

Journal Club

Editor's Note: These short, critical reviews of recent papers in the Journal, written exclusively by graduate students or postdoctoral fellows, are intended to summarize the important findings of the paper and provide additional insight and commentary. For more information on the format and purpose of the Journal Club, please see http://www.jneurosci.org/misc/ifa_features.shtml

Dissociating Guilt- and Inequity-Aversion in Cooperation and Norm Compliance

Hongbo Yu,¹ Bo Shen,¹ Yunlu Yin,¹ Philip R. Blume,² and Luke J. Chang¹

¹Center for Brain and Cognitive Sciences and Department of Psychology, Peking University, Beijing 100871, China, and ²Department of Psychology and Neuroscience, University of Colorado, Boulder, Colorado 80309

Social norms provide a set of expectations regarding context-specific appropriate behavior that aids in navigating social environments. Classic studies have demonstrated that people tend to conform to these norms even at the cost of their own interests and Henrich et al., 2001) and there are likely differing motivations for individuals to comply with these norms. For example, one motivation, consequentialism, emphasizes the outcome of an action as the sole measure of its moral worth (Mill, 1861/1998). From this philosophical perspective, one may avoid violating social norms simply because unfair and inequitable outcomes are bad for the greater good (e.g., distributive preferences). Alternatively, according to sentimentalism (Smith, 1759/2002), empathy with others "constitutes the moral approval. . . for agents and/or their actions" (Slote, 2010). This framework argues that people are motivated to comply with norms to avoid suffering from harm-ing another as a result of violating the trustee with multiple anonymous investors while undergoing fMRI. For each trial, trustees were given information about the investor's expectation and also the payoffs with their relative weights varying across individuals and contexts. Unfortunately, the majority of the research that uses social decision-making has been unable to effectively dissociate these two distinct motivations. This is likely a consequence of a peculiar convention in bargaining experiments to neither measure nor manipulate individuals' expectations. Thus, it has been unclear how much participants are motivated by distributive preferences (i.e., inequity-aversion) compared with disappointing a relationship partner (i.e., guilt-aversion). Fortunately, there has recently been a growing trend to both measure (Chang et al., 2011; Chang and Sanfey, 2013) and manipulate (Xiang et al., 2013) agents' expectations.

In a recent study published in *The Journal of Neuroscience* (Honsugi et al., 2015), we provided an important theoretical advance to dissociate the inequity- and guilt-aversion motivations in human norm compliance and identify the brain bases for each motivation. The experimenters used a modified trust game (Charness and Dufwenberg, 2006) in which participants initially decided as an investor whether or not to invest their endowment with an anonymous trustee and reported their belief about the likelihood of the trustee reciprocating. Participants then played the role of the trustee with multiple anonymous investors while undergoing fMRI. For each trial, trustees were given information about the investor's expectation and also the payoffs with their relative weights varying across individuals and contexts. Unfortunately, the majority of the research that uses social decision-making has been unable to effectively dissociate these two distinct motivations. This is likely a consequence of a peculiar convention in bargaining experiments to neither measure nor manipulate individuals' expectations. Thus, it has been unclear how much participants are motivated by distributive preferences (i.e., inequity-aversion) compared with disappointing a relationship partner (i.e., guilt-aversion). Fortunately, there has recently been a growing trend to both measure (Chang et al., 2011; Chang and Sanfey, 2013) and manipulate (Xiang et al., 2013) agents' expectations.

The basic framework for how these motivations were modeled was based on expected utility theory, which assumes that participants make decisions that maximize their expected payoff. Here, payoffs could be material (based on the amount of money each trustee receives) or psychological (based on concern for the investor's welfare) (Fehr and Camerer, 2007). The actually compared psychological payoffs arising from inequitable distributional outcomes (i.e., the absolute differ-

Received March 29, 2015; revised April 30, 2015; accepted May 5, 2015.

H.Y., B.S., Y.Y., and P.R.B. are supported by grants awarded to P.R.B. (NS055151) and L.J.C. (NS077011). We thank Professor Xiaolin Zhou for his helpful comments on our manuscript. We are also very grateful to Mr. Yin Wu and Mr. Shaofeng Yan without whose support this article would not have been possible.

Correspondence should be addressed to Hongbo Yu, Department of Psychology, Peking University, Beijing 100871, China. E-mail: hbyu101325@pku.edu.cn or hbyu101325@gmail.com.

DOI:10.1523/JNEUROSCI.1225-15.2015

Copyright © 2015 the authors 0270-6474/15/358973-03\$15.00/0

ence between the two players' payoffs) (Fehr and Schmidt, 1999) and feelings of guilt, which arose from disappointing a relationship partner by making a decision that resulted in the investor receiving a smaller payoff than he/she expected (i.e., the amount of money that the investor would have received had the trustee chosen to cooperate multiplied by the investor's estimated probability of the trustee's cooperation) (Battigalli and Dufwenberg, 2007). It is important to note that trustees had full information about the investor's expectations and each player's payoffs and thus their motivations can be inferred by how much they considered inequity or disappointing the investor when making their decision. A critical aspect of the experimental design was that the payoff matrix was constructed in such a way that the trustees' payoffs were uncorrelated with the amount of inequity between their partner's payoff, and both were uncorrelated with the investors' expectations about the likelihood the trustee would choose Cooperate. This allowed the experimenters to extend previous work (Chang et al., 2011) and disentangle these two otherwise intertwined motivations underlying human cooperation and norm compliance.

The authors found that the two motivations were associated with different neural circuitry. Controlling for guilt, inequity was positively associated with activation in the ventral striatum and amygdala. While other studies have implicated the ventral striatum in tracking inequity, it appears to go in the opposite direction, such that there is greater ventral striatal and amygdala activation associated with decreasing inequity (Fabiani et al., 2008; Tricomi et al., 2010). There are several possible reasons that can account for these discrepancies. First, these studies differed substantially in their design. In this study, the participants made decisions based on the inequity of the payoffs, while participants in the Tricomi et al. (2010)

consequentialism and sentimentalism considerations independently affect norm compliance and cooperation. Moreover, these motivations appear to be encoded in separate brain circuits. We believe that combining formal mathematical modeling, neuroscientific techniques, and social psychological theories will continue to further our insight into the material basis of our social nature.

References

- Battigalli P, Dufwenberg M (2007) Guilt in games. *Am Econ Rev* 97:170–176. [CrossRef](#)
- Bicchieri C (2006) *The grammar of society: the nature and dynamics of social norms*. Cambridge: Cambridge UP.
- Chang LJ, Sanfey AG (2013) Great expectations: neural computations underlying the use of social norms in decision-making. *Soc Cogn Affect Neurosci* 8:277–284. [CrossRef](#) [Medline](#)
- Chang LJ, Smith A, Dufwenberg M, Sanfey AG (2011) Triangulating the neural, psychological, and economic bases of guilt aversion. *Neuron* 70:560–572. [CrossRef](#) [Medline](#)
- Charness G, Dufwenberg M (2006) Promises and partnership. *Econometrica* 74:1579–1601. [CrossRef](#)
- Fehr E, Camerer CF (2007) Social neuroeconomics: the neural circuitry of social preferences. *Trends Cogn Sci* 11:419–427. [CrossRef](#) [Medline](#)
- Fehr E, Fischbacher U (2004) Third-party punishment and social norms. *Evol Hum Behav* 25:63–87. [CrossRef](#)
- Fehr E, Schmidt KM (1999) A theory of fairness, competition, and cooperation. *Q J Econ* 114:817–868. [CrossRef](#)
- Henrich J, Boyd R, Bowles S, Camerer C, Fehr E, Gintis H, McElreath R (2001) In search of Homo economicus: behavioral experiments in 15 small-scale societies. *Am Econ Rev* 91:73–78. [CrossRef](#)
- Knoch D, Schneider F, Schunk D, Hohmann M, Fehr E (2009) Disrupting the prefrontal cortex diminishes the human ability to build a good reputation. *Proc Natl Acad Sci U S A* 106:20895–20899. [CrossRef](#) [Medline](#)
- Koban L, Corradi-Dell'Acqua C, Vuilleumier P (2013) Integration of error agency and representation of others' pain in the anterior insula. *J Cogn Neurosci* 25:258–272. [CrossRef](#) [Medline](#)
- Mill JS (1861/1998) *Utilitarianism* (Crisp R, ed). New York: Oxford UP.
- Nihonsugi T, Ihara A, Haruno M (2015) Selective increase of intention-based economic decisions by noninvasive brain stimulation to the dorsolateral prefrontal cortex. *J Neurosci* 35:3412–3419. [CrossRef](#) [Medline](#)
- Ruff CC, Ugazio G, Fehr E (2013) Changing social norm compliance with noninvasive brain stimulation. *Science* 342:482–484. [CrossRef](#) [Medline](#)
- Sanfey AG, Stallen M, Chang LJ (2014) Norms and expectations in social decision-making. *Trends Cogn Sci* 18:172–174. [CrossRef](#) [Medline](#)
- Slote MA (2010) *Moral sentimentalism*. New York: Oxford UP.
- Smith A (1759/2002) *The theory of moral sentiments* (Haakonssen K, ed). Cambridge: Cambridge UP.
- Stagg CJ, Best JG, Stephenson MC, O'Shea J, Wylezinska M, Kincses ZT, Morris PG, Matthews PM, Johansen-Berg H (2009) Polarity-sensitive modulation of cortical neurotransmitters by transcranial stimulation. *J Neurosci* 29:5202–5206. [CrossRef](#) [Medline](#)
- Tabibnia G, Satpute AB, Lieberman MD (2008) The sunny side of fairness preference for fairness activates reward circuitry (and disregarding unfairness activates self-control circuitry). *Psychol Sci* 19:339–347. [CrossRef](#) [Medline](#)
- Tricomi E, Rangel A, Camerer CF, O'Doherty JP (2010) Neural evidence for inequality-averse social preferences. *Nature* 463:1089–1091. [CrossRef](#) [Medline](#)
- Xiang T, Lohrenz T, Montague PR (2013) Computational substrates of norms and their violations during social exchange. *J Neurosci* 33:1099–1108. [CrossRef](#) [Medline](#)
- Yoshida W, Dolan RJ, Friston KJ (2008) Game theory of mind. *PLoS Comput Biol* 4:e1000254. [CrossRef](#) [Medline](#)
- Yu H, Hu J, Hu L, Zhou X (2014) The voice of conscience: neural bases of interpersonal guilt and compensation. *Soc Cogn Affect Neurosci* 9:1150–1158. [CrossRef](#) [Medline](#)