

Neuroscience of Moral Decision Making

Yang Hu^{a,b,*}, Xiaoxue Gao^{a,b,*}, Hongbo Yu^c, Zhewen He^{b,d}, and Xiaolin Zhou^{a,b,e}

Department of Psychology, Beijing Normal University, Beijing 100875, China; Department of Psychology, Beijing Normal University, Beijing 100875, China; Department of Psychology, Beijing Normal University, Beijing 100875, China; Department of Psychology, Beijing Normal University, Beijing 100875, China; Department of Psychology, Beijing Normal University, Beijing 100875, China

Introduction	481
Moral Decision-Making in the Brain: A Multi-Stage Framework	482
V u -	483
Harm	483
Help	484
(Un)Fairness	485
(Dis)Honesty	486
betrayal	487
u	4
Direct Reciprocity	488
Indirect Reciprocity	489
r	40
Open Questions and Future Directions	490
Acknowledgments	491
References	491

Introduction

As moral agents, we human beings are equipped with the capability to make judgments about the moral appropriateness of the other's behaviors (i.e., moral judgment) (Baron, 2014; Malle et al., 2014; Wojciszke et al., 2015). However, a comprehensive and commonly agreed definition of morality (or moral domain), researchers often define morality from different perspectives (Bartels et al., 2014; Crockett, 2013; Haidt, 2007). Two of those are commonly adopted. The first approach highlights the moral domain as the behaviors that cause or prevent harm to others (e.g., lying, cheating, stealing, and harming others) (Baron, 2014). The second approach focuses on the behaviors that cause or prevent the other's suffering (Cushman, 2015). Accordingly, immorality refers to the behaviors

2012) and the underlying neural activation patterns (FeldmanHall et al., 2012; Gaspic et al., 2013) are different in hypothetical versus real contexts. Moreover, most of these previous studies were not designed to provide a mechanistic account for the moral behaviors,

the partner), which was measured by the no-provocation trials in the Taylor reaction-time aggression paradigm [Taylor, 1967](#)).

Notably, all studies mentioned above assumed that helping behaviors would surely reduce the other's suffering, which was not always true in real life. To address this issue, a study combining both fMRI and tDCS techniques developed a new paradigm in which

implying that higher demands in moral mentalizing are required in social decision-making when the decision to reject could not be readily justifi

revealing that the resting-state brain activity in the left ventral AI (as well as other regions) was correlated with the PIF response (Seu et al., 2015). Together, these findings suggest that the AI is not only engaged in signaling social norm violation during UG but also recruited in guiding subsequent adaptive behaviors (e.g., PIF response).

Learning

In real life, we not only make moral choices in one shot, but often need to form and update our beliefs about the moral trait of others, thereby guiding how we should get along with them in the future (Siegel et al., 2018). Although a substantial amount of evidence has revealed the neurocomputational mechanisms underlying how people learn through feedbacks under the general framework of reinforcement learning (O'Doherty et al., 2017), the neural underpinnings through which we infer the moral character of other people are still poorly understood. To investigate this issue, Hackel et al. (2015) performed a fMRI study in which participants were asked to learn how generous an anonymous partner was via trial-and-error learning based on the proportion of resources shared by the partner. As a control condition, participants also needed to learn which slot machine earned themselves more. Model-based analyses revealed that participants relied more on generosity information than on reward value during the task. Trial-wise prediction error (PE) of both types of information was commonly encoded in the right VS. However, the generosity prediction error recruited an additional network in association with the formation of social impression, including the ventral lateral prefrontal cortex (vlPFC), IPL, PCC extending to precuneus, as well as the right TPJ. Another study with a similar learning paradigm also found a signal of generosity PE in the PCC/precuneus (Stanley, 2010). Furthermore, our ability to infer others' moral character

The third issue is related to methodological approaches that should be taken to provide additional information from different viewpoints, thereby characterizing a panoramic view of the moral brain. Obviously, the current literature predominantly considers which parts of the brain (and the inter-regional connections) are associated with a specific form of moral decision using fMRI, supplemented by the causality methods such as brain lesion and non-invasive brain stimulation (e.g., TMS, tDCS). There have been several studies adopting the EEG technique (e.g., event-related potential, ERP) to explore the temporal features of moral decision

Garrett, N., Lazzaro, S.C., Ariely, D., Sharot, T., 2016. The brain adapts to dishonesty. *Nat. Neurosci.* 19 (12), 1727.

Garrigan, B., Adlam, A.L., Langdon, P.E., 2016. The neural correlates of moral decision-making: a systematic review and meta-analysis of social decision judgements. *Brain Cognit.* 109, 788.

Gert, B., 2004. *Common Morality: Deciding what to Do*. Oxford University Press.

Ginther, M.R., Bonnie, R.J., Hoffman, M.B., Shen, F.X., Simons, K.W., Jones, O.D., Marois, R., 2016. Parsing the behavioral and neural mechanisms of J. *Neurosci.* 36 (36), 9423-4.

Gneezy, U., 2005. Deception: the role of consequences. *Am. Econ. Rev.* 95 (1), 384.

Gospic, K., Sundberg, M., Maeder, J., Fransson, P., Petrovic, P., Isacsson, G., et al., 2014. Altruistic signals from amygdala. *Soc. Cognit. Affect Neurosci.* 9 (9), 1325-1332.

Graham, J., Haidt, J., Koleva, S., Motyl, M., Iyer, R., Wojcik, S.P., Ditto, P.H., 2013. Moral foundations theory: the pragmatic Axioms of Soc. Psychol. 47, 55-130.

Graham, J., Nosek, B.A., Haidt, J., Iyer, R., Koleva, S., Ditto, P.H., 2011. Mapping the moral domain. *J. Pers. Soc. Psychol.* 101 (2), 366.

Greccucci, A., Giorgetta, C., Waut, M., Bonini, N., Sanfey, A.G., 2012. Reappraising the ultimatum: an fMRI study of emotion regulation and decision making. *Cerebr. Cor.* (2), 399410.

Greene, J.D., 2015. The cognitive neuroscience of moral judgment and decision-making. In: Gazzaniga, M.S., Wheatley, T. (Eds.), *The Moral Brain: A m* pp. 197-220. Boston (Review).

Greene, J.D., Paxton, J.M., 2009. Patterns of neural activity associated with honest and dishonest moral decisions. *Proc Natl. Acad. Sci. U. S. A.* 106 (10), 12250-12255.

Guo, X., Zheng, L., Zhu, L., Li, J., Wang, Q., Dienes, Z., Yang, Z., 2013. Increased neural responses to unfairness in a loss event. *Neuroimage* 77, 246-250.

Gürgü, B., van den Bos, W., Rombouts, S.A., Crone, E.A., 2010. Unfair? It depends: neural correlates of fairness in social decisions. *PLoS One* 9 (9), e10799.

Haber, S.N., Knutson, B., 2009. The reward circuit: linking primate anatomy and human imaging. *Neuropsychopharmacology* 35 (1), 4-16.

Hackel, L.M., Amodio, D.M., 2018. Computational neuroscience approaches to social cognition. *Curr. Opin. Psychol.* 24, 92-100.

Hackel, L.M., Doll, B.B., Amodio, D.M., 2015. Instrumental learning of traits versus rewards: dissociable neural correlates. *Neurosci. Biobehav. Rev.* 59, 122-135.

Haidt, J., 2003. The moral emotions. In: Scherer, K.R., Goldsmith, H.H. (Eds.), *Handbook of Affective Sciences*, vol. 1. Oxford, pp. 85-97. University Press, O

Haidt, J., 2007. The new synthesis in moral psychology. *Science* 316 (5827), 998-1000.

Haidt, J., 2008. *Morality*. *Perspect. Psychol. Sci.* 3 (1), 65-72.

Harbaugh, W.T., Mayr, U., Burghart, D.R., 2007. Neural responses to taxation and voluntary giving reveal motives for charitable contributions. *Science* 316 (5811), 1621-1625.

Hare, T.A., Camerer, C.F., Knetsch, J.P., Rangel, A., 2010. Value computations in ventral medial prefrontal cortex during charitable decision making incorporate input from regions involved in social cognition. *J. Neurosci.* 30 (2), 583-590.

Hare, T.A., Schultz, W., Camerer, C.F., Knetsch, J.P., Rangel, A., 2011. Transformation of stimulus value signals into motor commands during simple choice. *Proc. Natl. Acad. Sci. U. S. A.* 108 (44), 18120-18125.

Haruno, M., Frith, C.D., 2010. Activity in the amygdala elicited by unfair divisions predicts social value orientation. *Nat. Neurosci.* 13 (2), 160-166.

Hauser, M., Lee, J., Huebner, B., 2010. The moral-conventional distinction in mature moral competence. *J. Cognit. Cult.* 10 (1), 1-26.

Hsu, M., Anen, C., Quartz, S.R., 2008. The right and the good: distributive justice and neural encoding of severity. *Cerebr. Cor.* 20 (5879), 1095-1102.

Hu, J., Hu, Y., Li, Y., Zhou, X., 2021. Computational and neurobiological substrates of prosocial helping decision. *J. Neurosci.* 41 (35), 3545-3555.

Hu, J., Li, Y., Yin, Y., Blue, P.R., Yu, H., Zhou, X., 2018. How do self-interest and other-need interact in the brain to determine altruistic behavior? *Neuroimage* 157, 59-69.

Hu, Y., He, L., Zhang, L., Wolk, T., Dreher, J.C., Weber, B., 2018. Spreading inequality: neural computations underlying giving and receiving. *J. Neurosci.* 38 (6), 575-589. <https://doi.org/10.1093/scan/nsy040>

Hu, Y., Pereira, A.M., Gao, X., Campos, B.M., Derrington, E., Corngnet, B., et al., 2021. Right temporoparietal junction trade-off of social decision spectrum disorder. *J. Neurosci.* 41 (8), 1799-1811. <https://doi.org/10.1523/JNEUROSCI.1237-20.2020>

Hu, Y., Scheele, D., Becker, B., Voos, G., David, B., Hurlmann, R., Weber, B., 2016. The effect of oxytocin on third-party altruistic decisions. *PLoS Stud. Sci. Rep.* 6, 20236. <https://doi.org/10.1038/srep20236>

Hu, Y., Strang, S., Weber, B., 2015. Helping or punishing strangers: neural correlates of altruistic decisions as third-party observers. *J. Neurosci.* 35 (24), 8411-8424. <https://doi.org/10.1523/JNEUROSCI.0002-15.2015>

Hutcherson, C., Bushong, B., Rangel, A., 2015. A neurocomputational model of altruistic choice and its implications. *Neuron* 87 (2), 451-462.

Hwang, C.-L., Lin, M.-J., 2012. *Group Decision Making under Multiple Criteria: Methods and Applications*, vol. 281. Springer Science & Business Media.

Ismayilov, H., Potters, J.J.J., 2015. Promises as Commitments.

Izuma, K., Saito, D.N., Sadato, N., 2010. Processing of the incentive for social approval in the ventral striatum duringitch. *Soc. Cognit. Affect Neurosci.* 5 (4), 630-639.

Kahneman, D., Knetsch, J.L., Thaler, R., 1986a. Fairness as a constraint on profit-making in the market. *Am. Econ. Rev.* 76, 728-741.

Kahneman, D., Knetsch, J.L., Thaler, R.H., 1986b. Fairness and the assumptions of economic theory. *J. Polit. Econ.* 94 (5), S102-S132. <http://www.jstor.org/stable/296367>

- Ma, Q., Hu, Y., Jiang, S., Meng, L., 2015. The undermining effect of facial attractiveness on brain responses to fairness in the Ultimatum Game: an ERP study. *Neurosci Lett*. 589, 77.
- Malle, B.F., Guglielmo, S., Monroe, A.E., 2014. A theory of blame. *Psychol. Bull.* 140, 216-237.
- Maréchal, M.A., Cohn, A., Ugazio, G., Ruff, C.C., 2017. Increasing honesty in humans with noninvasive brain stimulation. *PLoS One* 12, e0171363. U. S. National Library of Medicine.
- Masserman, J.H., Wechkin, S., Terris, W.A., 1964. Altruistic behavior in rhesus monkeys. *Am. J. Psychiatr.* 121, 584-588.
- McAuliffe, K., Blake, P.R., Steinbeis, N., Warneken, F., 2017. The developmental foundations of human fairness. *Nat. Human Behav.* 1 (2), 0042.
- McCullough, M.E., Kilpatrick, S.D., Emmons, R.A., Larson, D.B., 2001. Is gratitude a moral affect? *Psychol. Bull.* 127 (2), 249.
- Moll, J., De Oliveira-Souza, R., Zahn, R., 2008. The neural basis of moral cognition: sentiments, concepts, and values. *Annu. Rev. Psychol.* 59, 421-448.
- Moll, J., Krueger, F., Zahn, R., Pardini, M., de Oliveira-Souza, R., Grafman, J., 2006. Human mirror systems guide decisions about charitable donation. *Proc. Natl. Acad. Sci. U. S. A.* 103, 10348-10352.

- van Baar, J.M., Chang, L.J., Sanfey, A.G., 2019. The computational and neural substrates of moral strategies in social decisions. *Nat. Commun.* 10, 1416.
- van Den Bos, W., van Dijk, E., Westenberg, M., Rombouts, S.A., Crone, E.A., 2009. What motivates repayment? Neural correlates of reciprocity in the Trust Game. *Soc. Cogn. Affect Neurosci.* 4 (3), 294.
- van den Bos, W., van Dijk, E., Westenberg, M., Rombouts, S.A., Crone, E.A., 2011. Changing brains, changing perspectives: the reciprocity dilemma. *Psychol. Sci.* 22 (1), 670.
- Volz, K.G., Vogeley, K., Tittgemeyer, M., von Cramon, D.Y., Sutter, M., 2015. The neural basis of deception in strategic interactions. *Front. Behav. Neurosci.* 9, 27.
- Watanabe, T., Takezawa, M., Nakawake, Y., Kunimatsu, A., Yamasue, H., Nakamura, M., et al., 2014. Two distinct neural mechanisms underlie indirect reciprocity. *Proc. Natl. Acad. Sci. U. S. A.* 111 (11), 3995.
- Wojciszke, B., Parzuchowski, M., Bocian, K., 2015. Moral judgments and impressions. *Curr. Opin. Psychol.* 6, 50.
- Woo, C.-W., Koban, L., Kross, E., Lindquist, M.A., Banich, M.T., Ruzic, L., Andrews-Hanna, J.R., Wager, T.D., 2014. Separate neural systems for social rejection. *Nat. Commun.* 5, 5380.
- Wrangham, R.W., 2018. Two types of aggression in human evolution. *Proc. Natl. Acad. Sci. U. S. A.* 115 (2), 245.
- Wu, Y., Leliveld, M.C., Zhou, X., 2011. Social distance modulates recipient