

## Neurocomputational evidence that conflicting prosocial motives guide distributive justice

Yue Li<sup>a,b,1</sup>, Je Hu<sup>a,c,1,2</sup>, Christian C. Ruf<sup>c</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 motivesinequality 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 a ect 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). is 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 re ect di erent 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). is gives an example of the core motive con icts 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). is 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-o s between inequality aversion and at least two other motives that support the status quo—harm aversion (2, 18) and rank reversal aversion (8). Speci cally, 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 bene ts outweigh the harm (2, 18). is entails that people are reluctant to redistribute wealth by transferring money Resource allocation in human societies usually triggers discussions about fairness, but satisfactory solutions to distribution problems also involve other prosocial motives that may prescribe diferent actions. Here, we address how the human brain mitigates such conficts between multiple prosocial motives (fairness, harm aversion, and rank reversal aversion) during wealth distribution. Combining a experimental paradigm with fMRI and integrated neurocomputational modeling, we show that dif erent prosocial motives are separately represented and integrated into choices by neural activity in striatum and its interactions with dif erent brain regions. These fndings extend unidimensional economic theories of third-party social preferences, characterize biological bases for individual and contextual dif erences in resource distribution behavior, and have economic and political implications for the design of taxation policies.

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).

<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 or 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.

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.

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 con ict 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). us, 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 di erent motives, we have to uncouple them and examine how each of them contributes to redistribution behaviors in situations where they are in con ict. 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 ampli ed observed inequality aversion in these situations. Moreover, due to the limitations of previous paradigms and econometric models (25, 26), it is di cult to di erentiate harm aversion and rank reversal aversion from inequality aversion and to clarify how humans weigh between these motives to make redistribution decisions. e trade-o 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. Speci cally, 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 e ects on redistribution. We employ functional magnetic resonance imaging (fMRI) to clarify how information relevant for the di erent 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). us, 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 di erent signals interact with each other to a ect decisions in line with the latent motives. is allows us to di erentiate 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 di erent contexts, the connectivity of striatum with other brain regions has varied (6, 14).

erefore, one possible scenario is that equality signals are represented invariantly in the human brain, but conveyed di erentially to other systems during con icts with other motives (Scenario 1: Con ict 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). erefore, neural equality representations may vary in their strength when other motives conict with inequality aversion (Scenario 2: Con ict 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. ev were rst presented with a monetary distribution o er between two strangers (e.g., initial o er: Person A: ¥15, Person B: ¥3) and were told that these initial endowments were allocated randomly by a computer. ey could choose between two alternative o ers to reach a more equal distribution. Critically, we included two conditions: In the No Rank-reversal condition, the two alternative o ers were both more equal than the initial o er but maintained the payo ranking across the initially advantaged and disadvantaged person (e.g., O er 1: Person A: ¥14, Person B: ¥4; O er 2: Person A: ¥10, Person B: ¥8). In the Rank-reversal condition, by contrast, participants were presented with the same initial o er and the same more unequal alternative o er (e.g., O er 1: Person A: ¥14, Person B: ¥4), but with a di erent alternative o er (e.g., O er 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 o er more often regardless of whether or not the more equal o er will reverse the initially relative rankings. But if harm aversion and rank reversal aversion are at play, people will choose the more equal o er less often in the Rank-reversal condition than No Rank-reversal condition. is 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 o er matrix carefully so that the di erent motives were uncorrelated across trials, and our paradigm and model could capture the e ects of each motive (for details, see *SI Appendix* 

reduced the inequality level (e ect 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 (e) ect 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 o er was lower in the Rank-reversal condition than in the No Rank-reversal condition (ORE = 0.37, 95% CI [0.33 – 0.42],  $P_{No Rank-reversal}$  (Equal) = 0.78 ± 0.03 (MEAN ± SE),  $P_{Rank-reversal}$  (Equal) = 0.38 ± 0.04, t(56) = 8.88, P < 0.001, Fig. 1*C*, *Right*), demonstrating that rank reversal aversion in uences choices independently from inequality considerations (which were matched across the two conditions). Importantly, participants also chose the more equal o er less frequently when it entailed larger transfers of money from the advantaged to the disadvantaged party (e ect 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 a ected choices on top of rank reversal aversion. is 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 e ects 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 di erences in initial endowment and in transfer/inequality between the two alternative o ers, the e ects reported here are not confounded by the e ect of initial endowment (please see Fig. 1*D*, *Left* and *SI Appendix*, Fig. S1). ese post-hoc tests con rmed that harm aversion had a stronger e ect on redistribution for higher levels of inequality di erence (i.e.,  $\Delta$ Inequality = 8,  $t_{\Delta$ Transfer: low vs middle,  $\Delta$ Inequality = 8 (56) = 2.71,  $p_{\Delta$ Transfer: low vs middle,  $\Delta$ Inequality = 8 = .009;  $t_{\Delta}$ Transfer: low vs high,  $\Delta$ Inequality = 8 (56) = 2.36,  $p_{\Delta}$ Transfer: low vs high,  $\Delta$ Inequality = 8 = .022, Fig. 11 1 Tf0 Tw ( ) Tj, 1 1 Tf0 Tw ( 6 189.07

f

0

Model	Equation	Pa	arameters	LL	BIC	BF	Cross-validated prediction accuracy (Mean ± SE)
M1	α		αι	-2,567	5,151		0.478 ± 0.025
M2			$\alpha$ I $\delta$ ,	-2,505	5,035		$0.536 \pm 0.021$
M3a			αιβίδ,	-2,433	4,899		0.744 ± 0.025
M3b			$\alpha$ I $\beta$ ,	-2,458	4,940		$0.726 \pm 0.025$
M3c		Γ	αιβ[δ,	-2,454	4,942		$0.738 \pm 0.024$
M4a			$\alpha$ i $\beta$ l $\delta$ ,	-2,419	4,871	1	0.749 ± 0.025
M4b			$\alpha$ I $\beta$ ,	-2,457	4,938		$0.726 \pm 0.025$
M4c		[	<b>αι β[ δ</b> ,	-2,434	4,900		0.734 ± 0.024

 $\Box$  [, payof of the more unequal alternative of er for initially advantaged (disadvantaged) party;  $\Box$  [, payof of the more equal alternative of er for initially advantaged (disadvantaged) party;  $\Box$  [, payof of the more equal alternative of er for initially advantaged (disadvantaged) party;  $\Box$  [, payof of the more equal alternative of er for initially advantaged (disadvantaged) party;  $\Box$  [, payof of the more equal alternative of er for initially advantaged (disadvantaged) party;  $\Box$  [, payof of the more equal alternative of er for initially advantaged (disadvantaged) party;  $\Box$  [, dif erence in transfer amount between the two alternative of ers; , harm for the initially advantaged party; all models have inverse temperature parameter ;  $\alpha$ , 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 di erentiate 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.

e control model M1 only considered inequality aversion, whereas M2–M4 considered combinations of inequality aversion and the other motives.

e 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 quanti ed the additional e ect of rank reversal aversion for the more equal o er, on top of inequality concerns.

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

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 payo , which is the case in the Rank-reversal condition compared to the matched No Rank-reversal alternative. As shown in Fig. 1*B*, O er 2 in the No Rank-reversal condition (A: ¥10, B: ¥8) transferred ¥5, and O er 2 in the Rank-reversal condition (A: ¥8, B: ¥10) transferred ¥7 away from Person A's initial endowment.

erefore, the two types of O er 2 achieved the same equality level (i.e., absolute payo di erence 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 o ers between

the two conditions. Instead, they only needed to compare the more equal o er with a counterfactual o er in which the payo s are

ipped between the two parties. erefore, this "extra" transferred money equals the payo di erence in the more equal o er. 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 di erent 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 e ects 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 analyses rst 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 is showed that the three parameters in M4a could be analysis. recovered reliably and independently of each other (Fig. 2A), indicating that our paradigm and model could uncouple the e ect 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 con rming that di erent motives substantially a ect 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. is revealed that choice data generated by M4a were more accurately recovered by M4a (prediction accuracy:  $0.92 \pm 0.01$ , MEAN  $\pm$ SE) than by M3a (0.90  $\pm$  0.02, t(53) = 3.00, P = 0.004), whereas



Pank-reversal condition D - Parameter estimates

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 di erence in the accuracy with which choice data generated by M3a were recovered by M4a ( $0.84 \pm 0.02$ ) and M3a  $(0.82 \pm 0.02, t(53) = 1.66, P = 0.104).$ us, the winning model M4a was indeed able to predict and capture unique aspects of the data compared to the closest alternative model.

In line with the model-free analyses, modelbased analyses con rmed that participants' redistribution behaviors in the Rank-reversal condition were driven by inequality aversion, harm aversion, and rank reversal aversion: Participants weighed the inequality di erence between the two alternative o ers ( $\alpha = 0.51$  $\pm 0.06$ , t(56) = 8.90, P < 0.001, Cohen's d = 1.18), devalued the more equal o er 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 < 0.070.001, 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. 2*C*). 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 speci c form of the winning model M4a entails that people mainly consider o ers harmful if these entail taking more money than would be necessary to reach a given equality level.

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 clari ed how each of these motives (e.g., equality and harm signals) is represented in the brain. To do so, we de ned 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 e ects (psychophysiological interaction analyses PPIs, GLMs 3 and 4).

at 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 ese analyses were conducted at the whole-brain level, e ects. 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).

We rst 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 de ned equality signals as  $-\Delta F = |E_A - E_B| - |I_A - I_B|$  so that higher equality values corresponded to smaller di erences in inequality between the two alternative o ers. e rationale for this de nition 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 con ict 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 con rmed that activity in the striatum was related to equality. Speci cally, activity in bilateral caudate/ putamen (left peak MNI coordinates: [-18, 11, 1], voxel-wise p(FWE) = 0.048, t-value = 3.64, k = 111; right peak MNI coordinates: [15, 20, -5], voxel-wise p(FWE) = 0.064, t-value = 3.55, k = 76) varied parametrically with equality  $(-\Delta F)$  in the No Rank-reversal condition (Fig. 3Å), but not in the Rank-reversal condition. A comparison between conditions con rmed a more positive striatal parametric e ect of equality in the No Rankreversal than Rank-reversal condition (peak MNI coordinates: [6, 14, -5], voxel-wise p(FWE) = 0.032, t-value = 4.01, k = 45, Fig. 3B and SI Appendix, Fig. S7 for a visualization of this e ect). Note that this e ect was also con rmed in the subsequent whole-brain analysis (*SI Appendix*, Table S7). e 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 (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. is nding 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 (FWE-SVC) = 0.043, k = 6, ROI center MNI coordinates [-12, 10, -6]) was indeed associated with a signi cantly 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.

condition

No Rank-reversal B No Rank-reversal >

Rank-reversal

Α

Taken together, these ndings 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.

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 ese areas comprised dorsomedial prefrontal cortex/anterior o er. 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). us, these areas could either represent the strength of the harm aversion motive, or they could be involved in processing/resolving the con ict between concerns about inequality and harm. e latter interpretation may be in line with previous ndings that DMPFC/ACC, IFG, and MFG are often activated during cognitive control, con ict resolution, or behavioral adaptation (37, 38); and that TPJ is involved in mentalizing and perspective taking (39, 40). However, none of the neural e ects 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. is motivated us to further examine whether and how the strength of the di erent motives was represented by interactions between the di erent neural systems representing harm and equality.

We had observed weaker inequality aversion and dampened striatal sensitivity to equality in the Rank-reversal condition. ese ndings suggest that behaviorally relevant neural equality signals may not be represented invariably across di erent contexts, but may be modulated in situations where they con ict with harm signals. If this "con ict 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.

x = -3-

Parametric strength of D Parametric effect of inamine



0.6

equality

estimates of the signif cant 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 ef ects 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. Signif cant clusters are thresholded at voxelwise 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 e ects 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). e 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 quanties 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 payo s on the present trial), while simultaneously controlling for any main e ects of (i.e., simple correlations with) the striatum time course and the equality context (41).

ese 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. 4*C*, *Left*; note that this e ect was also present in control PPI analysis containing parametric inequality regressors; see SI Results). Importantly, the DMPFC region identi ed here largely overlapped with the DMPFC region involved in signaling harm to others (Fig. 4*C*, *Left*). A post-hoc comparison con rmed that this equality e ect 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. 4*C* right, Rank reversal absence vs. presence).

To assess whether the pattern of DMPFC-Striatum connectivity



Stronger DMPFC-Striatum connectivity associated with weaker neural equality signals in striatum and behavioral effects. (A) We focused the contextdependent 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 of ect 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 of ect 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 (*Ø*) (*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. •••, P < 0.001; ••, P < 0.01; ••, P < 0.05. Signif cant clusters

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 o er, again implying that harm-related activity in DMPFC and TPJ may deter more equal distributions, in particular for people who are averse to inequality.

Dif erent Motives Af ect Choice via Dif erential Patterns of Network Interactions. e patterns of results until now suggest that inequality and harm aversion are implemented by di erent 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 di erent 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. us, 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, re ecting 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

[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 o er (i.e., normalized  $\alpha$ , *tau* = 0.38, *P* < 0.001, Fig. 6 *A* and *B*, *Left*). suggests that the striatum interacts with IFG more strongly when more inequality-averse individuals choose the more unequal o er in contexts where the more equal o er 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 o er (i.e.,  $\delta$ , *tau* = 0.36, *P* < 0.001, Fig. 6 *A* and *B*, *Right*), suggesting that con icts 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 speci c functional mechanisms underlying these e ects. Note that although inequality aversion (i.e.,  $\alpha$ ) and rank reversal aversion (i.e.,  $\delta$ ) are

processing and equal choice in the No Rank-reversal condition,

but we found no such e ects 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 con icts, where striatal activity was not related to either equality processing or equal

choice. We thus de ned as ROI the striatum region involved in

equality processing and equal choice in the No Rank-reversal con-

dition (a sphere with 6-mm radius centered on peak MNI coor-

dinates 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. is revealed that the connectivity

strength between striatum and right IFG (peak MNI coordinates:



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 TPJvas enhanced when individuals chose the more unequal of er vs. more equal of er. (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 *a*) chose the more unequal of er, whereas activity in putamen was enhanced when more harm-averse (i.e., higher *b*) individuals chose the more unequal of er (*Right* panel). For visualization, neural estimates of the signif cant clusters were extracted, and scatter plots show the correlation patterns (*Bottom* panel). Signif cant clusters were thresholded at voxel-wise *P* < 0.001 uncorrected and cluster-wise FWE corrected *P* < 0.05.

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

We did not observe striatal connectivity speci cally 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 ndings emphasize that distinct neural pathways link di erent 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-o s 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 di erent frontal regions for more unequal choice was related to individuals' inequality aversion and rank reversal aversion in the Rank-reversal condition. is also implies that rank reversal aversion may interact with equality-related motives via striatal-prefrontal interactions during choices of (un)equal o ers.

## 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). erefore, the question of how

to achieve distributive justice has become an intensively studied issue among researchers in many elds, including economics, politics, philosophy, and psychology. Although in uential theories claim that fairness norms take precedence over other concerns (e.g., e ciency) 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 a ect wealth distribution, and most prevailing econometric models cannot explain why individuals can prefer greater inequality when di erent motives are in con ict (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 con icts between di erent 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 con ict with harm avoidance and rank preserving motives.

Our study extends economic theories of social preferences by highlighting the trade-o 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). e 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 con icted 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 st-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 di erent 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 o ers (47). us, 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 rst-party preferences thus requires further empirical study. In general, our ndings extend in uential theories of fairness norms (25, 26) which mainly focused on e ects of inequality aversion on distribution behaviors and emphasize the importance of considering other motives (i.e., harm aversion and



Neural networks linking dif erent motives to redistribution decisions. (A) In the Rank-reversal condition, Striatum-IFG connectivity strength was enhanced when more inequality-averse (higher *a*) individuals chose more unequal of ers vs. more equal of ers (*Left* panel), and Striatum-SFG connectivity strength was enhanced when more rank reversal-averse (i.e., higher *δ*) individuals chose more unequal of ers vs. more equal of ers (*Right* panel). Neural estimates of the signif cant clusters were extracted, and scatter plots show the correlation patterns (*B*). Signif cant 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 conicts between these di erent 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 bene t 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). erefore, 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 tness advantage by satisfying individuals' psychological need for order (50) and enhancing intragroup cooperation and productivity (51).

erefore, it is not surprising that people prefer to preserve rather than reverse preexisting hierarchy (8, 21). In line with these ndings, our results suggest that the reversal of initial rankings also contributes to the disutility of equality when rank preserving and equality seeking are in con ict. Together, we demonstrate that in contrast to inequality aversion, harm aversion and rank reversal aversion function as two di erent third-party prosocial preferences to deter more equal wealth redistribution.

Our neural results rst clari ed how equality-related information is represented. GLM results support the hypothesis that individuals are sensitive to equality signals in the absence of any con ict but will be less sensitive to equality and base their decisions more heavily on other motives when they con ict 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 speci c aspects of the distributions behavior engage the striatum and trigger the corresponding behavior—does it signal equality or other potentially rewarding aspects, such as e ciency or the other's outcomes? While stronger activity in putamen was related to higher e ciency (i.e., greater overall pro ts) (5), e ciency cannot account for the pattern of results in the current study since neither of the two alternative o ers changed the overall pro ts of the distributions. An alternative explanation is that striatum activity re ects 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 di erence between the two alternative o ers 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 nding that greater sensitivity to equality in putamen was related to higher probability of more equal choice, our results suggest that striatum not only re ects 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 con icts 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. ese

ndings help to clarify the neurocognitive mechanisms of the weighing processes of di erent 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 con ict monitoring, con ict resolution, and action selection in a variety of cognitive tasks (37, 38), which may support the resolution of con ict between di erent 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

ndings, 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 di cult 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 in uences on striatum are indeed causally involved in guiding redistribution behaviors under circumstances with con icts between multiple motives.

Our results also provide crucial evidence for frontostriatal circuitry in redistribution decisions. e 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 di erent sources for moral and prosocial decision-making (31, 59). Our ndings further re ne previous accounts of frontostriatal circuitry in moral decision-making by clarifying that different prosocial motives modulate redistribution decisions through di erential frontostriatal circuitries. Nevertheless, the speci c functional contributions (i.e., inhibitory or modulatory) of these interactions between the striatal and frontal regions still need to be clari ed in future studies.

Another critical contribution of our study is to clarify what neural processes underlie the modulations of di erent 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 di erent motives and which may di er between people with di erent motive strengths. Activity in both DMPFC and TPJ was stronger when more inequality-averse individuals chose the more unequal o er, and activity in putamen was stronger when more harm-averse individuals chose the more unequal o er. 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 di erentiated 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 speci cally involved in processing of context-dependent social information (64, 65). Although our ndings 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 speci c decisions. ese ndings 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 con icts between di erent motives.

In general, our ndings may have economic, political, and social e endowment e ect has been introduced for implications (66). 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 e ect (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 ndings 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 e ects of di erent motives are scienti cally 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 di erent income groups who pursue di erent motives on the other hand.

To conclude, the current study provides a neurocomputational account of the trade-o between multiple prosocial motives underlying resource distribution. Our ndings 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 o ers neural explanations for how di erent 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**

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° 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 Helsinski and was approved by the Ethics Committee of the Department of Psychology, Peking University.

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. 1 ). 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. 1). 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

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

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  $\times$  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, (Oxford University Press, 2005).
- J. Baron, Blind justice: Fairness to groups and the do-no-harm principle. 71–83 (1995).
- 3. J. Offer, R. Pinker, Eds.,
   (Bristol University Press ed. 1, 2017)

   10.2307/j.ctt22p7jvf.
   (Bristol University Press ed. 1, 2017)

8,

- 4. J. Rawls, (Harvard University Press, 1999), 10.2307/j.ctvkjb25m.
- M. Hsu, C. Anen, S. R. Quartz, The right and the good: Distributive justice and neural encoding of equity and efficiency. **320**, 1092-1096 (2008).
   F. Tircomi, A. Rannel, C. F. Camerer, J. P. O'Doherty, Neural evidence for inequality-averse social.
- E. Tricomi, A. Rangel, C. F. Camerer, J. P. O'Doherty, Neural evidence for inequality-averse social preferences. 463, 1089–1091 (2010).
- 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.
   9, 24 (2015).
- W. Xie, B. Ho, S. Meier, X. Zhou, Rank reversal aversion inhibits redistribution across societies. 1.0142 (2017).
- 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. 424–431 (2013).
- J. J. Jordan, M. Hoffman, P. Bloom, D. G. Rand, Third-party punishment as a costly signal of trustworthiness. 530, 473–476 (2016).
- 11. E. Lo Gerfo , The role of ventromedial prefrontal cortex and temporo-parietal junction in thirdparty punishment behavior. **200**, 501–510 (2019).
- 12. B. R. House , Social norms and cultural diversity in the development of third-party punishment. 287, 20192794 (2020).
- 13. E. Xiao, D. Houser, Emotion expression in human punishment behavior. **102**, 7398-401 (2005).
- R. Yu, A. J. Calder, D. Mobbs, Overlapping and distinct representations of advantageous and disadvantageous inequality. 35, 3290–3301 (2014).
- M. Iosifidi, N. Mylonidis, Relative effective taxation and income inequality: Evidence from OECD countries. 27, 57-76 (2017).
- M. Dimick, D. Rueda, D. Stegmueller, The Altruistic rich? Inequality and other-regarding preferences for redistribution in the US. 11, 385–439 (2016).
- 17. A. Argentiero, S. Casal, L. Mittone, A. Morreale, Tax evasion and inequality: Some theoretical and empirical insights. **22**, 309-320 (2021).
- 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. 111, 17320-17325 (2014).
- M. C. Leliveld, E. van Dijk, I. van Beest, Initial ownership in bargaining: Introducing the giving, splitting, and taking ultimatum bargaining game. 34, 1214-1225 (2008).
- 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?
   7, e39627 (2012).
- R. Fernandez, D. Rodrik, Resistance to reform: Status quo bias in the presence of individual-specific uncertainty. 81, 1146–155 (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. **102**, 98-115 (2012). 23. D. J. F. De Quervain , The neural basis of altruistic punishment. **305**, 1254-1258 (2004).
- A. Strobel , Beyond revenge: Neural and genetic bases of altruistic punishment. 54, 671-680 (2011).
- G. Charness, M. Rabin, Understanding social preferences with simple tests. 817-869 (2002).
- E. Fehr, K. Schmidt, A theory of fairness, competition and cooperation. (1999).
- 27. A. W. Cappelen , Equity theory and fair inequality: A neuroeconomic study. 111, 15368-15372 (2014).
- D. J. de Quervain, U. Fischbacher, V. Treyer, M. Schellhammer, The neural basis of altruistic punishment. **305**, 1–14 (2004).
- L. Glass, L. Moody, J. Grafman, F. Krueger, Neural signatures of third-party punishment: Evidence from penetrating traumatic brain injury. 11, 253–262 (2015).
- L. J. Chang, A. Smith, M. Dufwenberg, A. G. Sanfey, Triangulating the neural, psychological, and economic bases of guilt aversion. 70, 560-572 (2011).
- M. J. Crockett, J. Z. Šiegel, Z. Kurth-Nelson, P. Dayan, R. J. Dolan, Moral transgressions corrupt neural representations of value. 20, 879–885 (2017).
- T. Nihonsugi, A. Ihara, M. Haruno, Selective increase of intention-based economic decisions by noninvasive brain stimulation to the dorsolateral prefrontal cortex. 35, 3412–3419 (2015).
- X. Gao , Distinguishing neural correlates of context-dependent advantageous- and disadvantageous-inequity aversion.
   S. Lewandowsky, S. Farrell, (2010), 10.1017/
- 34. S. Lewandowsky, S. Farrell, CB09781107415324.004.
- 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. 76, 412-427 (2013).
- J. A. Clithero, A. Rangel, Informatic parcellation of the network involved in the computation of subjective value. 9, 1289–1302 (2014).

Author af liations: <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

- 37. C. R. Oehrn adaptation.
   , Neural communication patterns underlying conflict detection, resolution, and 34, 10438-10452 (2014).
- S. F. de Kloet , Bi-directional regulation of cognitive control by distinct prefrontal cortical output neurons to thalamus and striatum. 12, (2021).
- 39.F. Van Overwalle, Social cognition and the brain: A meta-analysis.<br/>(2009).30, 829-58
- 40. C. A. Hill , A causal account of the brain network computations underlying strategic social behavior. **20**, 1142–1149 (2017), 10.1038/nn.4602.
- K. J. Friston , Psychophysiological and modulatory interactions in neuroimaging. 6, 218-229 (1997).
- A. C. Loewke, A. R. Minerva, A. B. Nelson, A. C. Kreitzer, L. A. Gunaydin, Frontostriatal projections regulate innate avoidance behavior. 41, 5487–5501 (2021).
- K. E. Pickett, R. G. Wilkinson, Income inequality and health: A causal review. 316–326 (2015).
- B. K. Payne, J. L. Brown-lannuzzi, J. W. Hannay, Economic inequality increases risk taking. 114, 4643–4648 (2017).
- 45.
   E. Fehr, U. Fischbacher, Social norms and human cooperation.
   8, 185–190 (2004).

   46.
   E. Fehr, U. Fischbacher, Third-party punishment and social norms.
   25, 63–87
- (2004).
  47. J. Hu , Social status modulates the neural response to unfairness.
  11, 1-10 (2015).
- I. Ritov, J. Baron, Reluctance to vaccinate: Omission bias and ambiguity. 263–277 (1990).
- H. Yu, J. Hu, L. Hu, X. Zhou, The voice of conscience: Neural bases of interpersonal guilt and compensation. 9, 1150–1158 (2014).
- 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. 590–609 (2014).
- N. Halevy, E. Y. Chou, A. D. Galinsky, J. K. Murnighan, When hierarchy wins: Evidence from the national basketball association. 3, 398–406 (2012).
- J. Moll , Human fronto-mesolimbic networks guide decisions about charitable donation. 103, 15623–15628 (2006).
- W. T. Harbaugh, U. Mayr, D. R. Burghart, Neural responses to taxation and voluntary giving reveal motives for charitable donations. 316, 1622–1624 (2007).
- B. Knutson, C. M. Adams, G. W. Fong, D. Hommer, Anticipation of increasing monetary reward selectively recruits nucleus accumbens. 21, 1–5 (2001).
- M. Spitzer, U. Fischbacher, B. Herrnberger, G. Grön, E. Fehr, The neural signature of social norm compliance. 56, 185–196 (2007).
- 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. 85, 418-428 (2015).
- 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. (2014).
- J. W. Buckholtz, R. Marois, The roots of modern justice: Cognitive and neural foundations of social norms and their enforcement. 15, 655–61 (2012).
- J. Hu, Y. Hu, Y. Li, X. Zhou, Computational and neurobiological substrates of cost-benefit integration in altruistic helping decision. 41, 3545–3561 (2021).
- A. Ogawa, T. Kameda, Dissociable roles of left and right temporoparietal junction in strategic competitive interaction. 14, 1037–1048 (2020).
- 61. C. K. Kovach , Anterior prefrontal cortex contributes to action selection through tracking of recent reward trends. **32**, 8434–8442 (2012).
- 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. 33, 2242–2253 (2013).
- 63. J. X. O'Reilly , Dissociable effects of surprise and model update in parietal and anterior cingulate cortex. **110**, E3660–E3669 (2013).
- S. M. Lee, G. McCarthy, Functional heterogeneity and convergence in the right temporoparietal junction. 26, 1108–1116 (2016).
- A. Konovalov, C. Hill, J. Daunizeau, C. C. Ruff, Dissecting functional contributions of the social brain to strategic behavior. 109, 3323–3337.e5 (2021).
- B. Irlenbusch, M. C. Villeval, Behavioral ethics: How psychology influenced economics and how economics might inform psychology?
   6, 87–92 (2015).
- Z. Carmon, D. Ariely, Focusing on the forgone: How value can appear so different to buyers and sellers. 27, 360–370 (2000).
- D. Kahneman, J. L. Knetsch, R. H. Thaler, Anomalies: The endowment effect, loss aversion, and status quo bias. 5, 193–206 (1991).
- 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.
   45, 947–951 (2009).
- 70. R. A. Y. Weaver, S. Frederick, A reference price theory of the endowment effect. XLIX, 696-707 (2012).
- I. Kuziemko, R. W. Buell, T. Reich, M. I. Norton, Last-place aversion: Evidence and redistributive implications. 129, 105–149 (2014).

3,