behavioral sensitivity to equality concerns, we tested for the Rank-reversal condition whether across individuals, a stronger effect of equality signals on DMPFC-Striatum connectivity may relate to a weaker striatum response to equality and a dampened tendency for equal choice. To this end, we extracted an index of neural equality sensitivity (Beta (high  $-\Delta F$ ) – Beta (low  $-\Delta F$ )) from the independent striatum seed region shown in Fig. 4A. As already shown in the initial ROI analyses described above, this index confirmed that the striatum was sensitive to equality in the No Rank-reversal condition, but not in the Rank-reversal condition (Fig. 4B). In line with the conjecture that DMPFC may act to dampen these striatal equality representations, the index of DMPFC-Striatum connectivity (Fig. 4C) exhibited the opposite pattern: It was stronger during Rank-reversal and weaker during the No Rank-reversal condition. Importantly, this effect was not just present on average but also at an individual level, since the differences between the Rank-reversal and No Rank-reversal condition in equality-related DMPFC-Striatum connectivity correlated negatively with the corresponding differences in neural equality sensitivity in the striatum ( $\tau = -0.18$ ,  $P = 0.050$ ; robust regression:  $b = -0.16$ ,  $P = 0.051$ , Fig. 4D). Thus, during rank reversal, the stronger the equality-related DMPFC-Striatum connectivity, the weaker the neural equality sensitivity in the striatum. Similar relations were also observed for behavior, since greater equality-related DMPFC-Striatum connectivity was also associated with a lower probability of more equal choice in the Rank-reversal relative to the No Rank-reversal condition ( $\tau = -0.16$ ,  $P = 0.082$ , robust regression:  $b = -0.32$ ,  $P = 0.043$ , Fig. 4E, *Top*), and with stronger harm aversion ( $\beta$ ,  $\tau = 0.18$ ,  $P = 0.044$ , robust regression:  $b = 0.48$ ,  $P = 0.028$ , Fig. 4E, *Bottom*).

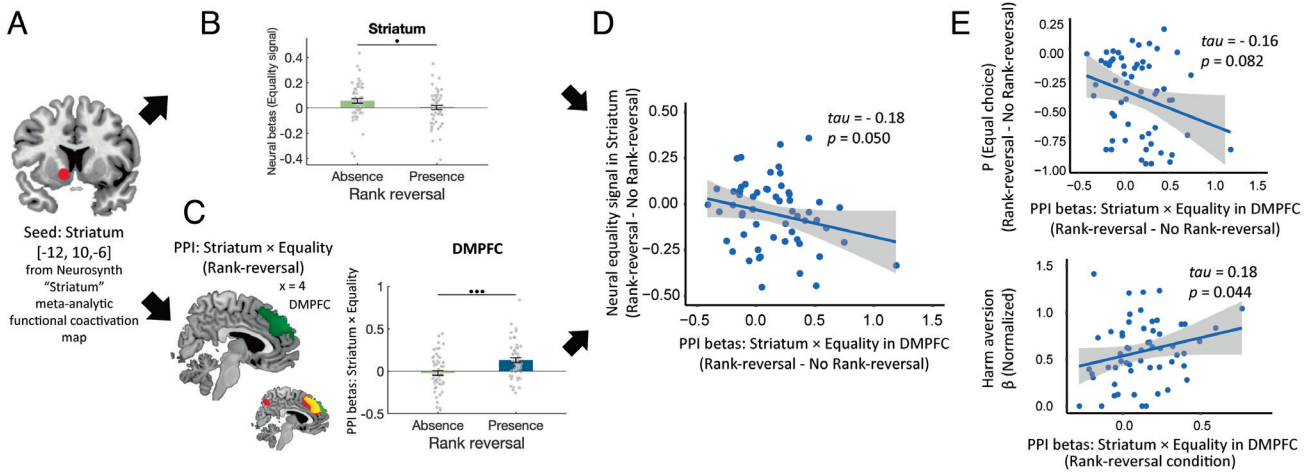
In sum, these findings suggest that harm aversion and equality concerns are balanced against each other in the brain by DMPFC-Striatum connectivity, since stronger connectivity relates to greater harm aversion, lower neural sensitivity to equality in the striatum, and less equal choices in situations where equity-seeking motives (i.e., inequality aversion) are in conflict with harm avoidance motives (i.e., harm aversion and rank reversal aversion). However, connectivity results alone do not allow unambiguous inferences about the directionality (from DMPFC to striatum or vice versa), modulatory type (inhibitory vs. modulatory), or causality of the functional interactions discovered here. Precise characterization of how the connectivity balances these motives thus requires further studies with other methods that offer more information and timing of such interactions.

**Activity in Cognitive Control and Harm-Related Areas during Unequal Choices.** We also examined how specific choice outcomes related to the trial-by-trial strength of neural motive representation, in terms of both regional activity and connectivity. This appears relevant for understanding what neural processes may lead individuals who are generally averse to inequality to nevertheless choose the more unequal offer on a specific trial. One possibility is that neural representations related to harm aversion or rank reversal aversion may be stronger for such choices. To address this issue, we modeled neural responses with respect to specific choices in each condition (i.e., four onset regressors: equal choice or unequal choice in the No Rank-reversal condition and Rank-reversal condition) and performed PPI analyses (GLM 4).

This revealed that in the No Rank-reversal condition, activity in MFG, IFG/Insula, ACC, MCC, TPJ, and IPL was increased when participants chose the more unequal offer (contrast: unequal choice > equal choice, Fig. 5A and *SI Appendix, Table S9*). No region was activated for the reverse contrast (contrast: equal choice > unequal choice) in the No Rank-reversal condition, or in any of the two contrasts in the Rank-reversal condition (*SI Appendix, Table S9*). Post-hoc interaction analyses between condition and choice (i.e., contrast: No Rank-reversal  $_{(\text{unequal choice} - \text{equal choice})} >$  Rank-reversal  $_{(\text{unequal choice} - \text{equal choice})}$ ) confirmed that activity in these regions was enhanced during unequal choices specifically in the No Rank-reversal condition, but not in the Rank-reversal condition (*SI Appendix, Table S10 and Fig. S9*). This suggests that overcoming equity concerns during unequal choices may involve high-level control and/or mentalizing processes implemented by frontal and parietal cortex.

An involvement of prefrontal and temporoparietal control processes in unequal choices was also suggested by correlation analyses of participants’ inequality and harm aversion parameters with brain activity related to unequal versus equal choice. This showed that the strength of activity in DMPFC ( $\tau = 0.29$ ,  $P = 0.003$ , robust regression:  $b = 0.09$ ,  $P = 0.002$ , Fig. 5B, *Left*) and TPJ ( $\tau = 0.36$ ,  $P < 0.001$ , robust regression:  $b = 0.12$ ,  $P < 0.001$ , Fig. 5B, *Middle*) during unequal choices was positively associated with inequality aversion (i.e.,  $\alpha$ ); and that the strength of activity in putamen was positively associated with harm aversion (i.e.,  $\beta$ ,  $\tau = 0.29$ ,  $P = 0.003$ , robust regression:  $b = 0.37$ ,  $P = 0.004$ , Fig. 5B, *Right* and *SI Appendix, Table S9*). Note that these results were robust to the exclusion of outliers ( $\tau(\text{DMPFC}-\alpha) = 0.24$ ,  $P = 0.018$ ;  $\tau(\text{TPJ}-\alpha) = 0.31$ ,  $P = 0.002$ ;  $\tau(\text{Putamen}-\beta) = 0.26$ ,  $P = 0.008$ ) and to statistical control for the effect of the other two parameters (see *SI Appendix, SI Results* for details). Activity in these three regions did not differ between unequal and equal choice in the Rank-reversal condition at group level (*SI Appendix, Fig. S10*).

Taken together, our analyses until now suggest that both DMPFC and TPJ are involved in harm signaling and that the same DMPFC region interacts with striatum to balance harm and equality motives



**Fig. 4.** Stronger DMPFC-Striatum connectivity associated with weaker neural equality signals in striatum and behavioral effects. (A) We focused the context-dependent analyses on a striatum region, with MNI coordinates  $[-12, 10, -6]$  which was derived from the “Striatum” mask at Neurosynth database. (B) We defined the neural equality signal as the difference in striatum BOLD signals between high  $-\Delta F$  (i.e.,  $-\Delta F = -2$  and  $-4$ ) and low  $-\Delta F$  (i.e.,  $-\Delta F = -6$  and  $-8$ ). These signals showed stronger equality sensitivity during absence of rank reversal (No Rank-reversal condition) than presence of rank reversal (Rank-reversal condition). (C) PPI analyses were performed to examine how connectivity with the striatum region in A changes with the contrast of “high  $-\Delta F >$  low  $-\Delta F$ .” These suggested a stronger DMPFC-Striatum connectivity effect of equality specifically in the Rank-reversal condition (Left panel, DMPFC in green), and this DMPFC region largely overlapped with the DMPFC region associated with harm signals (in red). The yellow area is the overlapping region. Post-hoc analyses confirmed a stronger effect of equality on PPI strength during the presence of rank reversal than absence of rank reversal. For visualization, we extracted the contrast value of the PPI regressors of the No Rank-reversal and Rank-reversal conditions within the significant cluster (Right panel). (D) Scatter plot shows that a stronger DMPFC-Striatum PPI strength of striatum\*equality is associated with a lower striatum neural sensitivity to equality in the Rank-reversal than No Rank-reversal condition. (E) Scatter plots show that stronger equality-related DMPFC-Striatum PPI connectivity is associated with a lower probability of more equal choice (Top panel) and with greater harm aversion ( $\beta$ ) (Bottom panel), in the Rank-reversal relative to No Rank-reversal condition. Each gray dot in (B) and (C) represents one participant, and error bars represent SEMs.  $\bullet\bullet\bullet$ ,  $P < 0.001$ ;  $\bullet\bullet$ ,  $P < 0.01$ ;  $\bullet$ ,  $P < 0.05$ . Significant clusters are thresholded at voxel-wise  $P < 0.001$  uncorrected and cluster-wise FWE corrected  $P < 0.05$ .

in the Rank-reversal condition. Congruent with these observations, we found that activity in DMPFC and TPJ was enhanced more strongly when more inequality-averse individuals chose the more unequal offer, again implying that harm-related activity in DMPFC and TPJ may deter more equal distributions, in particular for people who are averse to inequality.

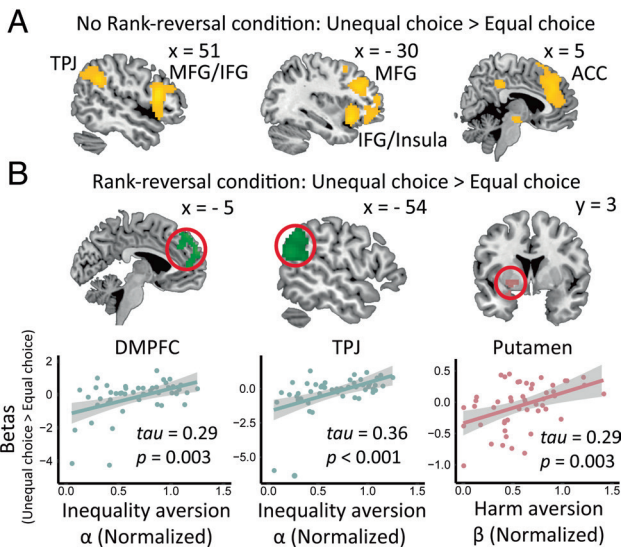
**Different Motives Affect Choice via Differential Patterns of Network Interactions.** The patterns of results until now suggest that inequality and harm aversion are implemented by different neural systems, which functionally interact with one another during redistribution choice. To test more directly for the relation between choice outcome and such network interactions, we performed PPI analyses focusing on the contrast between unequal choice and equal choice in the Rank-reversal condition and considered striatum (involved in equality processing) as the seed region. In particular, we examined how such network interactions may be expressed in individuals with strong behavioral expression of the different motives.

We examined two possibilities in this respect. First, for individuals with stronger inequality aversion to take unequal choices, harm- or rank-reversal-related neural activity may need to be recruited to interact with the striatum in a way that guides action selection according to context or individual preferences. Thus, in inequality-averse individuals, we should see stronger activity in harm- or rank-reversal-related neural systems and stronger connectivity with striatum during more unequal choices (see also refs. 31 and 42 for similar suggestions). Alternatively, individuals with strong harm and rank reversal aversion may exhibit more intense processing of the corresponding information and thus enhanced communication between the regions involved in these motives, reflecting more neural evidence about potential harm and rank reversal during more unequal choices.

In previous analyses, we have shown that the striatum (peak MNI coordinates  $[-18, 11, -2]$ ) was involved in equality

processing and equal choice in the No Rank-reversal condition, but we found no such effects in the Rank-reversal condition. In the current analysis, we thus explored whether this striatum region still interacted with other systems during unequal/equal choices in the Rank-reversal condition with motive conflicts, where striatal activity was not related to either equality processing or equal choice. We thus defined as ROI the striatum region involved in equality processing and equal choice in the No Rank-reversal condition (a sphere with 6-mm radius centered on peak MNI coordinates of  $[-18, 11, -2]$ ) and now examined with PPI analyses which areas show context-dependent connectivity with this area in the fully independent Rank-reversal condition, where equality was not neurally represented. This revealed that the connectivity strength between striatum and right IFG (peak MNI coordinates:  $[57, 23, 13]$ ,  $t$ -value = 5.08, cluster-wise  $p$  (FWE) = 0.046,  $k = 120$ , SI Appendix, Table S12) increased in people with greater inequality aversion when they chose the more unequal offer (i.e., normalized  $\alpha$ ,  $\tau = 0.38$ ,  $P < 0.001$ , Fig. 6 A and B, Left). This suggests that the striatum interacts with IFG more strongly when more inequality-averse individuals choose the more unequal offer in contexts where the more equal offer reverses ranks. Moreover, the connectivity strength between striatum and superior frontal gyrus (SFG, peak MNI coordinates:  $[-24, -1, 49]$ ,  $t$ -value = 5.35, cluster-wise  $p$  (FWE) = 0.041,  $k = 145$ , SI Appendix, Table S12) increased more strongly in people with greater rank reversal aversion when they chose the more unequal offer (i.e.,  $\delta$ ,  $\tau = 0.36$ ,  $P < 0.001$ , Fig. 6 A and B, Right), suggesting that conflicts between rank reversal aversion and equality-related motives during choice may be coordinated in the brain via neural connectivity between this SFG area and striatum. However, we note again that our connectivity analyses cannot provide conclusive evidence about directionality and modulatory nature of such interactions, preventing us from further speculation about the specific functional mechanisms underlying these effects. Note that although inequality aversion (i.e.,  $\alpha$ ) and rank reversal aversion (i.e.,  $\delta$ ) are





**Fig. 5.** Neural responses associated with more unequal choice link latent motives to behaviors. (A) In the No Rank-reversal condition, activity in MFG, IFG/Insula, ACC, and TPJ was enhanced when individuals chose the more unequal offer vs. more equal offer. (B) In the Rank-reversal condition, activity in DMPFC (Left panel) and TPJ (Middle panel) was enhanced when more inequality-averse individuals (i.e., higher  $\alpha$ ) chose the more unequal offer, whereas activity in putamen was enhanced when more harm-averse (i.e., higher  $\beta$ ) individuals chose the more unequal offer (Right panel). For visualization, neural estimates of the significant clusters were extracted, and scatter plots show the correlation patterns (Bottom panel). Significant clusters were thresholded at voxel-wise  $P < 0.001$  uncorrected and cluster-wise FWE corrected  $P < 0.05$ .

negatively correlated with each other, the findings that these two motives are related to differential connectivity patterns with striatum provide evidence that they function as two different motives that independently modulate neural circuitry underlying redistribution behaviors. The correlation patterns of the above networks also held after controlling for the effect of the other two model parameters (see *SI Appendix, SI Results* for details).

We did not observe striatal connectivity specifically associated with harm aversion in this analysis, but together with the observations of brain activity and connectivity associated with harm aversion shown in previous analyses, our findings emphasize that distinct neural pathways link different motives (inequality aversion, harm aversion, and rank reversal aversion) to redistribution behaviors, with striatum interacting with prefrontal areas in people with stronger aversion to inequality, harm, and rank reversal.

Together, our PPI results thus provide neural evidence that striatum connectivity is crucially involved in motive trade-offs from at least two perspectives. First, the strength of functional connectivity between the striatum (involved in equality processing) and DMPFC (involved in harm signaling) is associated with individuals' harm aversion, suggesting that this behavioral tendency relates to the functional communication between these two regions. Second, the striatum was related to equality responses and choices in the No Rank-reversal condition; and its connectivity with different frontal regions for more unequal choice was related to individuals' inequality aversion and rank reversal aversion in the Rank-reversal condition. This also implies that rank reversal aversion may interact with equality-related motives via striatal-prefrontal interactions during choices of (un)equal offers.

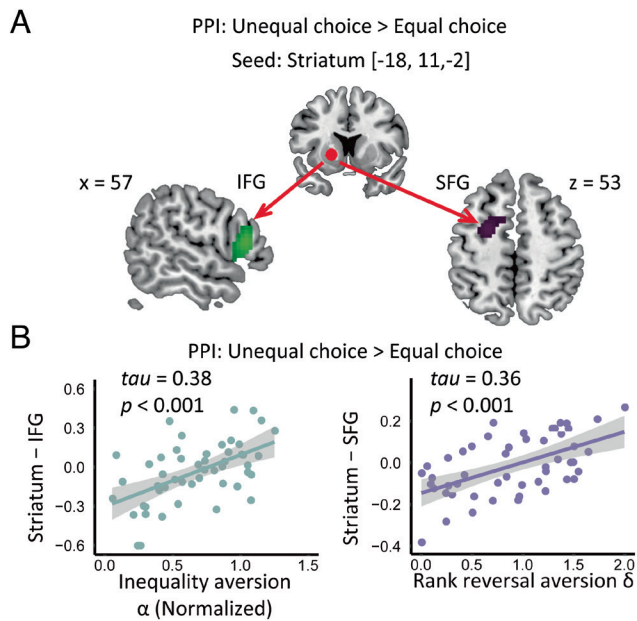
## Discussion

It is widely acknowledged that increased social inequality is associated with more risk-seeking behaviors, higher crime rate, and greater health problems (43, 44). Therefore, the question of how

to achieve distributive justice has become an intensively studied issue among researchers in many fields, including economics, politics, philosophy, and psychology. Although influential theories claim that fairness norms take precedence over other concerns (e.g., efficiency) underlying distributive justice (4), empirical evidence challenges this view and suggests that other motives can undermine fairness norms and deter equal distribution (5, 19). However, previous studies mainly focused on how self-interest motives may run counter to inequality concerns to affect wealth distribution, and most prevailing econometric models cannot explain why individuals can prefer greater inequality when different motives are in conflict (6, 25, 33). Although previous studies have demonstrated that harm aversion and rank reversal aversion are indeed involved in modulating moral decisions and redistribution decisions (8, 18, 31), it is still unclear how these motives interact with inequality aversion to bias individuals' choices.

Bridging these gaps, the current study establishes a redistribution paradigm and an integrated computational modeling approach to examine how conflicts between different prosocial motives bias individuals' preferences in wealth distribution. We demonstrate that harm aversion and rank reversal aversion can substantially interact with equality processing to prevent more equal distribution. Our neural results further suggest that the striatum serves as a hub for signaling equality and guiding decisions in line with equality concerns; and that striatal representations of equality may interact with other systems (e.g., frontal cortex) to drive choices when these are in conflict with harm avoidance and rank preserving motives.

Our study extends economic theories of social preferences by highlighting the trade-off between multiple prosocial motives in third-party wealth distribution and by exploring the boundaries within which inequality aversion determines wealth redistribution behavior. In the literature of third-party norms, theories often argue that people tend to punish norm violators in order to facilitate social norms (7, 45, 46). The current paradigm excludes the possibility of intentional violation of fairness norms, since the initially unequal distributions were generated from random draws. Given that participants still exhibit strong preferences for equal distribution in such situations, we suggest that inequality aversion, rather than motives to punish norm violation, drives redistribution behaviors as a core principle in wealth redistribution. However, we observed that people weighed equality less when it conflicted with preferences for harming others (i.e., harm aversion) or preserving initial rankings (i.e., rank reversal aversion), suggesting that equality-seeking motives (i.e., inequality aversion) are coordinated with other prosocial motives in wealth redistribution. Our results were gathered in the context of third-party preferences, so the question arises whether they would similarly apply to first-person contexts requiring people to allocate wealth between themselves and others. Previous studies suggest that similar mechanisms are at play in such contexts, but such studies have not yet clearly dissociated the different motives. For example, higher (lower) initial endowments will drive people to allocate more (less) wealth to themselves relative to others (19), and lower social ranking can also decrease individuals' inequality aversion strength and make them more willing to accept unfair offers (47). Thus, while people may also be averse to harm others or to reverse initial social ranking when making distributions for their own interests, these motives were often intertwined with self-interest and equality-seeking motives. Explicit evidence that our results would also apply to first-party preferences thus requires further empirical study. In general, our findings extend influential theories of fairness norms (25, 26) which mainly focused on effects of inequality aversion on distribution behaviors and emphasize the importance of considering other motives (i.e., harm aversion and



**Fig. 6.** Neural networks linking different motives to redistribution decisions. (A) In the Rank-reversal condition, Striatum-IFG connectivity strength was enhanced when more inequality-averse (higher  $\alpha$ ) individuals chose more unequal offers vs. more equal offers (Left panel), and Striatum-SFG connectivity strength was enhanced when more rank reversal-averse (i.e., higher  $\delta$ ) individuals chose more unequal offers vs. more equal offers (Right panel). Neural estimates of the significant clusters were extracted, and scatter plots show the correlation patterns (B). Significant clusters were thresholded at voxel-wise  $P < 0.001$  uncorrected and cluster-wise FWE corrected  $P < 0.05$ .

rank reversal aversion) in econometric models, especially since conflicts between these different motives are prevalent in real-life distribution decisions (e.g., taxation policy).

Harm aversion, as a critical type of moral virtue, drives people to achieve a more equal distribution by transferring as little money as possible between two parties. When making moral decisions, people typically conform to the “do-no-harm” principle and prefer not to benefit one party by harming another party (2, 18). Studies of morality suggest that people are not willing to take responsibility for others’ bad outcomes when making moral decisions (18, 48), as such moral responsibility will induce individuals’ anticipatory guilt emotion which proscribes people from harming others (30, 49). Therefore, taking more money away from others brings not only greater cost for the initially advantaged party but also greater cost of moral responsibility (i.e., harm aversion) for participants which will in turn dampen their motives to seek equality.

Moreover, we suggest that rank reversal aversion is another prosocial motive that discounts the utility of equality during wealth redistribution. A stable hierarchy can provide fitness advantage by satisfying individuals’ psychological need for order (50) and enhancing intragroup cooperation and productivity (51). Therefore, it is not surprising that people prefer to preserve rather than reverse preexisting hierarchy (8, 21). In line with these findings, our results suggest that the reversal of initial rankings also contributes to the disutility of equality when rank preserving and equality seeking are in conflict. Together, we demonstrate that in contrast to inequality aversion, harm aversion and rank reversal aversion function as two different third-party prosocial preferences to deter more equal wealth redistribution.

Our neural results first clarified how equality-related information is represented. GLM results support the hypothesis that individuals are sensitive to equality signals in the absence of any conflict but will be less sensitive to equality and base their decisions more heavily on other motives when they conflict with

inequality aversion. Although previous studies have proposed that the striatum signals rewarding aspects of equality-related distributions (5–7), it is still unclear which specific aspects of the distributions behavior engage the striatum and trigger the corresponding behavior—does it signal equality or other potentially rewarding aspects, such as efficiency or the other’s outcomes? While stronger activity in putamen was related to higher efficiency (i.e., greater overall profits) (5), efficiency cannot account for the pattern of results in the current study since neither of the two alternative offers changed the overall profits of the distributions. An alternative explanation is that striatum activity reflects dopaminergic responses in reward computation of social welfare, as it has been widely observed that stronger striatum activity is associated with charitable giving (52, 53), altruistic punishment to norm violation (23), and more equal wealth distributions (6, 7).

Moreover, striatum has been involved in arousal representations (54). For example, stronger striatal activation was related to greater motivation for norm compliance (55). In the current study, smaller equality difference between the two alternative offers may require participants to base their decisions more heavily on the evidence of equality signals and result in stronger motivation to comply with fairness norms for them, which is manifested by enhanced striatal activity. Together with the finding that greater sensitivity to equality in putamen was related to higher probability of more equal choice, our results suggest that striatum not only reflects processing of equality signals but also promotes fairness norm compliance.

Importantly, representations of equality in striatum were only observed in the No Rank-reversal condition, and this striatal signaling of equality was dampened in the context with conflicts between motives (i.e., Rank-reversal condition). Moreover, stronger DMPFC-Striatum connectivity was associated with lower equality sensitivity in striatum, less equal choice, and higher strength of harm aversion in the Rank-reversal condition. These findings help to clarify the neurocognitive mechanisms of the weighing processes of different motives, by providing a potential neural explanation for the weaker impact of equality on redistribution decisions in the Rank-reversal condition: DMPFC may process harm-related information, convey the harm aversion motive to striatum, interact with striatum, and dampen the tendency for more equal choice. Evidence from two lines of research supports such a modulating role of DMPFC. First, DMPFC, with adjacent regions ACC, is engaged in conflict monitoring, conflict resolution, and action selection in a variety of cognitive tasks (37, 38), which may support the resolution of conflict between different motives in the current paradigm. Second, DMPFC is also thought to be part of the mentalizing system that supports vicarious experiences of others’ pain or beliefs (39, 56), which may support harm signals in the current paradigm. In line with our findings, connectivity between prefrontal cortex and striatal value representations was also found to modulate individuals’ behaviors in other kinds of social and non-social decision-making (31, 57). However, despite the logical consistency of this interpretation, it is difficult to unambiguously infer the directionality and precise functional contributions of neural interactions from the results of PPI analyses. Future studies with brain stimulation may be needed to establish whether DMPFC influences on striatum are indeed causally involved in guiding redistribution behaviors under circumstances with conflicts between multiple motives.

Our results also provide crucial evidence for frontostriatal circuitry in redistribution decisions. The critical role of frontostriatal circuitry in decision-making has been highlighted in both social and non-social behaviors (31, 55, 57). In general, striatum is



suggested to receive inputs of goal-related representations from lateral prefrontal cortex and output value signals to guide response selection to maximize reward (58). In line with these suggestions, lateral prefrontal cortices are implicated in either modulating intuitive motivations or value representations that integrate information from different sources for moral and prosocial decision-making (31, 59). Our findings further refine previous accounts of frontostriatal circuitry in moral decision-making by clarifying that different prosocial motives modulate redistribution decisions through differential frontostriatal circuitries. Nevertheless, the specific functional contributions (i.e., inhibitory or modulatory) of these interactions between the striatal and frontal regions still need to be clarified in future studies.

Another critical contribution of our study is to clarify what neural processes underlie the modulations of different prosocial motives on redistribution decisions. Apart from processes involved in arbitrating between motives (i.e., DMPFC-Striatum connectivity), it is also important to identify processes that bias behavior on a trial-by-trial level in line with different motives and which may differ between people with different motive strengths. Activity in both DMPFC and TPJ was stronger when more inequality-averse individuals chose the more unequal offer, and activity in putamen was stronger when more harm-averse individuals chose the more unequal offer. One possibility suggested by the literature is that DMPFC and TPJ may support social cognitive processes such as mentalizing, perspective taking, inference, and learning about others' preferences (39, 56, 60). Recent studies further differentiated the roles of these two regions, by suggesting that while DMPFC is implicated in value-based action selection in a domain general manner (61–63), TPJ may be more specifically involved in processing of context-dependent social information (64, 65). Although our findings cannot provide a clear dissociation between DMPFC and TPJ, among all the regions involved in harm signaling, these two regions may be well-suited to link latent social motives to specific decisions. These findings also parallel the observation of stronger activity in TPJ for unequal choice vs equal choice in the No Rank-reversal condition, which may implicate the role of TPJ in social cognitive processing irrespective of whether there are conflicts between different motives.

In general, our findings may have economic, political, and social implications (66). The endowment effect has been introduced for decades to explain individuals' tendency to increase the subjective value of objects they own already (versus those they want to purchase) (67). Forgoing one's own good is seen as a kind of loss, and loss aversion will make it harder to give up the good (68, 69). In analogy to the endowment effect (70), our study highlights that people are inclined to maintain initial relative rankings and to take less money away from others in wealth redistribution, considering the reversal of initial rankings and others' monetary loss as a kind of third-party loss which proscribes actions to achieve higher equality (8). More generally, our findings may also explain resistance to reform policies that aim to promote social welfare or reduce income inequality (21, 71). For instance, rich people in regions with more equal income distribution, whose advantaged ranks can be more easily reversed, are less supportive of redistribution than those in regions with more unequal income distribution (16). Given that the effects of different motives are scientifically validated in the current study, this may help to develop better taxation policies by taking these motives into account when designing measures to reduce social inequality on the one hand and satisfy people in different income groups who pursue different motives on the other hand.

To conclude, the current study provides a neurocomputational account of the trade-off between multiple prosocial motives underlying resource distribution. Our findings suggest that in addition to inequality aversion, harm aversion and rank reversal aversion

work as two separate prosocial motives to modulate individuals' behaviors during wealth redistribution. Moreover, our study offers neural explanations for how different prosocial motives modulate redistribution behaviors, by highlighting a crucial role of striatum in equality processing and modulation of motives on ultimate decisions. Our approach improves our understanding of cognitive and neurobiological mechanisms underlying social preferences and distributive justice and may have implications for development of reform policies to promote fairness norms and social justice.

## Materials and Methods

**Participants.** Sixty-three right-handed healthy adults were recruited in the experiment. Six participants were excluded because of either making the same decision all the time or excessive head movement ( $> \pm 3$  mm in translation and/or  $> \pm 3^\circ$  in rotation). The remaining 57 participants were aged between 19 and 28 y (mean = 21.83 SD = 1.91; 31 female). No participant reported any history of psychiatric, neurological, or cognitive disorders. Informed written consent was obtained from each participant before the experiment. The study was carried out in accordance with the Declaration of Helsinki and was approved by the Ethics Committee of the Department of Psychology, Peking University.

**Experimental Procedure.** In the present study, we developed a redistribution task to assess individuals' preferences to redistribute unequal wealth allocations. In this task, participants were first presented with a monetary distribution scheme between two anonymous strangers. The initial endowment of each party was allocated unequally and randomly by computer, and participants had to choose between two redistribution options (i.e., alternative offers) which transferred a certain amount of money from the one with higher initial endowment (advantaged party) to the one with lower initial endowment (disadvantaged party, Fig. 1A). In the No Rank-reversal condition, both alternative offers were more equal than the initial offer and kept the same total payoffs and the same relative rankings between the two parties as the initial offer. While in the Rank-reversal condition, participants were presented with the same initial offer and the same more unequal alternative offer as the No Rank-reversal condition, but with a different more equal alternative offer that had the same inequality level as the more equal alternative offer in the No Rank-reversal condition but would reverse the initially relative advantageous/disadvantageous rankings of the two parties (Fig. 1B). There were 66 trials in each of the No Rank-reversal and Rank-reversal conditions and 15 trials in each of two filler conditions. The 162 trials were divided into three scanning sessions lasting ~15 min each. After the experiment, each participant received CNY 120 (~ USD 20) for compensation. For further details of the experimental paradigm, see *SI Appendix, SI Materials and Methods*.

**Computational Modeling Analyses.** To formalize different motives underlying redistribution behaviors, we performed model-based analyses by establishing four families of computational models to examine how inequality aversion, harm aversion, and rank reversal aversion affect individuals' redistribution behaviors in the Rank-reversal condition. For detailed modeling analyses, including model construction, estimation, comparison, and simulation, see *SI Appendix, SI Materials and Methods*.

**Neuroimaging Analyses.** We collected T2\*-weighted echo-planar images using a GE-MR750 3.0 T scanner with a standard head coil at Tongji University, China. The images were acquired in 40 axial slices parallel to the AC-PC line in an interleaved order, with an in-plane resolution of 3 mm  $\times$  3 mm, a slice thickness of 4 mm, an inter-slice gap of 4 mm, a repetition time of 2000 ms, an echo time of 30 ms, a flip angle of 90°, and a field of view of 200 mm  $\times$  200 mm. We used Statistical Parametric Mapping software SPM12 (Wellcome Trust Department of Cognitive Neurology, London, UK) which was run-through MATLAB (MathWorks) to preprocess the fMRI images, perform GLM analyses and PPI analyses. For detailed neuroimaging analyses, see *SI Appendix, SI Materials and Methods*.

**Data, Materials, and Software Availability.** Data (behavioral and fMRI) and customized MATLAB and R codes are available online (<https://osf.io/zd2tg/>).



**ACKNOWLEDGMENTS.** This study was supported by grants from the National Natural Science Foundation of China (31630034, 71942001). Dr. J.H. and Dr. C.C.R. also received funding from the European Research Council (ERC) under the European Union's Horizon 2020 research and innovation program (grant agreement no. 725355, BRAINCODES), from the UZH URPP AdaBD, and the SNSF (grant 100019L-173248).

1. K. Binmore, *Natural Justice* (Oxford University Press, 2005).
2. J. Baron, Blind justice: Fairness to groups and the do-no-harm principle. *J. Behav. Decis. Mak.* **8**, 71–83 (1995).
3. J. Offer, R. Pinker, Eds., *Social Policy and Welfare Pluralism* (Bristol University Press ed. 1, 2017) 10.2307/j.ctt22p7jvf.
4. J. Rawls, *A Theory of Justice* (Harvard University Press, 1999), 10.2307/j.ctvkjb25m.
5. M. Hsu, C. Anen, S. R. Quartz, The right and the good: Distributive justice and neural encoding of equity and efficiency. *Science* **320**, 1092–1096 (2008).
6. E. Tricomi, A. Rangel, C. F. Camerer, J. P. O'Doherty, Neural evidence for inequality-averse social preferences. *Nature* **463**, 1089–1091 (2010).
7. Y. Hu, S. Strang, B. Weber, Helping or punishing strangers: Neural correlates of altruistic decisions as third-party and of its relation to empathic concern. *Front. Behav. Neurosci.* **9**, 24 (2015).
8. W. Xie, B. Ho, S. Meier, X. Zhou, Rank reversal aversion inhibits redistribution across societies. *Nat. Hum. Behav.* **1**, 0142 (2017).
9. C. Corradi-Dell'Acqua, C. Civai, R. Rumiati, G. Fink, Disentangling self-and fairness-related neural mechanisms involved in the ultimatum game: An fMRI study. *Soc. Cogn. Affect. Neurosci.* **8**, 424–431 (2013).
10. J. J. Jordan, M. Hoffman, P. Bloom, D. G. Rand, Third-party punishment as a costly signal of trustworthiness. *Nature* **530**, 473–476 (2016).
11. E. Lo Gerfo *et al.*, The role of ventromedial prefrontal cortex and temporo-parietal junction in third-party punishment behavior. *Neuroimage* **200**, 501–510 (2019).
12. B. R. House *et al.*, Social norms and cultural diversity in the development of third-party punishment. *Proc. R. Soc. B* **287**, 20192794 (2020).
13. E. Xiao, D. Houser, Emotion expression in human punishment behavior. *Proc. Natl. Acad. Sci. U.S.A.* **102**, 7398–401 (2005).
14. R. Yu, A. J. Calder, D. Mobbs, Overlapping and distinct representations of advantageous and disadvantageous inequality. *Hum. Brain Mapp.* **35**, 3290–3301 (2014).
15. M. Iosifidi, N. Mylonidis, Relative effective taxation and income inequality: Evidence from OECD countries. *J. Eur. Soc. Policy* **27**, 57–76 (2017).
16. M. Dimick, D. Rueda, D. Stegmüller, The Altruistic rich? Inequality and other-regarding preferences for redistribution in the US. *Quart. J. Polit. Sci.* **11**, 385–439 (2016).
17. A. Argentiero, S. Casal, L. Mittone, A. Morreale, Tax evasion and inequality: Some theoretical and empirical insights. *Econ. Gov.* **22**, 309–320 (2021).
18. M. J. Crockett, Z. Kurth-Nelson, J. Z. Siegel, P. Dayan, R. J. Dolan, Harm to others outweighs harm to self in moral decision making. *Proc. Natl. Acad. Sci. U.S.A.* **111**, 17320–17325 (2014).
19. M. C. Leliveld, E. van Dijk, I. van Beest, Initial ownership in bargaining: Introducing the giving, splitting, and taking ultimatum bargaining game. *Pers. Soc. Psychol. Bull.* **34**, 1214–1225 (2008).
20. Y. Wu, J. Hu, E. van Dijk, M. C. Leliveld, X. Zhou, Brain activity in fairness consideration during asset distribution: Does the initial ownership play a role? *PLoS One* **7**, e39627 (2012).
21. R. Fernandez, D. Rodrik, Resistance to reform: Status quo bias in the presence of individual-specific uncertainty. *Am. Econ. Rev.* **81**, 1146–1155 (2004).
22. E. M. Zitek, L. Z. Tiedens, The fluency of social hierarchy: The ease with which hierarchical relationships are seen, remembered, learned, and liked. *J. Pers. Soc. Psychol.* **102**, 98–115 (2012).
23. D. J. F. De Quervain *et al.*, The neural basis of altruistic punishment. *Science* **305**, 1254–1258 (2004).
24. A. Strobel *et al.*, Beyond revenge: Neural and genetic bases of altruistic punishment. *Neuroimage* **54**, 671–680 (2011).
25. G. Charness, M. Rabin, Understanding social preferences with simple tests. *Q. J. Econ.* **117**, 817–869 (2002).
26. E. Fehr, K. Schmidt, A theory of fairness, competition and cooperation. *Q. J. Econ.* **114**, 817–868 (1999).
27. A. W. Cappelen *et al.*, Equity theory and fair inequality: A neuroeconomic study. *Proc. Natl. Acad. Sci. U.S.A.* **111**, 15368–15372 (2014).
28. D. J. de Quervain, U. Fischbacher, V. Treyer, M. Schellhammer, The neural basis of altruistic punishment. *World* **305**, 1–14 (2004).
29. L. Glass, L. Moody, J. Grafman, F. Krueger, Neural signatures of third-party punishment: Evidence from penetrating traumatic brain injury. *Soc. Cogn. Affect. Neurosci.* **11**, 253–262 (2015).
30. L. J. Chang, A. Smith, M. Dufwenberg, A. G. Sanfey, Triangulating the neural, psychological, and economic bases of guilt aversion. *Neuron* **70**, 560–572 (2011).
31. M. J. Crockett, J. Z. Siegel, Z. Kurth-Nelson, P. Dayan, R. J. Dolan, Moral transgressions corrupt neural representations of value. *Nat. Neurosci.* **20**, 879–885 (2017).
32. T. Nihonsugi, A. Ihara, M. Haruno, Selective increase of intention-based economic decisions by noninvasive brain stimulation to the dorsolateral prefrontal cortex. *J. Neurosci.* **35**, 3412–3419 (2015).
33. X. Gao *et al.*, Distinguishing neural correlates of context-dependent advantageous- and disadvantageous-inequity aversion. *Proc. Natl. Acad. Sci. U.S.A.* **115**, E7680–E7689 (2018).
34. S. Lewandowsky, S. Farrell, *Computational Modeling in Cognition* (2010), 10.1017/CBO9781107415324.004.
35. O. Bartra, J. T. McGuire, J. W. Kable, The valuation system: A coordinate-based meta-analysis of BOLD fMRI experiments examining neural correlates of subjective value. *Neuroimage* **76**, 412–427 (2013).
36. J. A. Clithero, A. Rangel, Informatic parcellation of the network involved in the computation of subjective value. *Soc. Cogn. Affect. Neurosci.* **9**, 1289–1302 (2014).
37. C. R. Oehrns *et al.*, Neural communication patterns underlying conflict detection, resolution, and adaptation. *J. Neurosci.* **34**, 10438–10452 (2014).
38. S. F. de Kloeet *et al.*, Bi-directional regulation of cognitive control by distinct prefrontal cortical output neurons to thalamus and striatum. *Nat. Commun.* **12**, (2021).
39. F. Van Overwalle, Social cognition and the brain: A meta-analysis. *Hum. Brain Mapp.* **30**, 829–58 (2009).
40. C. A. Hill *et al.*, A causal account of the brain network computations underlying strategic social behavior. *Nat. Neurosci.* **20**, 1142–1149 (2017), 10.1038/nn.4602.
41. K. J. Friston *et al.*, Psychophysiological and modulatory interactions in neuroimaging. *Neuroimage* **6**, 218–229 (1997).
42. A. C. Loewke, A. R. Minerva, A. B. Nelson, A. C. Kreitzer, L. A. Gunaydin, Frontostriatal projections regulate innate avoidance behavior. *J. Neurosci.* **41**, 5487–5501 (2021).
43. K. E. Pickett, R. G. Wilkinson, Income inequality and health: A causal review. *Soc. Sci. Med.* **128**, 316–326 (2015).
44. B. K. Payne, J. L. Brown-Iannuzzi, J. W. Hannay, Economic inequality increases risk taking. *Proc. Natl. Acad. Sci. U.S.A.* **114**, 4643–4648 (2017).
45. E. Fehr, U. Fischbacher, Social norms and human cooperation. *Trends Cogn. Sci.* **8**, 185–190 (2004).
46. E. Fehr, U. Fischbacher, Third-party punishment and social norms. *Evol. Hum. Behav.* **25**, 63–87 (2004).
47. J. Hu *et al.*, Social status modulates the neural response to unfairness. *Soc. Cogn. Affect. Neurosci.* **11**, 1–10 (2015).
48. I. Ritov, J. Baron, Reluctance to vaccinate: Omission bias and ambiguity. *J. Behav. Decis. Making* **3**, 263–277 (1990).
49. H. Yu, J. Hu, L. Hu, X. Zhou, The voice of conscience: Neural bases of interpersonal guilt and compensation. *Soc. Cogn. Affect. Neurosci.* **9**, 1150–1158 (2014).
50. J. P. Friesen, A. C. Kay, R. P. Eibach, A. D. Galinsky, Seeking structure in social organization: Compensatory control and the psychological advantages of hierarchy. *J. Pers. Soc. Psychol.* **106**, 590–609 (2014).
51. N. Halevy, E. Y. Chou, A. D. Galinsky, J. K. Murnighan, When hierarchy wins: Evidence from the national basketball association. *Soc. Psychol. Personal. Sci.* **3**, 398–406 (2012).
52. J. Moll *et al.*, Human fronto-mesolimbic networks guide decisions about charitable donation. *Proc. Natl. Acad. Sci. U.S.A.* **103**, 15623–15628 (2006).
53. W. T. Harbaugh, U. Mayr, D. R. Burghart, Neural responses to taxation and voluntary giving reveal motives for charitable donations. *Science* **316**, 1622–1624 (2007).
54. B. Knutson, C. M. Adams, G. W. Fong, D. Hommer, Anticipation of increasing monetary reward selectively recruits nucleus accumbens. *J. Neurosci.* **21**, 1–5 (2001).
55. M. Spitzer, U. Fischbacher, B. Herrnberger, G. Grön, E. Fehr, The neural signature of social norm compliance. *Neuron* **56**, 185–196 (2007).
56. M. M. Garvert, M. Moutoussis, Z. Kurth-Nelson, T. E. J. Behrens, R. J. Dolan, Learning-Induced plasticity in medial prefrontal cortex predicts preference malleability. *Neuron* **85**, 418–428 (2015).
57. W. van den Bos, C. A. Rodriguez, J. B. Schweitzer, S. M. McClure, Connectivity strength of dissociable striatal tracts predict individual differences in temporal discounting. *J. Neurosci.* **34**, 10298–10310 (2014).
58. J. W. Buckholtz, R. Marois, The roots of modern justice: Cognitive and neural foundations of social norms and their enforcement. *Nat. Neurosci.* **15**, 655–61 (2012).
59. J. Hu, Y. Hu, Y. Li, X. Zhou, Computational and neurobiological substrates of cost-benefit integration in altruistic helping decision. *J. Neurosci.* **41**, 3545–3561 (2021).
60. A. Ogawa, T. Kameda, Dissociable roles of left and right temporoparietal junction in strategic competitive interaction. *Soc. Cogn. Affect. Neurosci.* **14**, 1037–1048 (2020).
61. C. K. Kovach *et al.*, Anterior prefrontal cortex contributes to action selection through tracking of recent reward trends. *J. Neurosci.* **32**, 8434–8442 (2012).
62. E. D. Boorman, M. F. Rushworth, T. E. Behrens, Ventromedial prefrontal and anterior cingulate cortex adopt choice and default reference frames during sequential multi-alternative choice. *J. Neurosci.* **33**, 2242–2253 (2013).
63. J. X. O'Reilly *et al.*, Dissociable effects of surprise and model update in parietal and anterior cingulate cortex. *Proc. Natl. Acad. Sci. U.S.A.* **110**, E3660–E3669 (2013).
64. S. M. Lee, G. McCarthy, Functional heterogeneity and convergence in the right temporoparietal junction. *Cereb. Cortex* **26**, 1108–1116 (2016).
65. A. Konovalov, C. Hill, J. Daunizeau, C. C. Ruff, Dissecting functional contributions of the social brain to strategic behavior. *Neuron* **109**, 3323–3337.e5 (2021).
66. B. Irlenbusch, M. C. Villeval, Behavioral ethics: How psychology influenced economics and how economics might inform psychology? *Curr. Opin. Psychol.* **6**, 87–92 (2015).
67. Z. Carmon, D. Ariely, Focusing on the forgone: How value can appear so different to buyers and sellers. *J. Consum. Res.* **27**, 360–370 (2000).
68. D. Kahneman, J. L. Knetsch, R. H. Thaler, Anomalies: The endowment effect, loss aversion, and status quo bias. *J. Econ. Perspect.* **5**, 193–206 (1991).
69. C. K. Morewedge, L. L. Shu, D. T. Gilbert, T. D. Wilson, Bad riddance or good rubbish? Ownership and not loss aversion causes the endowment effect. *J. Exp. Soc. Psychol.* **45**, 947–951 (2009).
70. R. A. Y. Weaver, S. Frederick, A reference price theory of the endowment effect. *J. Mark. Res.* **XLIX**, 696–707 (2012).
71. I. Kuziemko, R. W. Buell, T. Reich, M. I. Norton, Last-place aversion: Evidence and redistributive implications. *Q. J. Econ.* **129**, 105–149 (2014).

Author affiliations: <sup>a</sup>School of Psychological and Cognitive Sciences and Beijing Key Laboratory of Behavior and Mental Health, Peking University, Beijing 100871, China; <sup>b</sup>PKU-IDG/McGovern Institute for Brain Research, Peking University, Beijing 100871, China; <sup>c</sup>Zurich Center for Neuroeconomics, Department of Economics, University of Zurich, Zurich 8006, Switzerland; <sup>d</sup>School of Business and Management, Shanghai International Studies University, Shanghai 200083, China; and <sup>e</sup>Shanghai Key Laboratory of Mental Health and Psychological Crisis Intervention and School of Psychology and Cognitive Science, East China Normal University, Shanghai 200062, China