



Arbitration between insula and temporoparietal junction subserves framing-induced boosts in generosity during social discounting

Manuela Sellitto^{a,*}, Susanne Neufang^b, Adam Schweda^a, Bernd Weber^c, Tobias Kalenscher^a

^a Comparative Psychology, Institute of Experimental Psychology, Heinrich Heine University Düsseldorf, Universitätsstraße 1, 40225 Düsseldorf, Germany

^b Department of Psychiatry and Psychotherapy, Medical Faculty, Heinrich Heine University Düsseldorf, Bergische Landstraße 2, 40629 Düsseldorf, Germany

^c Institute of Experimental Epileptology and Cognition Research, University Hospital Bonn, Sigmund-Freud Straße 25, 53127 Bonn, Germany

ARTICLE INFO

Keywords:

DCM
Framing effect
Insular cortex
Social discounting
TPJ
VMPFC

ABSTRACT

Generosity toward others declines across the perceived social distance to them. Here, participants chose between selfish and costly generous options in two conditions: in the gain frame, a generous choice yielded a gain to the other; in the loss frame, it entailed preventing the loss of a previous endowment to the other. Social discounting was reduced in the loss compared to the gain frame, implying increased generosity toward strangers. Using neuroimaging tools, we found that while activity in the temporoparietal junction (TPJ) and the ventromedial prefrontal cortex (VMPFC) was associated with generosity in the gain frame, the insular cortex was selectively recruited during generous choices in the loss frame. We provide support for a network-model according to which TPJ and insula differentially subserve generosity by modulating value signals in the VMPFC in a frame-dependent fashion. These results extend our understanding of the insula role in nudging prosocial behavior in humans.

1. Introduction

Most human societies are collaborative. Collaboration offers benefits to their members that they would not be able to achieve individually. However, societies can only function efficiently if their members are willing to contribute to causes whose beneficiaries are abstract and anonymous, such as public goods, and/or to causes whose beneficiaries are socially remote, as it is often the case with wealth redistribution for social welfare, public health insurance, or state pension systems (see also Kalenscher, 2014). Most people are indeed willing to sacrifice own resources for the welfare of others (Nowak, 2006; Rilling and Sanfey, 2011), but their generosity typically declines steeply with social distance between them and the recipients of help, a phenomenon called social discounting (Jones and Rachlin, 2006; Strombach et al., 2015). Hence, while people are ready to provide costly support to friends, relatives, and acquaintances, they are less inclined to help remote strangers.

The social discount function is idiosyncratic (Kalenscher 2017; Vekaria et al., 2017; Archambault et al., 2019), but it is far from stable within and across individuals. For instance, we and others have shown that participants from individualistic or collectivistic cultures (Strombach et al., 2014) differ in their attitude towards the welfare of socially close peers; that psychosocial stress (Margittai et al., 2015) and neurohormonal stress action (Margittai et al., 2018) can increase generosity towards socially close friends and acquaintances; and that the level of prosociality towards socially close others depends on gender

and cognitive load (Soutschek et al., 2017; Strombach et al., 2016). We further showed that disrupting the temporoparietal junction (TPJ) – a brain region we recently identified as a central hub orchestrating the balance between egocentric and other-regarding preferences in social discounting (Strombach et al., 2015), and which is also associated with perspective taking (Tusche et al., 2016) and theory of mind (Saxe and Kanwisher, 2003) – by means of transcranial magnetic stimulation increases the steepness of social discounting (Soutschek et al., 2016), thus lowering the willingness to support socially remote strangers.

This body of evidence suggests that the degree by which individuals value socially close and distant others' well-being is highly malleable. However, despite its paramount theoretical and societal significance, means to increase the inclination for costly support of socially remote beneficiaries are elusive.

Here, we provide behavioral and neural evidence for a simple manipulation that aims at significantly increasing individuals' willingness to costly support socially remote others. We make use of the observation that people are more sensitive to others' losses than gains (Bardsley, 2008; Dreu, 1997; Evans and Beest, 2017; Everett et al., 2015; Li et al., 2017; Liu et al., 2020; Schweda et al., 2020; Sip et al., 2015; Smith et al., 2015; Wang et al., 2017; Xiao et al., 2016; Zheng et al., 2010), and are consequently strongly reluctant to increase their own payoff at the expense of others' welfare (Baumeister et al., 1994; Chang et al., 2011; Chang and Sanfey, 2013; Crockett et al., 2014; List, 2007). We hypothesized that participants would be more altruistic

* Corresponding author.

E-mail address: manuela.sellitto@hhu.de (M. Sellitto).

<https://doi.org/10.1016/j.neuroimage.2021.118211>.

Received 11 November 2020; Received in revised form 29 April 2021; Accepted 9 May 2021

Available online 8 June 2021.

1053-8119/© 2021 The Author(s). Published by Elsevier Inc. This is an open access article under the CC BY-NC-ND license

(<http://creativecommons.org/licenses/by-nc-nd/4.0/>)

towards others, including socially remote strangers, if a costly generous choice was framed as preventing a monetary loss to others rather than granting them a gain, even if actual economic outcomes were equivalent. In other words, we expected that the way a prosocial decision problem was framed mattered for the shape of the social discount function.

To test this hypothesis, we elicited social preferences in a standard version of the social discounting task (gain frame; [Strombach et al., 2015](#)) as well as in a loss frame variant. In each trial, participants decided to share money with other individuals on variable social distance levels. They chose between a selfish option, yielding high own-payoff and zero other-payoff, and a generous option, yielding a lower own-payoff and a non-zero other-payoff. The main difference between conditions was the way the decision problem was described: in the gain frame, a costly generous choice would yield an equivalent gain to the other player, while, in the loss frame, it would imply preventing the loss of a previous endowment to the other player. Importantly, the payoff structure was mathematically identical across frame conditions, i.e., the choice alternatives in the loss frame yielded identical own- and other-payoffs to those in the gain frame. Participants were explicitly instructed that the other persons would only be informed about the final outcome, but not about their endowment, or the loss of it; hence, they knew about the economic equivalence across frames.

We show in a series of independent studies that participants were more reluctant to make a selfish choice if this implied a loss of the endowment to the other, resulting in a substantially flatter social discount function in the loss than the gain frame and, hence, higher generosity towards socially remote others.

To obtain further insights into the psychological and neural mechanisms underlying this framing effect on social discounting, we measured blood oxygen level-dependent (BOLD) responses while participants performed both frame conditions of the social discounting task. We hypothesized that the psychological motives underlying generosity were frame-dependent and dissociable on the neural level. Consistent with our previous work ([Strombach et al., 2015](#)), we predicted that generosity in the gain frame was vicariously rewarding and the result of the resolution of the conflict between selfish and altruistic motives. Specifically, generosity in ([Strombach et al., 2015](#)) was associated with activity in TPJ, which suggested facilitation in overcoming the egoism bias via the modulation of value signals in the ventromedial prefrontal cortex (VMPFC), a brain structure known to represent own and vicarious reward value ([Bartra et al., 2013](#); [Mobbs et al., 2009](#)), through the integration of other-regarding utility. In line with ([Soutschek et al., 2016](#); [Strombach et al., 2015](#)), we therefore expected that generous choices in the gain frame would elicit activation of the VMPFC along with TPJ. Conversely, in the loss frame, we expected that the disinclination to maximize own-gain at the expense of other-loss was motivated by the desire to comply to social norms, such as the respect of others' property rights, or the do-no-harm principle. We therefore hypothesized increased activity in brain regions that are implicated in the social sentiments that motivate individuals to comply to social norms, such as the negative emotions experienced during social norm transgressions, e.g., guilt and shame, as well as the aversive experience of unfairness and inequality ([Montague et al., 2007](#); [Xiang et al., 2013](#)). Such social sentiments have been consistently associated with the insular cortex ([Chang et al., 2011](#); [Chang and Sanfey, 2013](#); [Civai et al., 2012](#); [Corradi-Dell'Acqua et al., 2013](#); [Gu et al., 2015](#); [Lallement et al., 2013](#); [Oldham et al., 2018](#); [Samanez-Larkin et al., 2008](#); [Spitzer et al., 2007](#); [Tomasino et al., 2013](#); [Von Siebenthal et al., 2017](#); [Wang et al., 2017](#); [Yu et al., 2014](#)). Results support our main hypothesis that frame-dependent choice motives were associated with distinct neural signatures. During generous choice in the gain frame we found the involvement of VMPFC and TPJ ([Hutcherson et al., 2015](#); [Strombach et al., 2015](#)), while we identified the insular cortex as the core component of a network associated with generous choice in the loss frame.

2. Material and methods

2.1. Participants

2.1.1. Studies 1–3

Three separate behavioral studies were carried out to test the validity of our paradigm in different settings and with different compensation procedures. For these studies we did not calculate the sample size in advance as we were not aware of any previous similar manipulation of social discounting. Study 1 was run online ($n = 61$; seven participants later excluded from the analyses due to bad fitting; 28 females; mean age = 36 years, ± 11 standard deviation) and participants were paid a fixed allowance of €8.5. Study 2 was run online ($n = 36$; 32 females; mean age = 21 years, ± 2.6) and participants, all psychology students on campus, were reimbursed for their time with a fixed amount of university credits. Study 3 ($n = 39$; eight participants later excluded from the analyses; 20 females; mean age = 26 years, ± 6.0) was run in the laboratory and participants were paid out with the same fully incentive-compatible procedure as in the fMRI study 4 (see below). All three studies were conducted according to the Declaration of Helsinki and they were approved by the local ethics review board of the Heinrich-Heine University Düsseldorf. For studies 1 and 2 we did not collect informed consent, as this was allowed by the local ethics committee for online studies, which were fully anonymized, whereas we collected written informed consent in study 3, in the laboratory.

2.1.2. Study 4

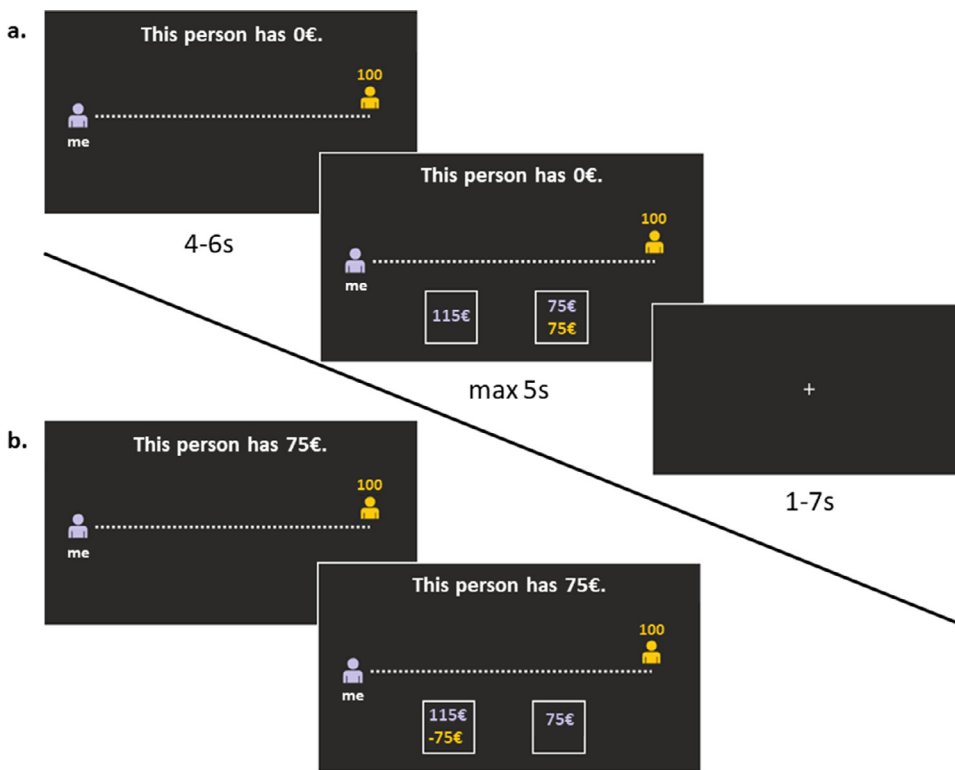
After having replicated our results across the three behavioral studies with a within-subject design (see Results), for the fMRI study we estimated, via G*Power, assuming a medium-to-large effect size, that the sample size necessary to achieve a power of 0.95 was $n = 23$. Considering frequent participants' drop out during long scanning sessions as ours, or due to excessive movement, we opted for $n = 40$. Forty healthy young volunteers were therefore recruited at the Life&Brain Research Center in Bonn for an fMRI study. All participants met MR-compatible inclusion criteria in addition to no self-reported current or history of neurological or psychiatric disorder, as well as no current use of medication affecting the central nervous system. Due to excessive head motion during measurements (>4 mm, $>4^\circ$ rotation, as computed through Artrepair Toolbox; Stanford Psychiatric Neuroimaging Laboratory, see ([Cho et al., 2013](#); [Strombach et al., 2015](#); [Wendelken et al., 2011](#)), 10 participants were excluded from all analyses. Thus, the final sample included 30 subjects (21 females; mean age = 25 years, ± 4.6 , range: 19–35 years) with high-education level (mean education = 14 years, ± 1.9 , range: 12–18 years, from high school to university master degree). Fifteen participants had a net monthly income between €0 and €499, eight between €500 and €999, five between €1000 and €1499, one between €1599 and €1999, and one larger than €2500.

All participants were fluent German speakers, right-handed, and had normal or corrected-to-normal vision. As reimbursement, they were paid €20 as participation fee, plus earnings from the social discounting task. Therefore, participants' payoff ranged from €27.5 up to €35.5 (see Social discounting task).

The study was conducted according to the Declaration of Helsinki and it was approved by the local ethics review board of the Universitätsklinikum Bonn. All volunteers gave written informed consent to participate in the study.

2.2. Social discounting task

In this task (adapted from [Strombach et al., 2015](#)), participants were first asked to imagine people from their social environment represented on a scale ranging from 1 (the person socially closest to them) to 100 (a random stranger), where a person at rank 50 was described as a person that the subject had seen several times without knowing the name. They were instructed to select six real persons located at social distances of 1,



a new trial started.

5, 10, 20, 50, and 100 (with no need of specifying the name and their social relationship for the social distances 50 and 100). Participants were encouraged to avoid thinking of people that they felt negatively toward and people they shared a bank account or household with. Each trial began with the display of the social distance level of the partner the participant was playing with. Social distance was represented with a ruler scale consisting of 101 icons. The left-most icon, highlighted in purple, depicted the participant. One of the remaining 100 other icons was highlighted in yellow, indicating the social distance of the partner. Furthermore, social distance information was additionally indicated as a number on top of the highlighted yellow icon to prevent perceptual inaccuracies in estimating social distance (cf. Fig. 1 for an example on a partner on social distance 100).

We included two experimental conditions, a *gain frame* and a *loss frame*. The *gain frame* manipulation was near-identical to the task used in (Strombach et al., 2015). Briefly, after presenting the social distance information as described above, participants were instructed that, in this trial, the experimenter gave an initial endowment of €0 to the other person (“This person has €0”; Fig. 1a). Participants were explicitly and repeatedly instructed that the other person was not aware of her zero endowment, she would only be informed of the final payoff after implementing the participant’s choice. Then, two monetary options appeared, a selfish and a generous option. The selfish option (on the left in Fig. 1a, in purple letters) indicated the reward magnitude for the participant, if chosen (e.g., €115 to the participant and no other-reward to other). The generous option contained a smaller own-reward to the participant (€75) and an other-reward to the other person (€75) (on the right in Fig. 1a). Own-rewards were always indicated in purple and other-rewards were always indicated in yellow. Participants indicated their choice of the selfish or the generous alternative by a left or right button press.

In the *loss frame*, participants were informed, after the social distance presentation, that the other person has received an initial endowment of €75 (“This person has €75”; Fig. 1b). As before, participants were explicitly and repeatedly instructed that the other person was not aware of her initial endowment, or the potential loss of it. On the next screen,

Fig. 1. Social discounting task (fMRI study 4). a. Trial example of the *gain frame*. b. Trial example of the *loss frame*. Each trial started with the presentation of a ruler-based representation of social distance to the other-person, with a left-most purple icon representing the participant and a yellow icon indicating the social distance of the other-person in the current trial (100 in this example). Additionally, participants received information as to the endowment of the other person, i.e., “This person has €0” for the gain frame (a), or “This person has €75” for the loss frame (b) (4–6 s). Afterwards, the two choice options appeared. The selfish alternative was displayed in purple fonts, indicating the own-reward magnitude to the participant (here: €115). Selfish choices implied a null gain for the partner in the gain frame (a) or the loss of the initial €75-endowment (in yellow) for the other person in the loss frame (b). The generous alternative was displayed in yellow fonts, and always yielded an equal €75 own-reward and €75 other-reward split in the gain frame (a), or a €75 own-reward gain and the possibility to keep the €75 other-endowment in the loss frame (b). As soon as the two choice options appeared, participants had 5 s to choose one of the two alternatives. After a choice was made, or after the 5 s had passed, a blank screen with a fixation cross appeared (1–7 s), and then

a selfish (on the left in Fig. 1b) and a generous alternative (on the right in Fig. 1b) appeared. When choosing the selfish alternative, the participant received the own-reward amount indicated in purple (here, €115), and the other person lost her initial endowment, as indicated in yellow (–€75), thus leaving her empty-handed. When choosing the generous alternative, the participant received a smaller own-reward indicated in purple (€75), implying that the other person would keep her endowment.

In addition to the framing (gain frame, loss frame) and the social distance levels of the other (1, 5, 10, 20, 50, 100), in each condition, we manipulated the magnitude of the own-reward across trials: we used nine selfish reward amounts per frame condition, ranging from €75 to €155 in steps of €10. The generous alternative’s payoff was invariant, always yielding €75 own-reward and €75 other-reward in all conditions and trials.

Thus, in the gain frame condition, the other person always had a €0 endowment, the selfish alternative always yielded a variable own-reward and no reward for the other, and the generous alternative invariably yielded an equal €75/€75 split between participant and other person. In the loss frame condition, the other-endowment was always €75, the selfish alternative yielded a variable own-reward accompanied by the loss of the €75 endowment to the other, and the generous alternative always yielded €75 own-reward and had no financial consequences for the other, i.e., she could keep her initial endowment of €75.

To summarize the logic of the task, both frames were mathematically equivalent, i.e., they yielded identical final payoff states to the participant and the other person (in the example in Fig. 1: both frames yield an own-reward gain of €115 to the participant and €0 gain to the other person after a selfish choice, or €75 own-reward and €75 other-reward after a generous choice). The only difference between conditions was that a €0 other-reward outcome was framed as a loss of the initial endowment in the *loss frame* vs. a null-gain in the *gain frame*, and a €75 other-reward was framed as keep-endowment in the *loss frame* vs. a €75 gain in the *gain frame*.

The order of frame conditions, selfish-reward presentations, as well as the left or right screen- position of the selfish and generous alternative

were randomized and counterbalanced across trials. The task of studies 1–3 had a total of 108 self-paced trials (54 trials per each frame). The task of study 4 had a total of 216 trials as each trial type was repeated twice to allow for full left/right position counterbalancing. For events duration of study 4, please refer to Fig. 1.

2.2.1. Incentivization procedure

In studies 3 and 4 the social discounting task was fully incentive-compatible. At the end of the session, one of the participant's choices was randomly drawn and 10% of the own-reward amount was paid out, as well as, in case of a generous choice, 10% of the other-reward amount was paid out to the other person in that trial, either via cheque, for the other-persons indicated by the participants at social distance 1, 5, 10, 20, or in cash to a random person on site in the case of other-persons at social distance 50 and 100. Note that the recipients of other-reward were only notified in case of a positive payoff, but not in case of a zero payoff or in case a trial was randomly chosen that did not consider them; in addition, they were not informed beforehand about this experiment, and, thus, had no prior outcome expectations. Hence, our incentivization procedure made it logically impossible for the other persons to know about their endowment, or the loss of it.

2.3. General procedure

2.3.1. Studies 1–3

All participants performed a social discounting task and, at the end, they completed a questionnaire assessing social desirability (see Supplementary material and methods). Although the social discounting task of studies 1 and 2 was not incentivized, participants were strongly encouraged to think as if they were making decisions for real. In studies 1 and 2, participants were instructed about the social discounting task, and then, after answering comprehension questions, they assigned other persons (i.e., name and personal relationship with them) from their social environment to the social distances 1, 5, 10, and 20, completed the task (see Social discounting task), and finally filled out a questionnaire (see Supplementary material and methods) through Unipark online survey software (Unipark questback). Participants were provided with a web link to do so, after being recruited via flyers and advertisements on social platforms. Monetary payment for study 1 was implemented via Clickworker (GmbH), whereas university credits reimbursement was carried out on campus for study 2. In study 3, after recruitment, participants were invited to the laboratory, they were instructed on the social discounting task along with the comprehension questions. They then completed the task, implemented in Matlab R2016a (MathWorks) and Cogent toolbox (www.vislab.ucl.ac.uk), and filled out a questionnaire on a laptop. Finally, they were reimbursed for participation contingent on their choices, identical to the incentivization procedure in study 4.

2.3.2. Study 4

Upon arrival, participants received instructions about the social discounting task and then, after applying comprehension questions to check for full understanding of the task, they assigned other persons (i.e., name and personal relationship with them) from their social environment to the social distances 1, 5, 10, and 20 via paper and pencil. Afterwards, participants performed a few sample trials to familiarize with the task structure and they were subsequently cleared for the scanning session. At the end of the scanning session, they answered control questions concerning the social discounting task, and filled out a demographic questionnaire as well as questionnaires assessing social desirability and empathy (see Supplementary material and methods). Finally, participants were debriefed and received their monetary allowance.

2.4. Studies 1–4

2.4.1. Behavioral data analysis

Hyperbolic model: Similar to previous studies (Jones and Rachlin, 2006; Margittai et al., 2015, 2018; Soutschek et al., 2016;

Strombach et al., 2015), we approximated the participants' decay in generosity across social distance with a hyperbolic function:

$$v = \frac{V}{(1 + k * SD)} \quad (1)$$

where v represents the discounted value of generosity, SD represents social distance, k represents the degree of discounting, and V is the intercept at social distance 0, thus the origin of the social discount function (Jones and Rachlin, 2006; Margittai et al., 2015; Soutschek et al., 2016; Strombach et al., 2015). While V can be considered an indicator of generosity towards socially close others (Margittai et al., 2015, 2018; Strang et al., 2017), k describes the discount rate, i.e., the steepness by which the social discount function decays across social distance. We estimated k and V for each participant separately, depending on her individual choice pattern.

To estimate V , we titrated the selfish amount to determine, at each social distance, the point at which the subject was indifferent between the selfish and generous options (i.e., indifference point; see Supplementary results). Logistic regression, implemented in Matlab R2016a (MathWorks), was used to determine the indifference points where the likelihood of choosing the selfish and the generous options was 50% (Soutschek et al., 2016; Strombach et al., 2015). Across the four studies, V ranged between 10 and 99, 95% CI [76, 82] for the gain frame, and between 10 and 98, 95% CI [73, 79] for the loss frame. Across the four studies, the median R^2 of the estimated V parameters equalled 0.99, range = 1 for the gain frame, and 0.93, range = 1 for the loss frame.

To fit Eq. (1) and estimate k , we modeled trial-by-trial choices via a softmax function to compute the probability P of choosing the selected option o_i over the other option o_{ii} on a given trial:

$$P_{o_i} = \frac{1}{1 + \exp(-1 * m * (v_{o_i} - v_{o_{ii}}))} \quad (2)$$

given the subjective values v (based on the current selfish amount and social distance) of the current available options o_1 (v_{o_1}) and o_2 (v_{o_2}) as in Eq. (1). The nuisance parameter m reflects the stochasticity of individual performance. The larger the m , the less noisy the choice pattern. Individual discount rates were defined by the respective k value that yielded the best prediction of the observed choice probabilities by applying maximum-likelihood estimation using nonlinear optimization procedures (fminsearch function), implemented in Matlab R2016a (MathWorks). To this end, we minimized the log-likelihood of the choice probabilities to obtain the best-fitting k and m parameter estimates, by summing across trials, given a specific set of model parameters k and m , the logarithm of $P(o_i)$. Across the four studies, k ranged between $5E-12$ and 0.69, 95% CI [0.03, 0.06] for the gain frame, and between $5E-13$ and 0.48, 95% CI [0.01, 0.03] for the loss frame. Across the four studies, the median log-likelihood of estimated k parameters equalled -21 , range = 60 for the gain frame, and -20 , range = 58 for the loss frame.

We additionally performed *parameter recovery simulation* to check that the fitting procedure had generated meaningful parameter values. Based on the procedures described in Wilson and Collins (2019), we used obtained individual discount parameters k and their respective noise parameter m to create synthetic participants, computed 10 simulations of responses of these synthetic participants, fitted the simulated data with our model (see Eqs. (1) and (2)), and compared the mean values of the obtained recovered parameters from the simulations against the inputted parameters of all four studies collapsed (see also Studer et al., 2019). Participants with null discounting (no variance in choice) were excluded from the sample as model parameters could not be estimated for them. Four simulations (out of a total of 2440) were excluded as they led to a k value of the order of $E+13$. The parameter recovery simulations showed adequate recovery of the k parameters, with a Pearson correlation between inputted and mean recovered k parameter estimates of $r = 0.95$, $p < 0.001$, Cohen's $d = 6.08$ for the gain frame and of $r = 0.96$, $p < 0.001$, Cohen's $d = 6.86$ for the loss frame (see Fig. S1a,b). Moreover,

we correlated the difference in k values between gain and loss frame for each real participant with the difference in k values between gain and loss frame for their respective average simulation. This correlation was $r = 0.89$, $p < 0.001$, Cohen's $d = 3.90$, a result that is corroborated by comparing simulated k parameters for the gain and the loss frame: once again, k values for the loss frame were significantly smaller than k values for the gain frame (median $k_{\text{gain}} = 0.02$, range $k_{\text{gain}} = 0.86$ vs. median $k_{\text{loss}} = 0.008$, range $k_{\text{loss}} = 0.52$; Wilcoxon matched-pair test: $Z = 7.10$, $p < 0.001$, $r = 0.70$), matching our results. Additionally, the recovery of the noise parameter m led to a correlation of $r = 0.66$, $p < 0.001$, Cohen's $d = 1.76$ for the gain frame and of $r = 0.41$, $p < 0.001$, Cohen's $d = 0.90$ for the loss frame. Note that the recovery was compromised by low-noise participants (with high m values) because their noise parameters were likely overestimated. Excluding those participants (9 out of 140 for the gain frame and 7 out of 105 for the loss frame) improved the recovery of the noise parameter m to $r = 0.95$, $p < 0.001$ for both the gain and the loss frame, without altering the recovery of the discount parameter k ($r = 0.95$, $p < 0.001$ for the gain frame and of $r = 0.96$, $p < 0.001$ for the loss frame) (see Fig. S1c,d; please note that the excluded participants are not reported in these two figures as they were masking the representation of the rest of the data). Across the four studies collapsed, the median log-likelihood of recovered k parameters was -20 , range = 56 for the gain frame, and -19 , range = 57 for the loss frame. Thus, in summary, our simulations showed reliable and adequate recovery of the hyperbolic discount parameter k and the noise parameter m .

Both estimated variables V and k were analyzed using non-parametric statistics as they were, in most of the cases across the four studies, not normally distributed even after log-transformation (Kolmogorov-Smirnov, all $d_s > 0.20$, all $p_s < 0.05$). When participants did not discount at all (i.e., they always chose the generous or the selfish option), k was set to 0 and V was set to 80 (i.e., maximum reward amount foregone = maximum selfish amount 155 – generous amount 75) for all generous choices or to 0 for all selfish choices. All behavioral analyses were run in Statistica 12 (StatSoft). For additional analyses (i.e., indifference points, area under the curve, reaction times, and questionnaires) please see the Supplementary information.

2.5. Study 4

2.5.1. fMRI procedures

Magnetic resonance images were collected on a 3T whole-body scanner (Magnetom Trio, Siemens Medical Systems, Erlangen, Germany) with an 8-channel head coil. For functional imaging, gradient-echo echo-planar images (EPI) were acquired at TR = 2500 ms (TE = 30 ms; number of slices = 37; slice thickness = 3 mm; distance factor = 10%; FoV = 192 mm × 192 mm; matrix size = 64 × 64; flip angle = 90°). Slices (voxel size = 2 × 2 × 3 mm) were sampled in transversal orientation covering all of the brain, including the midbrain. The scanning session started with a brief localizer acquisition. Afterwards, functional data were acquired in 3 separate runs of ~421 volumes each, to allow for brief resting periods in between. In order to get information for B_0 distortion correction of the acquired EPI images, a gradient echo field map sequence (TR = 392 ms; TE 1 = 4.92 ms; TE 2 = 7.38 ms; number of slices = 37; voxel size 3 × 3 × 3 mm) was recorded before each functional run. Structural images were collected at the end (~5 min), using a T1-weighted sequence (rapid acquisition gradient echo sequence; 208 sagittal images; voxel size = 0.8 × 0.8 × 0.8 mm; 0.8 mm slice thickness).

Head movements were minimized by the use of foam pads and scanner noise was reduced with earplugs. When necessary, vision was corrected-to-normal via fMRI compatible goggles. The social discounting task was programmed via an in-house software and presented via a mirror that projected a screen lying behind the participant, who made their choices via a left and a right button boxes.

2.5.2. fMRI preprocessing

Imaging data were preprocessed and analyzed with Statistical Parametric Mapping (SPM12, Wellcome Trust Centre for Neuroimaging, University College London, UK) implemented in Matlab R2016a (MathWorks). After checking raw data quality for each participant using the SPM Check Reg function (Stanford Psychiatric Neuroimaging Laboratory), all images were preprocessed by reorienting them according to the EPI SPM template and coregistered to the fieldmap via FieldMap toolbox. After the functional images were realigned and unwarped to the middle volume and all volumes for participants' motion correction by using phase correction, ArtRepair toolbox (Mazaika et al., 2009) was run in order to identify bad volumes. Bad volumes of participants included in the final sample were not repaired. However, we modeled these bad volumes as regressors of no-interest in the statistical analyses (see fMRI analyses). Finally, functional and structural images were coregistered and the images were spatially normalized based on segmentation of the anatomical image with resampling to 2 × 2 × 2 mm, and spatially smoothed using a 6 mm FWHM Gaussian kernel. High-pass temporal filtering (using a filter width of 128 s) was also applied to the data.

2.5.3. fMRI analyses

At the first-level analysis, trial-related activity for each participant was modeled by delta functions convolved with a canonical hemodynamic response function to model the effects of interest, as well as six covariates capturing residual motion-related artifacts, and a temporal derivative for each regressor of interest to account for slice timing differences.

For each participant, relevant contrasts were computed for each general linear model (GLM) (see below for details) and entered into second-level random effect analysis. The following variables were considered in the analyses: the loss frame condition; the gain frame condition; generous choices; selfish choices. Comparisons were run via one-way Analyses of Variance (ANOVAs), within subject, and via one-sample t-tests, where appropriate.

GLM1 searched for differences in BOLD activations between frame conditions during generous choices, where the onset of a generous choice was defined as the participant's button press to choose the generous option after the monetary options had appeared on the screen (see Fig. 1). It included an unmodulated regressor of all generous choices made in the loss frame condition and an unmodulated regressor of all generous choices made in the gain frame condition. Additionally, the selfish amount magnitude (see Social discounting task) was included as trial-by-trial parametric modulator of all main regressors, separately. In the main manuscript, we additionally considered the reward foregone as a parametric trial-by-trial regressor. Note that the reward foregone is a linear transformation of, and thus collinear with, the selfish reward magnitude; neural activations identified by this parametric regressor might therefore reflect selfish amount or reward foregone (see main text). Reaction times (RTs) were used as duration to account for differences between gain and loss frames (see Supplementary behavioral results). Additionally, missed trials were included as regressors of no-interest and modeled with duration = 5 s, i.e., the maximum time allowed to respond.

Please note that at the level of choice, where the choice onset was defined as the participant's button press after the release of the monetary options at each trial, a full model including separate regressors for both frames (gain and loss) and both types of choice (selfish and generous), as well as the trial-by-trial selfish amount as parametric regressor, was possible only for sixteen participants. This was due to participants who had to be excluded because they never, or only very rarely (not in all experimental runs) chose the selfish alternative in the loss frame. To address potential statistical power concerns associated with small sample size and to attend to potential selective sampling biases, in addition to generous choice being our main focus, we ran instead the above-mentioned model. Nevertheless, results of this full model (**GLMS1**), as

well as of a model including only selfish choices (*GLMS2*), are reported in the Supplementary information for completeness.

GLM2 tested for the effect of frame condition, and therefore included an unmodulated regressor of the onsets of the loss frame condition and an unmodulated regressor of the onsets of the gain frame condition. The frame onset was defined as the trial start (see Fig. 1). The social distance was included as trial-by-trial parametric modulator of the frame onsets, separately for the gain and the loss frames. A stick function was used as duration.

To address potential statistical concerns relative to having modelled, separately, the two main events of our task (i.e. frame onset and participants' response onset), we ran an additional analysis (*GLMS3*) including all main regressors and all parametric regressors of both *GLM1* and *GLM2*. We replicated results of both models, indicating that the frame regressors do not compete for variance with the choice regressors. Thus, the anterior insula activation has been correctly attributed to generous choice in the loss frame, making our original interpretation plausible (see Supplementary information for analysis details and results).

All whole-brain level results as well as ROI-based (see below) results were initially thresholded at $p < 0.001$ (uncorrected), minimum cluster size = 5 voxels, and then corrected at the cluster level for multiple comparisons ($p < 0.05$, family-wise error rate [FWE]). Bad volume onsets (as measured via ArtRepair toolbox; (Mazaika et al., 2009)), modeled with a stick function, were included as regressors of no-interest in all the above GLMs.

We additionally conducted, where relevant (see main text) ROI analyses for VMPFC, TPJ, and insular cortex by using anatomical bilateral masks from the Harvard-Oxford Atlas (Jenkinson et al., 2012), and the SPM Anatomical Automatic Labeling Toolbox, version 3 (Rolls et al., 2020), via SPM12 in Matlab R2016a. The probability maps of the SPM Anatomy Toolbox, version 3.0 (Eickhoff et al., 2007), and Neurosynth (<http://neurosynth.org>) were used for double checking region localization throughout GLMs.

Dynamic causal modeling (DCM): We used DCM analysis as implemented in SPM12. This analysis focused on the interplay between insula and VMPFC and between TPJ and VMPFC, addressing both (i) regions endogenous connectivity and (ii) condition specific modulation of the regions (driving inputs) and their connections (modulatory inputs). We therefore constructed a hierarchical model with regressors defining both frame conditions activations against the total baseline activation. Thus, we entered in the DCM: a regressor of no-interest for baseline connectivity ('all trials', used to correct for global activation) including onsets of the screen presenting the framing information and the social distance, and the onsets of the screen presenting the monetary options, at all trials; a regressor ('all loss trials') including onsets of the screen presenting the framing information and the social distance, and the onsets of the screen presenting the monetary options, for the loss frame trials; a regressor ('all gain trials') including onsets of the screen presenting the framing information and the social distance, and the onsets of the screen presenting the monetary options, for the gain frame trials.

Subject-specific coordinates were guided by ROI-based group activation maxima in the three network regions from the univariate, group-level results (see Results section). Volumes of interest (VOI) spheres, with a radius of 6 mm, were built around the posterior insula (*GLM2*, [34, -16, 8]), rTPJ (*GLM1*, [50, -66, 36]), and VMPFC (*GLM2*, [2, 50, -8]). Note that we focused our DCM analysis on the posterior insula cluster only, as we were interested in a baseline frame activation; including the anterior insula cluster, specific for generous choice within the loss frame (see Results section), might have biased the results in favor of our hypotheses. Also note that we found TPJ in the choice-related analysis (*GLM1*) only: in our opinion, it was still preferable to opt for this experimentally driven ROI rather than including an ROI taken from the literature. Regional time series were extracted as the first eigenvariate of the three network regions for 'all trials' and mean-corrected for the effect of interest F-contrast at a liberal threshold of $p = 0.1$. This

threshold was lowered for some participants until all regions could be detected (Zeidman et al., 2019a, 2019b).

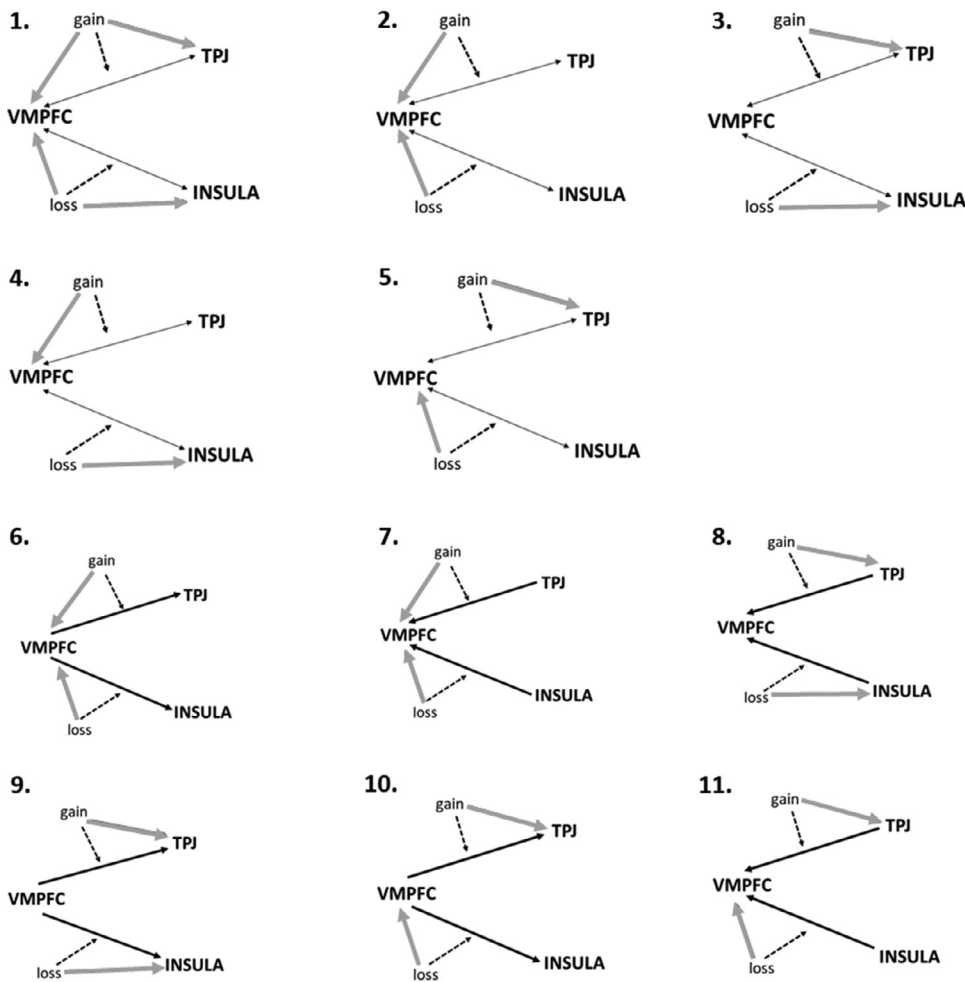
Based on our univariate results, we constructed bilinear models where the endogenous connectivity across the three regions was always assumed. We specified models with nodes reciprocally connected, where the gain and loss frame were allowed to modulate all connections (Li et al., 2015). The resulting 15 models were grouped in two families: A and B. In family A, both condition-specific driving inputs and condition-specific modulatory inputs were assumed. In family B, only condition-specific driving inputs were assumed.

Family A included eleven models (Fig. 2). In model 1 (sum of the log-evidence $SF = -4.0797E+05$, exceedance probability $xp = 0.1282$), we assumed that the gain frame condition had direct inputs on VMPFC and TPJ, and a modulatory input on their connections; the loss frame condition had direct inputs on VMPFC and insula, and a modulatory input on their connections. In model 2 ($SF = -4.0864E+05$, $xp = 0.1294$), the gain frame had a driving input on VMPFC, and a modulatory input on its connectivity with TPJ; the loss frame had a driving input on VMPFC and a modulatory input on its connectivity with the insula. In model 3 ($SF = -4.0829E+05$, $xp = 0.005$), the gain frame had a driving input on TPJ and a modulatory input on its connectivity with VMPFC; the loss frame had a driving input on the insula and a modulatory input on its connectivity with VMPFC. In model 4 ($SF = -4.0826E+05$, $xp = 0.049$), the gain frame had a driving input on VMPFC and a modulatory input on its connectivity with TPJ; the loss frame had driving input on the insula and a modulatory input on its connectivity with VMPFC. In model 5 ($SF = -4.0786E+05$, $xp = 0.6595$), the gain frame had a driving input on TPJ and a modulatory input on its connectivity with VMPFC; the loss frame had a driving input on VMPFC and a modulatory input on its connectivity with the insula. Therefore, connectivity between regions in model 1 to 5 is assumed to be bidirectional. Additionally, in model 6 ($SF = -4.0892E+05$, $xp = 0.0048$), the gain frame had a driving input on VMPFC and a modulatory input on its connectivity to TPJ; the loss frame had a driving input on VMPFC and a modulatory input on its connectivity to the insula. In model 7 ($SF = -4.0953E+05$, $xp = 0.0009$), the gain frame had a driving input on VMPFC and a modulatory input on the connectivity from TPJ to VMPFC; the loss frame had a driving input on VMPFC and a modulatory input on the connectivity from the insula to VMPFC. In model 8 ($SF = -4.0867E+05$, $xp = 0$), the gain frame had a driving input on TPJ and a modulatory input on its connectivity to VMPFC; the loss frame had a driving input on the insula and a modulatory input on its connectivity to VMPFC. In model 9 ($SF = -4.0890E+05$, $xp = 0.001$), the gain frame had a driving input on TPJ and a modulatory input on the connectivity from VMPFC to TPJ; the loss frame had a driving input on the insula and a modulatory input on the connectivity from VMPFC to the insula. In model 10 ($SF = -4.0861E+05$, $xp = 0$), the gain frame had a driving input on TPJ and a modulatory input on the connectivity from VMPFC to TPJ; the loss frame had a driving input on VMPFC and a modulatory input on its connectivity to the insula. In model 11 ($SF = -4.0851E+05$, $xp = 0.0221$), the gain frame had a driving input on TPJ and a modulatory input on its connectivity to VMPFC; the loss frame had a driving input on VMPFC and a modulatory input on the connectivity from the insula to VMPFC.

Family B included four models (Fig. 2). In model 12 ($SF = -4.0886E+05$, $xp = 0$), the gain frame had driving input on VMPFC and TPJ; the loss frame had driving inputs on VMPFC and the insula. In model 13 ($SF = -4.0936E+05$, $xp = 0$), the gain frame had a driving input on TPJ and the loss frame had a driving input on insula. In model 14 ($SF = -4.0999E+05$, $xp = 0.0001$), both frame conditions' driving inputs were on VMPFC. In model 15 ($SF = -4.0893E+05$, $xp = 0$), the gain and the loss frame had driving inputs on all three regions, insula, TPJ, and VMPC, to check whether at increased number of connections, the model fitted the data better.

All the hypothesized models were entered into Bayesian Model Selection (BMS), as implemented in SPM, to determine the best-fit family and model. The inference method used to compare the models across

Family A



Family B

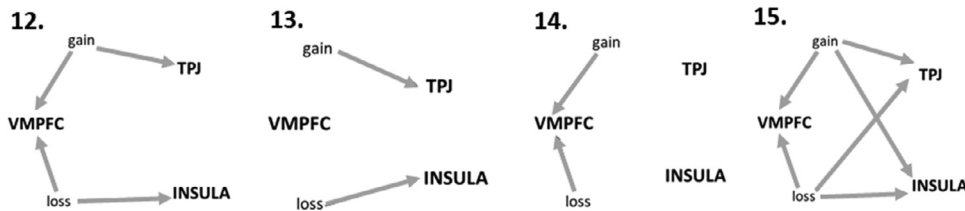


Fig. 2. DCM models. Fifteen models were hypothesized to describe the data. Gray lines represent driving inputs. Dashed black lines represent modulatory inputs. Thin black lines represent bidirectional connectivity. Thick black lines represent unidirectional connectivity. Since endogenous connectivity is always assumed between all three regions in all models, it is not represented here. *Family A*, which assumed both condition-specific driving inputs and condition-specific modulatory inputs, includes models 1 to 11. *Family B*, which assumed only condition-specific driving inputs, includes models 12 to 15.

subjects and session was random effects (2nd-level, RFX). Bayesian Model Averaging (BMA) was used for model comparison. Once the optimal model was selected, the participant-specific parameters for the two frame conditions were averaged across the three runs and entered into group analysis with one-sample and paired-sample t-tests, where appropriate. This allowed us to summarize the consistent findings from the subject-specific DCMs using classical statistics (Cho et al., 2013; Li et al., 2015; Neufang et al., 2016; Wiehler et al., 2017; Zhang et al., 2018).

Mediation analysis: This analysis was run via Hayes’s PROCESS-macro (Hayes, 2017) as implemented in the IBM Statistical Package for the Social Sciences (SPSS). The analysis aimed at testing the idea that the gain and the loss frame had an effect on generous behavior through the mediating influence of condition-specific neural activations. The frame

condition was included as binary independent variable X (dummy variable: 1 = gain; 2 = loss), the proportion of generous choices (gain frame and loss frame) was entered as dependent variable Y, and the neural activations were entered as mediators. Specifically, beta estimates for the posterior insula [34, -16, 8; GLM2], VMPFC1 [2, 50, -8; GLM2], the anterior insula [42, 4, -4; GLM1], TPJ [50, -66, 36; GLM1], and VMPFC2 [0, 54, 14; GLM1] were extracted, at the single-subject level, for both frames and included in the model, via MarsBaR region of interest toolbox for SPM12 (Brett et al., 2002). Neural activations across both frames were treated as parallel mediators (model template 4, Hayes, 2017). Partially standardized values are reported, and 95% biased-corrected CIs are adopted. Number of bootstrap samples was set to 5000. To determine the statistical power for mediation, the online tool MedPower was

used (<https://davidakenny.shinyapps.io/MedPower/>) using effects of X on mediator (M) (path a), of M on Y (path b), and the direct effect of X on Y (path c'), at $\alpha=0.05$. Total achieved power was ~ 0.60 .

3. Results

3.1. Social discounting is flatter in the loss than the gain frame

First, in a set of behavioral experiments, we established that our framing manipulation affected generosity towards socially distant others. In a within-subject design, we elicited social preferences in a standard version of the social discounting task (the gain frame; [Strombach et al., 2015](#)) as well as in a loss frame variant (see [Fig. 1](#)), interleaved in a trial-by-trial fashion. In the gain frame, participants played with other persons at variable social distance levels, and made choices between a selfish alternative, yielding higher monetary payoff to the participant and zero payoff to the other, and a generous alternative, always yielding a lower own-payoff of €75 along with a payoff of €75 to the other. In the loss frame, participants were first informed that the other person had received an initial endowment of €75. The selfish alternative yielded a variable, higher own-payoff as well as the loss of the €75 endowment to the other person, hence, resulting in a zero payoff to the other; the generous alternative yielded a fixed €75 payoff to the participant, and no further consequence to the other, thus leaving her with the initial €75 endowment. Crucially, the payoff structure was mathematically equivalent across both frame conditions, i.e., the choice alternatives in the loss frame yielded identical own- and other-payoffs to those in the gain frame. The main difference between conditions was that, in the gain frame, a generous choice would imply a gain of €75 to the other, while in the loss frame, a generous choice would imply preventing the loss of the previous €75 endowment. Importantly, participants were repeatedly instructed that the other person was unaware of her initial endowment, or the loss of it, and that she would only be informed about the final outcome of the payoff after implementing the participant's choice at the end of the experiment. Task comprehension, in particular regarding participants' understanding that the other person would only be informed about the final outcome, but not about her endowment, or loss of it, was further stressed during the explanation of the incentivization procedure as well as assessed in post-hoc structured interviews (see Material and methods). All participants understood the task well.

3.1.1. Study 1

In a first study, data collection was done online and the task was not incentive-compatible; participants ($n = 54$) were paid a fixed allowance of €8.5. In the social discounting task, participants can either make a selfish choice or a generous choice, in each frame condition. We adopted the hyperbolic discount model to describe the effect of framing on our participants' behavior as it is the best-documented model to investigate social discounting, with demonstrated better goodness-of-fits in comparison to other, e.g., exponential, models (e.g. [Jones and Rachlin, 2006](#)). Specifically, to reconstruct the individual social discount functions, separately for the two frame conditions, we fit a standard hyperbolic model (see [Eq. \(1\)](#) ([Jones and Rachlin, 2006](#); [Strombach et al., 2015](#)); Material and methods) to trial-by-trial binary choices (i.e., either selfish or generous) to estimate the parameter k , a measure of the steepness of the social discount function. Additionally, we determined, for each participant and each social distance level, and separately for the two frame conditions, the point at which the participant was indifferent between the selfish and the generous alternative using logistic regression ([Strombach et al., 2015](#)). The difference in reward magnitudes for the participant between the two alternatives at the indifference points (see Supplementary results) represented the amount of money a subject was willing to forego (i.e., reward amount foregone) to increase the wealth of another person at a given social distance, and could be construed as a social premium, that is, the price tag participants put on increasing

the wealth of the other. We took the estimated parameter V , the intercept at social distance 0, thus the origin of the social discount function ([Jones and Rachlin, 2006](#); [Margittai et al., 2015, 2018](#); [Soutschek et al., 2016](#); [Strombach et al., 2015](#)), as an indicator of generosity towards socially close others ([Soutschek et al., 2016](#); [Strombach et al., 2015](#)).

Participants' generosity dropped much less steeply in the loss compared to the gain frame (median $k_{\text{gain}} = 0.022$, range $k_{\text{gain}} = 0.18$ vs. median $k_{\text{loss}} = 4.74\text{E-}11$, range $k_{\text{loss}} = 0.10$; Wilcoxon matched pairs test: $Z = 4.97$, $p < 0.001$; $r = 0.68$; see supplementary [Fig. S2a](#)). The difference in social discount functions between frames was most pronounced at high social distance levels, indicating that participants were substantially more generous towards strangers in the loss than the gain frame. We also found a significant difference in V between frame conditions (median $V_{\text{gain}} = 87$, range $V_{\text{gain}} = 88$ vs. median $V_{\text{loss}} = 77$, range $V_{\text{loss}} = 68$; $Z = 2.78$, $p < 0.01$; $r = 0.38$) that, however, disappeared when removing all participants with zero discounting from the analysis (see Supplementary results). These data suggest that participants were strongly more generous towards socially distant others in the loss than the gain frame.

3.1.2. Study 2

In a second study, we replicated the results of our first experiment. Data collection was done online and participants ($n = 36$) were reimbursed for their time with a fixed amount of university credits. We again found that participants had flatter social discounting in the loss than the gain frame (median $k_{\text{gain}} = 0.020$, range = 0.62 vs. median $k_{\text{loss}} = 0.0005$, range = 0.10; $Z = 4.81$, $p < 0.001$; $r = 0.80$; see supplementary [Fig. S2b](#)), and we found no difference in V between frame conditions (median $V_{\text{gain}} = 81$, range $V_{\text{gain}} = 66$ vs. median $V_{\text{loss}} = 80$, range $V_{\text{loss}} = 50$; $Z = 0.38$, $p = 0.71$). Again, these results held when excluding participants with null discounting.

3.1.3. Study 3

Studies 1 and 2 were not incentive-compatible. To determine whether hypothetical versus real payoffs made a difference in the frame effect on social discounting ([Vlaev, 2012](#)), we ran a third fully incentive-compatible study in a laboratory setting ($n = 31$). Payoff was contingent on the participants' choices, and was paid out to self and other, identical to ([Strombach et al., 2015](#)) and to the fMRI study 4 (see next paragraph and Material and methods). Once again, we could replicate the frame effect on k (median $k_{\text{gain}} = 0.022$, range $k_{\text{gain}} = 0.69$ vs. median $k_{\text{loss}} = 4\text{E-}07$, range $k_{\text{loss}} = 0.48$; $Z = 3.71$, $p < 0.001$; $r = 0.67$; see supplementary [Fig. S2c](#)), and there was no difference in the V parameter between frame conditions (median $V_{\text{gain}} = 81$, range $V_{\text{gain}} = 99$ vs. median $V_{\text{loss}} = 80$, range $V_{\text{loss}} = 88$; $Z = 1.41$, $p = 0.16$). These results held when excluding participants with null discounting.

Additionally, we plot for all three studies the proportion of generous choices, averaged across participants, as a function of the selfish amount to highlight the flatter decay in generous choices in the loss frame compared with the gain frame, especially at remote social distances ([Fig. S3a,b,c](#)). Furthermore, the distributions of individual k value differences and V value differences between frames are shown in [Fig. S4a,b,c](#) for all three studies.

Moreover, social desirability, as measured via the Social Desirability Scale (SDS-17; [Stöber, 2001](#)), did not explain the frame effect on social discounting parameters (see Supplementary material).

The result of increased generosity, especially at larger social distances, in the loss compared to the gain frame in the three behavioral studies was also corroborated via a model-free measure, i.e. the area under the curve (AUC), as well as via an analysis of the indifference points (the selfish reward magnitude at which participants were indifferent between the selfish and the generous alternative at each social distance level and in each frame condition; see Supplementary analyses and results).

Overall these results suggest that, while generosity to socially close others was comparable between frame conditions, it decayed signifi-

cantly less steeply across social distance in the loss than in the gain frame, indicating that participants were considerably more generous towards socially distant others in the loss frame.

3.2. Neural mechanisms underlying the frame effect on social discounting

To obtain more substantial insights into the psychological and neural mechanisms underlying this framing effect on social discounting, in study 4 we measured BOLD responses while participants performed both frame variants of the social discounting task. The fundamental premise of our study is that the decision motives and their neural correlates differ between gain and loss frames. Specifically, we reasoned that generosity in the gain frame was mainly stimulated by other-regarding considerations. Conversely, we predicted that generous decisions in the loss frame were motivated by the desire to comply to social norms, such as the *do-no-harm* principle, or the respect of others' property rights (Sethi et al., 1996), infringements of which are associated with negative social sentiments of guilt and shame. To test this idea, we focused on two main hypotheses. We, first, expected that generosity in the gain frame recruited a network of structures, including VMPFC and TPJ (Hutcherson et al., 2015; Strombach et al., 2015), known to represent vicarious reward value and prosocial behavior. Second, we hypothesized that brain areas implicated in negative social sentiments of social norm transgressions, such as the insular cortex (e.g. Paulus et al., 2003; Chang et al., 2011; Chang and Sanfey 2013; Lallemand et al., 2013; Gu et al., 2015; Seara-Cardoso et al., 2016; Sethi and Somanathan 2016; Siebenthal et al. 2017; Wang et al., 2017; Huggins et al., 2018), would be selectively recruited during generous choices in the loss, but not the gain frame.

We first replicated, once more, the behavioral framing effect on social discounting ($n = 30$). As before, the drop in generosity across social distance was pronouncedly flatter in the loss than the gain frame (median $k_{\text{gain}} = 0.021$, range $k_{\text{gain}} = 0.16$ vs. median $k_{\text{loss}} = 0.003$, range $k_{\text{loss}} = 0.28$; Wilcoxon matched pairs test: $Z = 3.69$, $p < 0.001$; $r = 0.67$; Fig. 3a), but, again, there was no difference in the V parameter between conditions (median $V_{\text{gain}} = 80$, range $V_{\text{gain}} = 57$ vs. median $V_{\text{loss}} = 80$, range $V_{\text{loss}} = 53$; $Z = 0.88$, $p > 0.37$; the results remained identical when excluding participants with null discounting). Additionally, we plot the proportion of generous choices, averaged across participants, as a function of the selfish amount (Fig. 3b) to illustrate the flatter decay in generous choices in the loss compared to the gain frame, especially at remote social distances. Furthermore, the individual distributions of k value differences and V value differences between frames are shown in Fig. S4d, as well as we plot the choice probability as a function of the difference in value between the generous and the selfish option (Fig. S5).

Moreover, neither social desirability (SDS-17; Stöber 2001), nor perspective taking, empathic concern, personal distress, or fantasy (as measured via the Interpersonal Reactivity Index; IRI; (Davis, 1983; Paulus, 2009) explained the frame effect on social discounting parameters (see Supplementary material).

The above results were corroborated, once again, also via an analysis of the AUC, as well as via an analysis of the indifference points at each social distance (see Supplementary analyses and results).

Our first hypothesis predicted activity in brain structures known to represent vicarious reward value and prosocial behavior in the gain frame (similar to Soutschek et al., 2016; Strombach et al., 2015). Our results (GLM1; see Material and methods) indeed revealed clusters located in VMPFC (0, 54, 14, whole-brain $p_{\text{FWE-corr}} < 0.001$) as well as right TPJ (rTPJ; 50, -66, 36, whole-brain $p_{\text{FWE-corr}} < 0.035$) to be selectively activated, in addition to other prefrontal regions, when participants made generous choices in the gain frame relative to generous choices in the loss frame. ROI analyses confirmed significant clusters of activation in both VMPFC ($p_{\text{FWE-corr}} < 0.001$) and rTPJ ($p_{\text{FWE-corr}} = 0.01$). Thus, consistent with (Hutcherson et al., 2015; Strombach et al., 2015), a network comprising VMPFC and rTPJ seems to underlie the motivation for costly generosity in the gain frame (Fig. 4; see supplementary Table S1). Additionally, the selfish amount magnitude, included as trial-by-trial re-

gressor, did not parametrically modulate activity in VMPFC and rTPJ (GLM1; see Material and methods).

Our second hypothesis predicted that generosity in the loss frame was motivated by social norm compliance rather than other-regarding considerations; generosity should, consequently, go along with a different neural activation pattern in the loss than the gain frame. In a first step, we attempted to isolate frame-dependent neural correlates, independent of participants' choices. To this end, we searched for differential neural activity at trial onset, i.e., when participants learned about the social distance level of the other person and which frame was relevant in the current trial (see Fig. 1), by contrasting neural activity between the two frames (GLM2; see Material and methods). We found significant activation in the right posterior insula (34, -16, 8, whole-brain $p_{\text{FWE-corr}} = 0.007$) in the loss vs. gain frame contrast, which was accompanied by significant activations in frontal regions, including VMPFC (2, 50, -8, whole-brain $p_{\text{FWE-corr}} = 0.001$), as well as temporal regions (Fig. 5; see supplementary Table S2 for a complete list of activations). ROI analyses confirmed significant clusters of activation in the right insula ($p_{\text{FWE-corr}} = 0.03$) as well as in VMPFC ($p_{\text{FWE-corr}} = 0.02$). The opposite contrast, gain frame vs. loss frame, did not reveal any significant activation. Social distance information, included as trial-by-trial regressor, did not parametrically modulate neural activity in any of these contrasts (GLM2; see Material and methods), suggesting that the activations in insula and VMPFC reflected frame but not social distance information.

In support of this conclusion, we found that the right anterior insula (42, 4, -4, ROI analysis, $p_{\text{FWE-corr}} < 0.02$; GLM1, see Material and methods), was selectively activated during generous choices in the loss frame relative to generous choices in the gain frame (Fig. 6; see supplementary Table S1). The location within the insula mask was slightly anterior to the peak activation we found at trial onset.

Our analysis so far suggests that insula activation reflects the psychological motives underlying generous choice in the loss frame. However, other explanations of our insula finding are conceivable, too. For instance, participants made more generous choices overall in the loss than the gain frame; i.e., they forewent more own-payoff in the loss than the gain frame, and insula activation might reflect the higher level of reward foregone in the loss frame. Yet, the trial-by-trial regressor of reward amount foregone (GLM1; see Material and methods) revealed no parametric modulation of insula activity, nor of activity in any other brain region, during generous choices in either frame condition. Additionally, insula activity is unlikely to reflect the own-reward component of the generous alternative because it was fixed (always €75) and, thus, invariant across trials in both frames.

3.3. Frame-dependent modulation of VMPFC activation by rTPJ and insula

We previously provided empirical support for a network model according to which, in a task similar to our gain frame condition, TPJ would facilitate generous decision-making by modulating basic reward signals in the VMPFC, incorporating other-regarding preferences into an otherwise exclusive own-reward value representation, thus computing the vicarious value of a reward to others (Strombach et al., 2015). Here, we expand on this idea and propose that, in addition to the TPJ-VMPFC connectivity in the gain frame, frame-related information in the loss frame would activate insula, which in turn would down-regulate own-value representations in VMPFC, thus promoting generous choices by decreasing the attractiveness of own-rewards. Hence, in brief, we predicted a complex, frame-dependent pattern of connectivity between insula, TPJ, and VMPFC that reflects the different motives underlying generosity in the gain and the loss frame.

To identify the relations between those regions, we estimated their effective connectivity via DCM analysis (Friston et al., 2003). More specifically, we tested the idea that the frame information at the beginning of each trial would drive increased insula activation selectively in the loss frame, and increased TPJ activation selectively in the gain frame. Additionally, we expected increased endogenous connectivity as

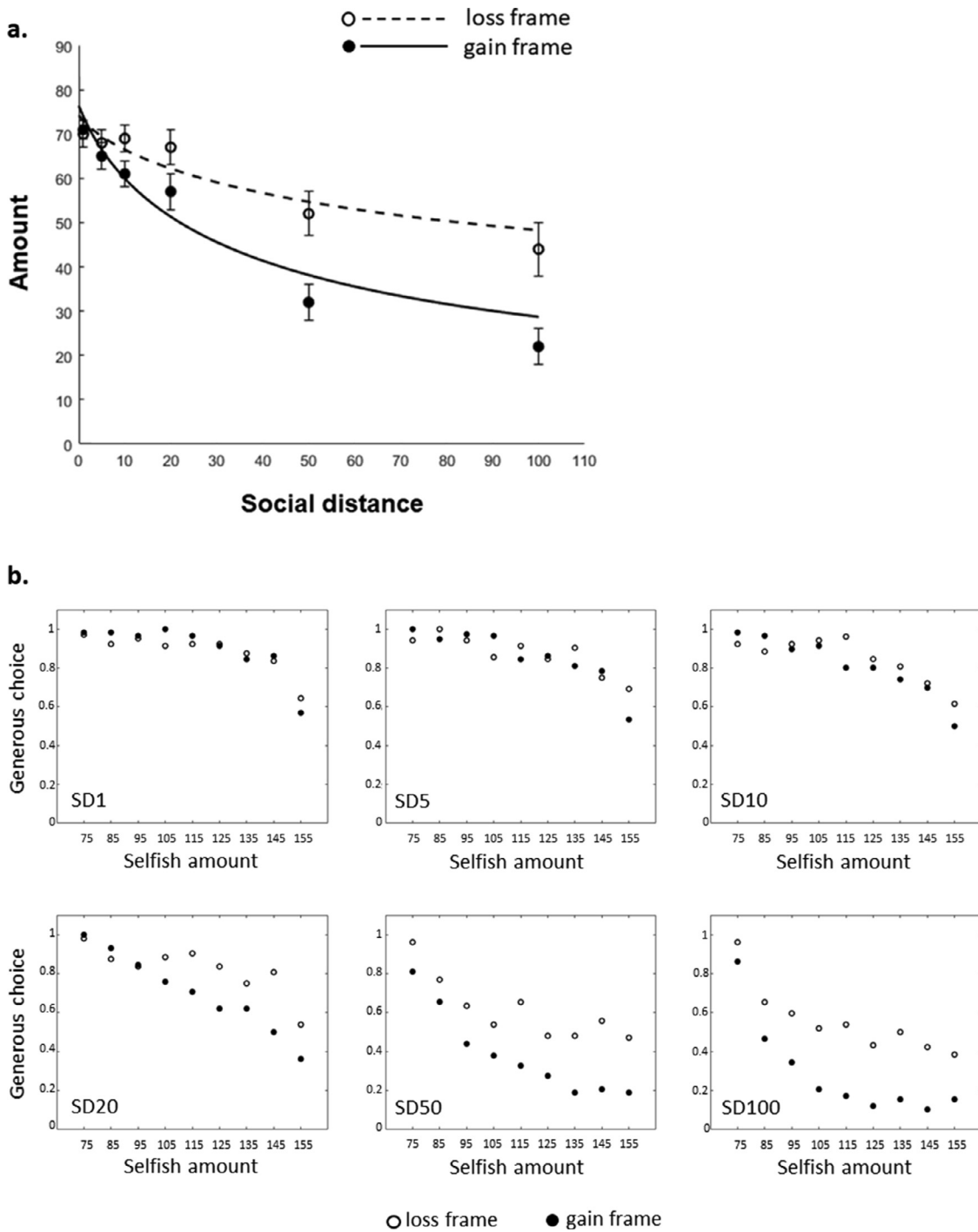


Fig. 3. Hyperbolic discount function fit and proportion of generous choices at each social distance level (fMRI study 4). (a) The change in generosity across social distances was captured by a hyperbolic discount model (see main text for details). The figure shows the mean of the participants' individual best-fitting hyperbolic functions, along with the mean amounts foregone (i.e. the social-distance-dependent reward amount that participants were willing to pay to increase the wealth of another person by €75; see main text) at each social distance (i.e. 1, 5, 10, 20, 50, 100), computed separately for the gain frame and the loss frame. The social discounting curve for the loss frame (dashed line) was significantly flatter than the social discounting curve for the gain frame (solid line). Circles represent the amounts foregone for the loss frame, dots represent the amounts foregone for the gain frame. Error bars represent the standard error of the mean. (b) Descriptive proportion of generous choices, averaged across participants, as a function of the selfish amount for the loss (circles) and the gain (dots) frame, separately for each social distance (SD).

well as condition-specific modulation between each respective region with VMPFC. Note that we focused our DCM analysis on the posterior insula cluster only, as we were interested in a baseline frame activation; including the anterior insula cluster, specific for generous choice within

the loss frame (see above), might have biased the results in favor of our hypotheses.

In total we defined 15 models (see Fig. 2), grouped into two model families: A, which assumed both condition-specific driving inputs and

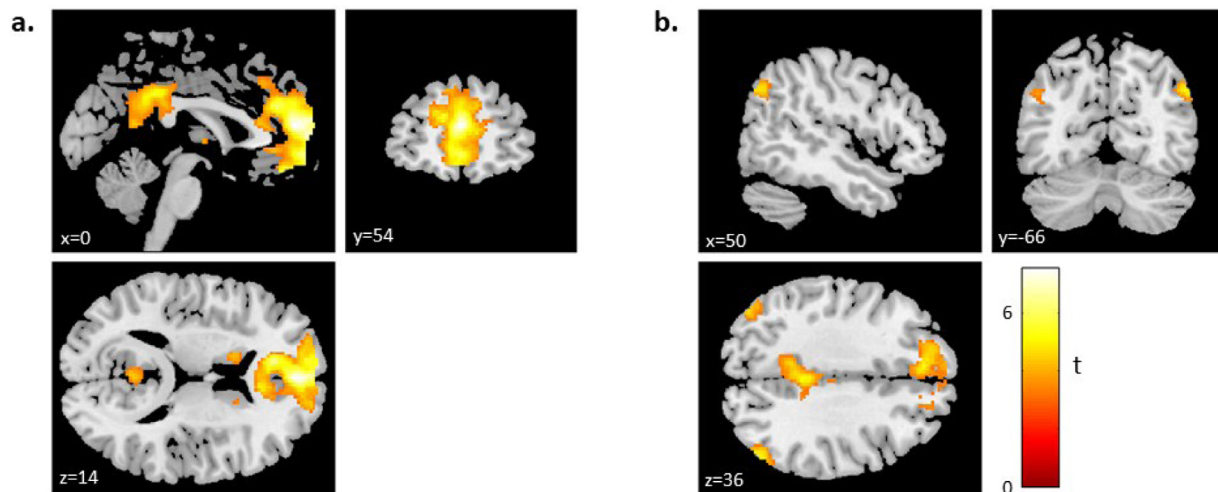


Fig. 4. Generous choices in the gain frame correlate with VMPFC and rTPJ activity. VMPFC (MNI peak [0, 54, 14]) (a) as well as right TPJ [50, -66, 36] (b) were selectively activated during [generous choice in gain frame > generous choice in loss frame] (GLM1; $p < 0.05$ FWE whole-brain corrected at the cluster level; for illustration purposes, activations are displayed at $p < 0.001$, uncorrected, minimum cluster size ≥ 5). Color bar indicates T-value.

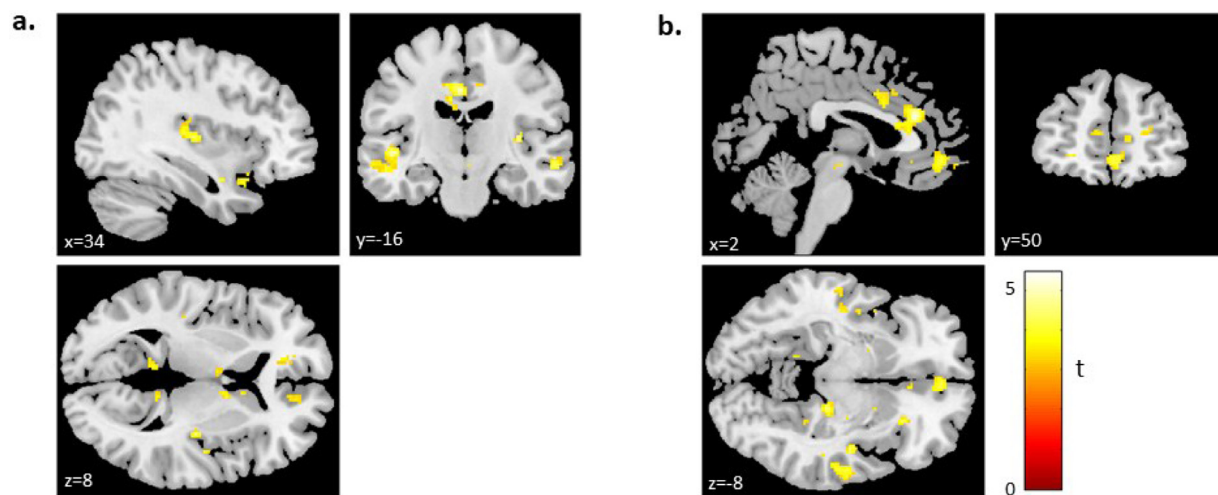


Fig. 5. The loss frame information recruits the insula and VMPFC. Insula [34, -16, 8] (a) as well as VMPFC [2, 50, -8] (b) were selectively activated during [loss frame > gain frame] onset. (GLM2; $p < 0.05$ FWE-corrected at the cluster level; for illustration purposes, activations are displayed at $p < 0.001$, uncorrected, minimum cluster size ≥ 5). Color bar indicates T-value.

condition-specific modulatory inputs; B, which assumed only condition-specific driving inputs.

Among the two model families tested, model comparison favored family A, i.e., the family of models that assumed condition-specific effects at the level of both driving input and modulatory input (family A expected posterior probability: 0.9678 vs. family B expected posterior probability: 0.0322). The winning model was model number 5 (sum of the log-evidence $SF = -4.0786E+05$, exceedance probability $x_p = 0.6595$), which assumed that the gain frame had an effect on the TPJ and its connectivity with the VMPFC, while the loss frame had an effect on the VMPFC and its connectivity with the insula (i.e., connectivity between regions is assumed to be bidirectional).

Concerning the driving inputs, we compared the average activity in TPJ in the gain frame against 0, and the average activity in VMPFC in the loss frame against 0 (we checked, beforehand, that no effect of repetition across runs was present; all $p_s > 0.18$), but none of the driving inputs was significantly different from 0 (all $p_s > 0.26$; Table 1).

Next, when addressing the modulatory inputs, the only significant difference was found in the loss frame for modulatory activity from the insula to VMPFC against the endogenous connectivity from the insula

Table 1

DCM estimated parameters of the winning model and statistics. Values are expressed as mean \pm standard error (s.e.). Statistics refer to paired t-tests between the modulatory activity and the respective endogenous connectivity, and to one-sample t-tests against 0 for driving inputs. $t = t$ -value; $p = p$ -value; subscript numbers are degrees of freedom; $*$ = $p < 0.05$; endo = endogenous connectivity; mod = modulatory connectivity; drivInp = driving input; Gain = gain frame; Loss = loss frame; TPJ = temporoparietal junction; INS = insula; VMPFC = ventromedial prefrontal cortex. Arrows indicate connectivity direction.

DCM estimated parameters	mean \pm s.e.	Statistics
endo: TPJ \rightarrow VMPFC	0.05 \pm 0.03	-
endo: VMPFC \rightarrow TPJ	0.02 \pm 0.03	-
endo: INS \rightarrow VPMFC	0.05 \pm 0.02	-
endo: VMPFC \rightarrow INS	0.01 \pm 0.02	-
mod_Gain: TPJ \rightarrow VMPFC	0.03 \pm 0.09	$t_{29} = -0.13, p = 0.90$
mod_Gain: VMPFC \rightarrow TPJ	-0.03 \pm 0.08	$t_{29} = -0.53, p = 0.60$
mod_Loss: INS \rightarrow VMPFC	-0.22 \pm 0.09	$t_{29} = -2.56, p = 0.02^*$
mod_Loss: VMPFC \rightarrow INS	-0.02 \pm 0.07	$t_{29} = -0.30, p = 0.80$
drivInp_Gain: TPJ	0.04 \pm 0.03	$t_{29} = 1.31, p = 0.20$
drivInp_Loss: VMPFC	0.00 \pm 0.03	$t_{29} = 0.10, p = 0.92$

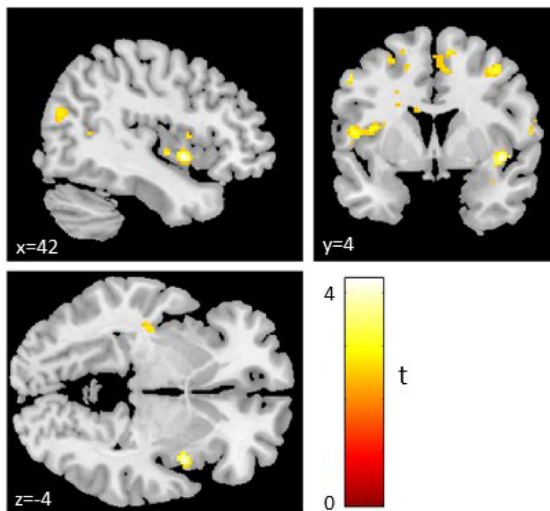


Fig. 6. Insula activation underlies generous choices in the loss frame. Insula [42, 4, -4] was selectively activated during [generous choice in loss frame > generous choice in gain frame] (GLM1; $p < 0.05$ ROI FWE-corrected at the cluster level; for illustration purposes, activations are displayed at $p < 0.01$, uncorrected, minimum cluster size ≥ 5). Color bar indicates T-value.

to VMPFC ($M_{\text{modulatory}} = -0.2158$ vs. $M_{\text{endogenous}} = 0.04672$, $p = 0.016$, Bonferroni corrected), reflecting a significant modulation of endogenous connectivity by the loss frame information (all other p s > 0.60 ; Table 1). In addition, the modulatory input was negative, hinting towards an inhibitory influence of insula on VMPFC in the loss frame (as before, there was no effect of repetition across runs in neither modulatory activity nor endogenous connectivity; all p s > 0.13).

3.4. The mediating role of the insula in the frame effect on social discounting

To provide further support to our idea that the frame effect on social discounting was brought about by condition-specific neural activity patterns, we ran a mediation analysis on the relation between frame information, generous behaviour, and neural activation in these regions. More specifically, frame was entered as independent variable X (gain and loss), the proportion of generous choices (gain frame and loss frame) was entered as dependent variable Y, and the neural activations were entered as mediators. We focused on a model where neural activations across both frames were treated as parallel mediators. Neural activations included the posterior insula and the anterior insula, TPJ, and VMPFC (both clusters in GLM1 and GLM2) (see Material and methods for details). Frame condition significantly correlated with all neural activations (all p s < 0.05) with the exception of VMPFC (GLM2) ($p = 0.052$). Additionally, while the direct effect of the frame condition on the proportion of generous choices was not significant ($p = 0.26$), the indirect effect of anterior insula on it was significant, indicating that it influenced frame-specific generosity (partially standardized $B = 0.15$, $SE = 0.09$, 95% biased-corrected CI 0.003 to 0.36) (Fig. 7).

In conclusion, our results suggest that the frame effect on social discounting was mediated by the interplay between insula and VMPFC in the loss frame, and between TPJ and VMPFC in the gain frame. Thus, we maintain that the most parsimonious explanation of insula activation and its negative modulatory interplay with VMPFC is indeed a frame-dependent downregulation of own-reward values in the valuation network during social discounting, thus decreasing participants' selfishness, while TPJ-VMPFC coupling in the gain frame reflects the upregulation of vicarious reward value signals in VMPFC, hence promoting altruism by increasing the attractiveness of the generous option. Thus, in brief, the different motives underlying generosity in the gain and the loss frame

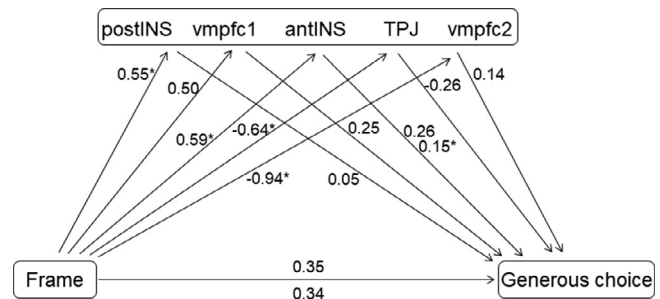


Fig. 7. Mediation analysis. A mediation model was built to clarify the effect of frame conditions (X) on generous choice (Y). Neural activations were entered as parallel potential mediators. Numbers are partially standardized effects. * refers to significant effects ($p < 0.05$). Where two numbers for the same path are reported, the one on the top refers to the direct effect and the one on the bottom refers to the indirect effect (i.e. when the mediators are included in the model).

are reflected by differential, frame-dependent activation and connectivity patterns in the brain.

4. Discussion

We provide behavioral and neural evidence for a simple nudge that aims at increasing individuals' willingness to provide costly support to socially remote others. We adapted a social discounting task where participants chose between a selfish option – a high gain to self and zero-gain to the other – and a generous option – a lower gain to self and equal gain to the other (Soutschek et al., 2016; Strombach et al., 2015). Based on previous evidence that people are strongly reluctant to increase their own payoff at the expense of others' welfare (Baumeister et al., 1994; Chang et al., 2011; Chang and Sanfey, 2013; Crockett et al., 2014), we framed the generous option either as a monetary gain to the other (gain frame), or as the prevention of the loss of a previous monetary endowment to the other (loss frame) (Everett et al., 2015; Li et al., 2017; List, 2007; Liu et al., 2020; Sip et al., 2015; Smith et al., 2015; Wang et al., 2017; Xiao et al., 2016; Zheng et al., 2010). Crucially, between frames, the choice alternatives differed only in the description of the decision problem, but not with regard to their actual economic consequences. In a series of four independent studies, we show that the social discount function was significantly flatter in the loss than the gain frame, indicating that participants were more generous towards socially remote others if a personal gain implied the other's loss of their previous endowment. Notably, our incentivization procedure made it logically impossible for the other persons to know about their endowment, or the potential loss of it, and participants were explicitly instructed about this; all that mattered was the final positive payoff to self and others. Yet, the fact that our participants were still reluctant to inflict losses to others suggests that they had internalized the social norm of not taking away money from others to such a degree that it prevailed even in the absence of any real economic consequences for others.

We hypothesized that the frame-dependent motives underlying generosity are dissociable on the neural level. Consistent with our previous work (Strombach et al., 2015), we found that generosity in the gain frame recruited a network of structures, including VMPFC and rTPJ, known to represent vicarious reward value and prosocial behavior. By contrast, in the loss frame, we expected that the reluctance to maximize own-gain at the expense of other-loss would be ideally mediated by social norm compliance and associated social sentiments, such as the negative emotions experienced during social norm transgressions, e.g., guilt and shame, as well as the aversive experience of unfairness and inequality. We therefore hypothesized that increased activity in brain regions associated with such social sentiments, specifically the insular cortex, would be associated with generous choices in the loss frame specifically (Bellucci et al., 2018; Canessa et al., 2017, 2013; Civai et al., 2012;

Corradi-Dell'Acqua et al., 2013; Huggins et al., 2018; Lallement et al., 2013; Lamm et al., 2011; Montague et al., 2007; Oldham et al., 2018; Paulus et al., 2003; Samanez-Larkin et al., 2008; Siebenthal et al., 2017; Singer et al., 2006; Sokol-Hessner et al., 2013; Sokol-Hessner and Rutledge, 2019; Spitzer et al., 2007; Tomasino et al., 2013; Wagner et al., 2011; Wang et al., 2017; Yu et al., 2014). We indeed found that the anterior insula was significantly more activated when participants made generous choices in the loss frame, relative to the gain frame. Extending these findings, we found that also the posterior part of the insula seemed to be involved in these processes, specifically supporting the representation of the loss frame information even before the decision was made (see also Drouman et al., 2015). Building upon this evidence, we further explored how both activation clusters mediated frame-specific social discounting behavior. We propose and provide empirical support for a network model that predicts that the frame effect on social discounting was associated with a frame-dependent neural connectivity pattern between insula and VMPFC in the loss frame, and TPJ and VMPFC in the gain frame. More specifically, DCM confirmed that posterior insula activation at loss frame onset exerted a negative modulatory effect on VMPFC. It is tempting to speculate that a frame-dependent downregulation of own-reward values in the valuation network during social discounting might lie at the core of the enhanced generosity observed in the loss frame. By contrast, the same analyses confirmed TPJ-VMPFC coupling in the gain frame, consistent with our previous finding (Strombach et al., 2015) that altruism in the gain frame is promoted by increasing the attractiveness of the generous option through TPJ-related upregulation of vicarious reward value signals in the valuation network. Overall, these results call for the idea that the motives behind generosity are likely qualitatively different in the gain and the loss frame, and dissociable on the neural level.

Our analyses revealed two separate clusters within insula; while a more posterior cluster was activated in response to general loss frame information, the more anterior cluster was specific to generous choices in the loss frame. This topographic dissociation within insula is consistent with previous findings suggesting a regional gradient in representing the level of abstraction of social sentiments during moral decision-making (e.g. Drouman et al., 2015; Ying et al., 2018). This pattern of result is in line with the idea that anterior and posterior insula may not subservise qualitatively different functions, but rather reflect different aspects of the same function, such as the interoceptive and visceral aspects of social sentiments in response to vicarious feelings of potential loss (posterior cluster), and their relevance for choice selection (anterior cluster). In addition to this, it is worth noting that, unlike insula, TPJ activity was only found at the decision stage but not at the frame information stage at trial onset. This time difference in activation allows for some speculation on the frame-dependent choice dynamics: it is conceivable that, in the loss frame, the frame information at trial onset signals the frame context, and, hence, prompts the tendency to make generous choices largely independent of social distance or selfish reward magnitude. Thus, the decision to be generous in the loss frame would be determined at trial onset already, and it would not be influenced by subsequently presented social distance or reward information. By contrast, in the gain frame, participants trade off selfish (own-payoff maximization) with other-regarding motives (granting others a gain) in a social-distance-dependent way (Strombach et al., 2015). This conflict between selfish and other-regarding motives can only be resolved when all information on frame type, social distance, and own-reward magnitude is available, that is, at the decision stage. Therefore, it is interesting as well as plausible to speculate that, while TPJ might be involved only at a later stage of the decision process, posterior insula might signal the frame context, and, hence, prompt generosity already at trial onset. Future research needs to clarify the specific functional differentiation of anterior versus posterior insula, as well as TPJ, contributions to social economic decision-making, by using, for instance, finite impulse response (FIR) models or mental chronometry approach (e.g. Menon, 2012; Schilbach et al., 2008).

Our findings expand on previous evidence that preventing harm to others is a great motivator of prosocial performance (Everett et al., 2015; Wang et al., 2017; Xiao et al., 2016; Zhang et al., 2017; Zheng et al., 2010). However, while others have found that harm prevention was particularly pronounced in a public context (Everett et al., 2015) and dependent on social feedback (Sip et al., 2015; Smith et al., 2015), we show that similar cognitive mechanisms can strongly boost generosity even in a private context and in the absence of social feedback, thus independent of reputational concerns, judgment by social peers, or third-party punishment threats. This suggests that other-harm prevention might be an internalized motive that works unconditionally and universally across contexts, regardless of social consequences. In addition, previous experiments on harm prevention did not manipulate, or provide information on, social distance between donor and recipient (Bardsley, 2008; Crockett et al., 2014; Everett et al., 2015; Li et al., 2017; Liu et al., 2020; Xiao et al., 2016). Hence, while the effects of the resource allocation mode on social discounting were elusive so far, our findings imply that it matters: harm-prevention motives in the loss frame were less dependent on social distance than other-regarding considerations in the gain frame, thus resulting in flatter social discounting.

A recent study used a similar framing manipulation and also reported TPJ involvement (Liu et al., 2020). However, their study differed from ours in several important ways. First, the task in Liu et al. (2020) involved trading off own-wealth maximization with avoiding electric shocks to others. However, their task did not involve social distance information about the recipients of shocks. Second, Liu et al. (2020) did not reveal any insula recruitment, or insula-VMPFC connectivity - the core finding in our study - related to generosity or task framing. Most importantly, perhaps, while Liu et al. (2020) identified TPJ-VMPFC connectivity to be relevant for their frame-related increase in costly harm-prevention, we found instead that insula-VMPFC connectivity was associated with the frame-related boost in generosity during social discounting. This suggests that Liu et al. (2020) most likely studied different framing-related cognitive and neural mechanisms than the ones investigated here.

Our results are consistent with the idea that certain costly altruistic behaviors are not motivated by genuinely other-regarding considerations, but instead by compliance to internalized social norms. But what impels participants to comply to social norms? Here, we propose, along with previous evidence (Chang et al., 2011; Spitzer et al., 2007), that compliance to social norms might be linked to anticipated feelings of guilt, shame, and remorse, and accompanied by insula activation (see also Belfi et al., 2015; Sellitto et al., 2016), which ultimately sustain prosocial behavior. According to this view, insula would reflect the negative sentiment associated with social norm transgressions as they occur when being responsible for someone else's loss (i.e. vicarious loss experience). Our data show that this social sentiment and accompanying neural signature can be elicited even when the others' outcomes are merely described as losses, thus, in the absence of real losses to others.

The success of our framing manipulation in increasing generosity came at a methodological cost: because participants rarely made selfish choices in the loss frame, the analyses on selfish choices were underpowered. Hence, results involving selfish decisions, and how they map on insula, TPJ, or VMPFC, have to be interpreted with caution. To shed more light on the neural correlates of selfish choices in the loss frame, future studies should replicate our experiment with a slightly less effective nudge that would allow for more selfish choices.

5. Conclusions

The acceptance and support of the principle of a caring society, and the attitude towards the welfare of socially remote strangers, is central for a civilization to function well. It seems vital for societies to successfully meet current challenges, such as integrating refugees, addressing economic inequality, acceding the trials and promises of a globalized world (Kalenscher, 2014), or managing the public health implications of

the current COVID-19 pandemic. Here, we present a simple behavioral framing manipulation that boosts generosity towards socially remote others: framing a selfish choice as a loss to others can motivate prosocial behavior, even if the framing of the choice options is irrelevant for the actual payoff to others. Our neuroimaging data identify insula as the core component in a network associated with this enhanced generosity in the loss frame. Our results imply that prosocial attitudes towards others are highly malleable and strongly depend on the architecture of the decision problem. The insights gained in this study might, thus, help in designing policies aimed at increasing the acceptance and support of the principle of a caring society, and to change the attitude towards the welfare of socially remote strangers.

Credit authorship contribution statement

Manuela Sellitto: Conceptualization, Visualization, Investigation, Formal analysis, Writing – original draft, Writing – review & editing. **Susanne Neufang:** Writing – original draft, Writing – review & editing. **Adam Schweda:** Visualization, Investigation, Writing – review & editing. **Bernd Weber:** Visualization, Writing – review & editing. **Tobias Kalenscher:** Conceptualization, Writing – original draft, Writing – review & editing.

Acknowledgments

This work was supported by Deutsche Forschungsgemeinschaft (DFG) grant no. KA 2675/4–3 (to T.K.).

Data and code availability statement

Digital Imaging and Communications in Medicine (DICOM) images reported in this paper have been deposited in XNAT Central (<https://central.xnat.org/>) under the project name 'Framesocdisc'. Behavioral datasets have been supplied in Figshare (<https://figshare.com/>) under the doi: 10.6084/m9.figshare.10265309. The code used for fMRI analyses was retrieved from the statistical parametric mapping (SPM) platform.

Supplementary materials

Supplementary material associated with this article can be found, in the online version, at doi:10.1016/j.neuroimage.2021.118211.

References

- Archambault, C., Kalenscher, T., De Laat, J., 2019. Generosity and livelihoods: experimental evidence on the multidimensional nature of sharing among the Kenyan Maasai. *J. Behav. Decis. Mak.* 1–12.
- Bardsley, N., 2008. Dictator game giving: altruism or artefact? *Exp. Econ.* 11, 122–133. doi:10.1007/s10683-007-9172-2.
- Bartra, O., McGuire, J.T., Kable, J.W., 2013. The valuation system: a coordinate-based meta-analysis of BOLD fMRI experiments examining neural correlates of subjective value. *Neuroimage* 76, 412–427. doi:10.1016/j.neuroimage.2013.02.063.
- Baumeister, R.F., Stillwell, A.M., Heatherton, T.F., 1994. Guilt: An Interpersonal Approach. *Psychol. Bulletin.* 115, 243–267.
- Belfi, A.M., Kosciak, T.R., Tranel, D., 2015. Damage to the insula is associated with abnormal interpersonal trust. *Neuropsychologia* 71, 165–172. doi:10.1016/j.neuropsychologia.2015.04.003.
- Bellucci, G., Feng, C., Camilleri, J., Eickhoff, S.B., Krueger, F., 2018. The role of the anterior insula in social norm compliance and enforcement: evidence from coordinate-based and functional connectivity meta-analyses. *Neurosci. Biobehav. Rev.* 92, 378–389. doi:10.1016/j.neubiorev.2018.06.024.
- Canessa, N., Crespi, C., Baud-BOVY, G., Dodich, A., Falini, A., Antonellis, G., Cappa, S.F., 2017. Neural markers of loss aversion in resting-state brain activity. *Neuroimage* 146, 257–265. doi:10.1016/j.neuroimage.2016.11.050.
- Canessa, N., Crespi, C., Motterlini, M., Baud-BOVY, G., Chierchia, G., Pantaleo, G., Tettamanti, M., Cappa, S.F., 2013. The functional and structural neural basis of individual differences in loss aversion. *J. Neurosci.* 33, 14307–14317. doi:10.1523/JNEUROSCI.0497-13.2013.
- Chang, L.J., Sanfey, A.G., 2013. Great expectations: neural computations underlying the use of social norms in decision-making. *Soc. Cognit. Affect. Neurosci.* 8, 277–284. doi:10.1093/scan/nsr094.

- Chang, L.J., Smith, A., Dufwenberg, M., Sanfey, A.G., 2011. Article triangulating the neural, psychological, and economic bases of guilt aversion. *Neuron* 70, 560–572. doi:10.1016/j.neuron.2011.02.056.
- Cho, Y.T., Fromm, S., Guyer, A.E., Detloff, A., Pine, D.S., Fudge, J.L., Ernst, M., 2013. Nucleus accumbens, thalamus and insula connectivity during incentive anticipation in typical adults and adolescents. *Neuroimage* 66, 508–521. doi:10.1016/j.neuroimage.2012.10.013.
- Civai, C., Crescentini, C., Rustichini, A., Rumiati, R.I., 2012. Equality versus self-interest in the brain: differential roles of anterior insula and medial prefrontal cortex. *Neuroimage* 62, 102–112. doi:10.1016/j.neuroimage.2012.04.037.
- Corradi-Dell'Acqua, C., Civai, C., Rumiati, R.I., Fink, G.R., 2013. Disentangling self- and fairness-related neural mechanisms involved in the ultimatum game: an fMRI study. *Soc. Cognit. Affect. Neurosci.* 8, 424–431. doi:10.1093/scan/nss014.
- Crockett, M.J., Kurth-nelson, Z., Siegel, J.Z., Dayan, P., Dolan, R.J., Crockett, M.J., Kurth-nelson, Z., Siegel, J.Z., Dayan, P., Dolan, R.J., 2014. Harm to others outweighs harm to self in moral decision making. *Proc. Natl. Acad. Sci.* 111, 17320–17325. doi:10.1073/pnas.1424572112.
- Davis, M.H., 1983. A multidimensional approach to individual differences in empathy. *J. Pers. Soc. Psychol.* 44, 113–126. doi:10.1037/0022-3514.44.1.113.
- Dreu, D., 1997. Gain-loss frames and cooperation in two-person social dilemmas: a transformational analysis. *Sth. J. Person. Soc. Psychol.* pp. 1093–1106.
- Droutman, V., Bechara, A., Read, S.J., 2015. Roles of the different sub-regions of the insular cortex in various phases of the decision-making process. *Front. Behav. Neurosci.* 9, 309. doi:10.3389/fnbeh.2015.00309.
- Eickhoff, S.B., Paus, T., Caspers, S., Grosbras, M.H., Evans, A.C., Zilles, K., Amunts, K., 2007. Assignment of functional activations to probabilistic cytoarchitectonic areas revisited. *Neuroimage* 36, 511–521. doi:10.1016/j.neuroimage.2007.03.060.
- Evans, A.M., Beest, I.V., 2017. Gain-loss framing effects in dilemmas of trust and reciprocity. *J. Exp. Soc. Psychol.* 73, 151–163. doi:10.1016/j.jesp.2017.06.012.
- Everett, J.A.C., Faber, N.S., Crockett, M.J., 2015. The influence of social preferences and reputational concerns on intergroup prosocial behaviour in gains and losses contexts. *R. Soc. Open Sci.* 2, 150546. doi:10.1098/rsos.150546.
- Friston, K.J., Harrison, L., Penny, W., 2003. Dynamic causal modelling. *Neuroimage* 19, 1273–1302. doi:10.1016/S1053-8119(03)00202-7.
- Gu, X., Wang, X., Hula, X.A., Wang, S., Xu, S., Lohrenz, T.M., Knight, R.T., Gao, Z., Dayan, P., Montague, P.R., 2015. Necessary, yet dissociable contributions of the insular and ventromedial prefrontal cortices to norm adaptation: computational and lesion evidence in humans. *J. Neurosci.* 35, 467–473. doi:10.1523/JNEUROSCI.2906-14.2015.
- Hayes, A.F., 2017. *Introduction to mediation, moderation, and Conditional Process Analysis: A regression-Based Approach*. Guilford Publications, New York.
- Huggins, A.A., Belleau, E.L., Miskovich, T.A., Pedersen, W.S., Larson, C.L., 2018. Moderating effects of harm avoidance on resting-state functional connectivity of the anterior insula. *Front. Hum. Neurosci.* 12, 1–9. doi:10.3389/fnhum.2018.00447.
- Hutcherson, C.A., Bushong, B., Rangel, A., 2015. A neurocomputational model of altruistic choice and its implications. *Neuron* 87, 451–463. doi:10.1016/j.neuron.2015.06.031.
- Jenkinson, M., Beckmann, C.F., Behrens, T.E.J., Woolrich, M.W., Smith, S.M., 2012. FSL. *Neuroimage* 62, 782–790. doi:10.1016/j.neuroimage.2011.09.015.
- Jones, B., Rachlin, H., 2006. Social discounting. *Psychol. Sci.* 17, 283–286. doi:10.1111/j.1467-9280.2006.01699.x.
- Kalenscher, T., 2017. Social psychology: love thy stranger as thyself. *Nat. Hum. Behav.* 1, 0108. doi:10.1038/s41562-017-0108.
- Kalenscher, T., 2014. Attitude toward health insurance in developing countries from a decision-making perspective. *J. Neurosci. Psychol. Econ.* 7:3, 174–193.
- Lallement, J.H., Kuss, K., Trautner, P., Weber, B., Falk, A., Fliessbach, K., 2013. Effort Increases Sensitivity to Reward and Loss Magnitude in the Human Brain. *Soc. Cogn. Affect. Neurosci.* pp. 342–349.
- Lamm, C., Decety, J., Singer, T., 2011. NeuroImage Meta-analytic evidence for common and distinct neural networks associated with directly experienced pain and empathy for pain. *Neuroimage* 54, 2492–2502. doi:10.1016/j.neuroimage.2010.10.014.
- Li, R., Smith, D.V., Clithero, J.A., Venkatraman, V., Carter, R.M., Huettel, S.A., 2017. Reason's enemy is not emotion: engagement of cognitive control networks explains biases in gain/loss framing. *J. Neurosci.* 37. doi:10.1523/JNEUROSCI.3486-16.2017, 3486–16.
- Li, Z., Yan, C., Xie, W.Z., Li, K., Zeng, Y.W., Jin, Z., Cheung, E.F.C., Chan, R.C.K., 2015. Anticipatory pleasure predicts effective connectivity in the mesolimbic system. *Front. Behav. Neurosci.* 9, 1–8. doi:10.3389/fnbeh.2015.00217.
- List, J., 2007. On the interpretation of giving in dictator games. *J. Polit. Econ.* 115, 482–493. doi:10.1086/519249.
- Liu, J., Gu, R., Liao, C., Lu, J., Fang, Y., Xu, P., Luo, Y., Cui, F., 2020. The neural mechanism of the social framing effect: evidence from fMRI and tDCS studies. *J. Neurosci.* 40. doi:10.1523/JNEUROSCI.1385-19.2020, JN-RM-1385-19.
- Margittai, Z., Strombach, T., van Wingerden, M., Joëls, M., Schwabe, L., Kalenscher, T., 2015. A friend in need: time-dependent effects of stress on social discounting in men. *Horm. Behav.* 73, 75–82. doi:10.1016/j.yhbeh.2015.05.019.
- Margittai, Z., Van Wingerden, M., Schnitzler, A., Joëls, M., Kalenscher, T., 2018. Dissociable roles of glucocorticoid and noradrenergic activation on social discounting. *Psychoneuroendocrinol.* doi:10.1016/j.psyneuen.2018.01.015.
- Brett, M., Anton, J.L., Romain Valabregue, J.B.P., 2002. Region of interest analysis using an SPM toolbox. In: *Proceedings of the 8th International Conference on Functional Mapping of the Human Brain* doi:10.1007/978-3-540-75387-2_17.
- Mazaika, P.K., Hoefft, F., Glover, G.H., Reiss, A.L., 2009. *Methods and Software for fMRI Analysis for Clinical Subjects*, 47:1. *NeuroImage*, p. S58.
- Menon, R.S., 2012. Mental chronometry. *Neuroimage* 62, 1068–1071. doi:10.1016/j.neuroimage.2012.01.011.

- Robbs, D., Yu, R., Meyer, M., Passamonti, L., Seymour, B., Calder, A.J., Schweizer, S., Frith, C.D., Dalglish, T., 2009. A key role for similarity in vicarious reward. *Science* 324, 900. doi:10.1126/science.1170539.
- Montague, P.R., Lohrenz, T., Humphrey, H., 2007. To detect and correct: norm violations and their enforcement. *Neuron* 14–18. doi:10.1016/j.neuron.2007.09.020.
- Neufang, S., Akhrif, A., Herrmann, C.G., Drepper, C., Homola, G.A., Nowak, J., Waider, J., Schmitt, A.G., Lesch, K.P., Romanos, M., 2016. Serotonergic modulation of “waiting impulsivity” is mediated by the impulsivity phenotype in humans. *Transl. Psychiatry* 6. doi:10.1038/tp.2016.210.
- Nowak, M.A., 2006. Five rules for the evolution of cooperation. *Science*, 314:5805, 1560–1563. doi:10.1126/science.1133755.
- Oldham, S., Murawski, C., Fornito, A., Youssef, G., Lorenzetti, V., 2018. The anticipation and outcome phases of reward and loss processing: a neuroimaging meta-analysis of the monetary incentive delay task. *Hum. Brain Mapp.* 3398–3418. doi:10.1002/hbm.24184.
- Paulus, C., 2009. Der Saarbrücker Persönlichkeitsfragebogen (SPF-IRI) zur Messung von Empathie. 1–11.
- Paulus, M.P., Rogalsky, C., Simmons, A., Feinstein, J.S., Stein, M.B., 2003. Increased activation in the right insula during risk-taking decision making is related to harm avoidance and neuroticism. *NeuroImage* 19, 1439–1448. doi:10.1016/S1053-8119(03)00251-9.
- Rilling, J.K., Sanfey, A.G., 2011. *The Neuroscience of Social Decision-Making*, 62. *The Annual Review of Psychology*, pp. 23–48.
- Rolls, E.T., Huang, C.C., Lin, C.P., Feng, J., Joliet, M., 2020. Automated anatomical labelling atlas 3. *Neuroimage* 206, 116189. doi:10.1016/j.neuroimage.2019.116189.
- Samanez-Larkin, G.R., Hollon, N.G., Carstensen, L.L., Knutson, B., 2008. *Psychol. Sci.* 19, 320–323. doi:10.1111/j.1467-9280.2008.02087.x.Individual.
- Saxe, R., Kanwisher, N., 2003. People thinking about thinking people: the role of the temporo-parietal junction in “theory of mind. *NeuroImage* 19, 1835–1842. doi:10.1016/S1053-8119(03)00230-1.
- Schilbach, L., Eickhoff, S.B., Mojisich, A., Vogeley, K., 2008. What’s in a smile? Neural correlates of facial embodiment during social interaction. *Soc. Neurosci.* 3, 37–50. doi:10.1080/17470910701563228.
- Schweda, A., Margittai, Z., Kalenscher, T., 2020. Acute stress counteracts framing-induced generosity boosts in social discounting in young healthy men. *Psychoneuroendocrinol.* 121, 104860. doi:10.1016/j.psyneuen.2020.104860.
- Seara-Cardoso, A., Sebastian, C.L., Mccroy, E., Foulkes, L., Buon, M., Roiser, J.P., Viding, E., 2016. Anticipation of guilt for everyday moral transgressions: the role of the anterior insula and the influence of interpersonal psychopathic traits. *Sci. Rep.* 6, 36273. doi:10.1038/srep36273.
- Sellitto, M., Ciaramelli, E., Mattioli, F., di Pellegrino, G., 2016. Reduced sensitivity to sooner reward during intertemporal decision-making following insula damage in humans. *Front. Behav. Neurosci.* 9, 367. doi:10.3389/fnbeh.2015.00367.
- Singer, T., Seymour, B., Doherty, J.P.O., Stephan, K.E., Dolan, R.J., Frith, C.D., 2006. Empathic neural responses are modulated by the perceived fairness of others. *Nature* 439, 10–13. doi:10.1038/nature04271.
- Sethi, R., Somanathan, E., E., 1996. *The evolution of social norms in common property resource use. Am. Econ. Rev.* 86 (4), 766–788 (Sep., 1996)Published by : American Econom 86, 766–788.
- Sip, K.E., Smith, D.V., Porcelli, A.J., Kar, K., Delgado, M.R., 2015. Social closeness and feedback modulate susceptibility to the framing effect. *Soc. Neurosci.* 10, 35–45. doi:10.1080/17470919.2014.944316.
- Smith, D.V., Sip, K.E., Delgado, M.R., 2015. Functional connectivity with distinct neural networks tracks fluctuations in gain/loss framing susceptibility. *Hum. Brain Mapp.* 36, 2743–2755. doi:10.1002/hbm.22804.
- Sokol-Hessner, P., Camerer, C.F., Phelps, E.A., 2013. Emotion Regulation Reduces Loss Aversion and Decreases Amygdala Responses to Losses, 8:13. *Soc. Cogn. Affect. Neurosci.*, pp. 341–350.
- Sokol-Hessner, P., Rutledge, R.B., 2019. *The Psychological and Neural Basis of Loss Aversion*, 28:1. *Curr. Dir. Psychol. Sci.*, pp. 20–27.
- Soutschek, A., Ruff, C.C., Strombach, T., Kalenscher, T., Tobler, P.N., 2016. Brain stimulation reveals crucial role of overcoming self-centeredness in self-control. *Sci. Adv.* 2, e1600992. doi:10.1126/sciadv.1600992, –e1600992.
- Soutschek, A., Ugazio, G., Crockett, M.J., Ruff, C.C., Kalenscher, T., Tobler, P.N., 2017. Binding oneself to the mast: stimulating frontopolar cortex enhances precommitment. *Soc. Cognit. Affect. Neurosci.* 1–8. doi:10.1093/scan/nsw176.
- Spitzer, M., Fischbacher, U., Herrnberger, B., Grön, G., Fehr, E., 2007. The neural signature of social norm compliance. *Neuron* 56, 185–196. doi:10.1016/j.neuron.2007.09.011.
- Stöber, J., 2001. The Social Desirability Scale-17 (SDS-17): convergent validity, discriminant validity, and relationship with age. *Eur. J. Psychol. Assess.* 17, 222–232. doi:10.1027//1015-5759.17.3.222.
- Strang, S., Gerhardt, H., Marsh, N., Oroz, S., Hu, Y., Hurlmann, R., Park, S.Q., 2017. Short communication a matter of distance — the effect of oxytocin on social discounting is empathy-dependent. *Psychoneuroendocrinol.* 78, 229–232. doi:10.1016/j.psyneuen.2017.01.031.
- Strombach, T., Jin, J., Weber, B., Kenning, P., Shen, Q., Ma, Q., Kalenscher, T., 2014. Charity begins at home: cultural differences in social discounting and generosity. *J. Behav. Decis. Mak.* 27, 235–245. doi:10.1002/bdm.1802.
- Strombach, T., Margittai, Z., Gorczyca, B., Kalenscher, T., 2016. Gender-specific effects of cognitive load on social discounting. *PLoS One* 11, 1–15. doi:10.1371/journal.pone.0165289.
- Strombach, T., Weber, B., Hangebrauk, Z., Kenning, P., Karipidis, I.I., Tobler, P.N., Kalenscher, T., 2015. Social discounting involves modulation of neural value signals by temporoparietal junction. *Proc. Natl. Acad. Sci.* 112, 1619–1624. doi:10.1073/pnas.1414715112.
- Studer, B., Koch, C., Knecht, S., Kalenscher, T., 2019. Conquering the inner couch potato: precommitment is an effective strategy to enhance motivation for effortful actions. *Philos. Trans. R. Soc. Lond. B Biol. Sci.* 374, 20180131. doi:10.1098/rstb.2018.0131.
- Tomasino, B., Lotto, L., Sarlo, M., Civai, C., Rumiati, R., Rumiati, R.I., 2013. Framing the ultimatum game: the contribution of simulation. *Front. Hum. Neurosci.* 7, 1–16. doi:10.3389/fnhum.2013.00337.
- Tusche, X.A., Bo, A., Kanske, X.P., Trautwein, X.F., Singer, T., 2016. Decoding the charitable brain: empathy, perspective taking, and attention shifts differentially predict altruistic giving. *J. Neurosci.* 36, 4719–4732. doi:10.1523/JNEUROSCI.3392-15.2016.
- Vekaria, K.M., Brethel-Haurwitz, K.M., Cardinale, E.M., Stoycos, S.A., Marsh, A.A., 2017. Social discounting and distance perceptions in costly altruism. *Nat. Hum. Behav.* 1, 0100. doi:10.1038/s41562-017-0100.
- Vlaev, I., 2012. How different are real and hypothetical decisions? Overestimation, contrast and assimilation in social interaction. *J. Econ. Psychol.* 33, 963–972. doi:10.1016/j.joep.2012.05.005.
- Von Siebenthal, Z., Boucher, O., Rouleau, I., Lassonde, M., Lepore, F., Nguyen, D.K., De, C., 2017. Decision-making impairments following insular and medial temporal lobe resection for drug-resistant epilepsy. *Soc. Cogn. Affect. Neurosci.* 128–137. doi:10.1093/scan/nsw152.
- Wagner, U., N’Diaye, K., Ethofer, T., Vuilleumier, P., 2011. Guilt-specific processing in the prefrontal cortex. *Cereb. Cortex* 21, 2461–2470. doi:10.1093/cercor/bhr016.
- Wang, X.T., Rao, L.-L., Zheng, H., 2017. Neural substrates of framing effects in social contexts: a meta-analytical approach. *Soc. Neurosci.* 12, 268–279. doi:10.1080/1080/17470919.2016.1165285.
- Wendelken, C., O’Hare, E.D., Whitaker, K.J., Ferrer, E., Bunge, S.A., 2011. Increased functional selectivity over development in rostralateral prefrontal cortex. *J. Neurosci.* 31, 17260–17268. doi:10.1523/JNEUROSCI.1193-11.2011.
- Wiehler, A., Petzschner, F.H., Stephan, K.E., Peters, J., 2017. Episodic Tags Enhance Striatal Valuation Signals during Temporal Discounting in Pathological Gamblers, 4.
- Wilson, R.C., Collins, A.G.E., 2019. Ten simple rules for the computational modeling of behavioral data. *Elife* 8, 1–33. doi:10.7554/eLife.49547.
- Xiang, T., Lohrenz, T., Montague, P.R., 2013. Computational substrates of norms and their violations during social exchange. *J. Neurosci.* 33, 1099–1108. doi:10.1523/JNEUROSCI.1642-12.2013.
- Xiao, F., Zheng, Z., Zhang, H., Xin, Z., Chen, Y., Li, Y., 2016. Who are you more likely to help? the effects of expected outcomes and regulatory focus on prosocial performance. *PLoS One* 11, 1–15. doi:10.1371/journal.pone.0165717.
- Ying, X., Luo, J., Chiu, C.yue, Wu, Y., Xu, Y., Fan, J., 2018. Functional dissociation of the posterior and anterior insula in moral disgust. *Front. Psychol.* 9, 1–10. doi:10.3389/fpsyg.2018.00860.
- Yu, H., Hu, J., Hu, L., Zhou, X., 2014. The voice of conscience: neural bases of interpersonal guilt and compensation. *Soc. Cognit. Affect. Neurosci.* 1150–1158. doi:10.1093/scan/nst090.
- Zeidman, P., Jafarian, A., Corbin, N., Seghier, M.L., Razi, A., Price, C.J., Friston, K.J., 2019a. A guide to group effective connectivity analysis, part 1: first level analysis with DCM for fMRI. *Neuroimage* 200, 174–190. doi:10.1016/j.neuroimage.2019.06.031.
- Zeidman, P., Jafarian, A., Seghier, M.L., Litvak, V., Cagnan, H., Price, C.J., Friston, K.J., 2019b. A guide to group effective connectivity analysis, part 2: second level analysis with PEB. *Neuroimage* 200, 12–25. doi:10.1016/j.neuroimage.2019.06.032.
- Zhang, X., Liu, Yi, Chen, X., Shang, X., Liu, Yongfang, 2017. Decisions for others are less risk-averse in the gain frame and less risk-seeking in the loss frame than decisions for the self. *Front. Psychol.* 8, 1–10. doi:10.3389/fpsyg.2017.01601.
- Zhang, Y.Y., Xu, L., Liang, Z.Y., Wang, K., Hou, B., Zhou, Y., Li, S., Jiang, T., 2018. Separate neural networks for gains and losses in intertemporal choice. *Neurosci. Bull.* 34, 725–735. doi:10.1007/s12264-018-0267-x.
- Zheng, H., Wang, X.T., Zhu, L., 2010. Framing effects: behavioral dynamics and neural basis. *Neuropsychologia* 48, 3198–3204. doi:10.1016/j.neuropsychologia.2010.06.031.