

Toward a Brain-Based Bio-Marker of Guilt

H. b Y^{U1}, L. K ba ^{2,3}, M. J. C. ⁴, Xa Z ^{U5,6,7,8,9} a. T. D. Wa . ^{2,10}

¹Department of Psychological and Brain Sciences, University of California Santa Barbara, Santa Barbara, California, United States of America. ²Institute of Cognitive Science, University of Colorado, Boulder, CO, USA. ³Department of Psychology and Neuroscience, University of Colorado, Boulder, CO, USA. ⁴Department of Psychology, Yale University, New Haven, Connecticut, United States of America. ⁵School of Psychological and Cognitive Sciences, Peking University, Beijing, China. ⁶Beijing Key Laboratory of Behavior and Mental Health, Peking University, Beijing, China. ⁷PKU-IDG/McGovern Institute for Brain Research, Peking University, Beijing, China. ⁸Institute of Psychological and Brain Sciences, Zhejiang Normal University, Zhejiang, China. ⁹Key Laboratory of Applied Brain and Cognitive Sciences, School of Business and Management, Shanghai International Studies University, Shanghai, China. ¹⁰Department of Psychological and Brain Sciences, Dartmouth College, Hanover, New Hampshire, United States of America.

Neuroscience Insights Volume 15: 1–3 © The Author(s) 2020 Article reuse guidelines: sagepub.com/journals-permissions DOI: 10.1177/2633105520957638

\$SAGE

ABSTRACT: Guilt is a quintessential emotion in interpersonal interactions and moral cognition. Detecting the presence and measuring the intensity of guilt-related neurocognitive processes is crucial to understanding the mechanisms of social and moral phenomena. Existing neuroscience research on guilt has been focused on the neural correlates of guilt states induced by various types of stimuli. While valuable in their own right, these studies have not provided a sensitive and specific bio-marker of guilt suitable for use as an indicator of guilt-related neurocognitive processes in novel experimental settings. In a recent study, we identified a distributed Guilt-Related Brain Signature (GRBS) based on 2 independent functional MRI datasets. We demonstrated the sensitivity of GRBS in detecting a critical cognitive antecedent of guilt, namely one's responsibility in causing harm to another person, across participant populations from 2 distinct cultures (ie, Chinese and Swiss). We also showed that the sensitivity of GRBS did not generalize to other types of negative affective states (eg, physical and vicarious pain). In this commentary, we discuss the relevance of guilt in the broader scope of social and moral phenomena, and discuss how guilt-related biomarkers can be useful in understanding their psychological and neurocognitive mechanisms underlying these phenomena.

KEYWORDS: Guilt, biomarker, function MRI, social cognition, morality

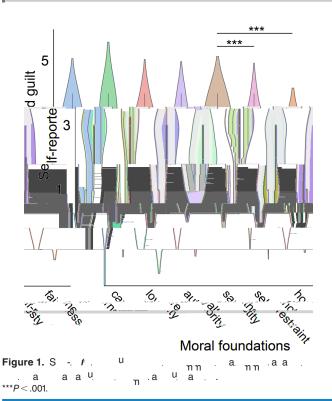
Guilt as a multifaced concept

Guilt, like many other social emotions, is a multifaced psychological construct and is often used equivocally in everyday life. Hurting an innocent person is a paradigmatic scenario in which people feel and express guilt.^{1,2} However, even in this case, we may not be dealing with one single kind of guilt—it is an open question whether an initial intention to harm influences the quality and magnitude of the guilt an agent later experiences.3 When we shift our focus to non-social use of the term "guilt," we will see even more diversity and complexity.⁴ For example, guilt appeal has been used as an advertising strategy for healthy diets. Some snack brands, instead of using label such as "reduced fat" or "reduced calories" for high-fat, highcalorie food products, directly label them as "reduced guilt" in order to ease customers' worries about the healthfulness of those products.⁵ We feel and express guilt when we fail to live up to our personal goals that are not directly related to other individuals or moral norms, such as keeping a healthy diet, working hard for an exam, and physical exercise. Indeed, people report experiencing guilt in their everyday life over almost all the domains of moral violations proposed in the Moral

Foundations Theory,^{6,7} including harm, unfairness, disloyalty, subversion, degradation, dishonesty, and lack of self-restraint. In fact, violation of self-restraint elicits stronger guilty feelings (on a 5-point Likert scale) than violation of fairness (mean difference = 1.40, SE = 0.21, z-ratio = 6.69, P<.001) and violation of honesty principles (mean difference = 0.87, SE = 0.17, z-ratio = 5.22, P<.001) (Figure 1), according to a large-scale experience sampling survey.⁷

What strategies should we take to investigate the neurocognitive basis of guilt in the face of its conceptual complexity? An analogy with pain, another multifaceted concept, may be useful to illustrate the approach we are introducing here. In social neuroscience, it has been debated whether physical pain, bodily sensation induced by external nociceptive stimuli, and social "pain" – psychological anguish elicited by social isolation, rejection or empathy – share the same neurocognitive basis. Studies examining the blood-oxygen-level-dependent (BOLD) signals that correlates with each phenomenon have consistently shown overlapping brain areas elicited by physical pain and social "pain." However, the relatively low spatial resolution of BOLD signal hinders the inference from overlapping activations to

2 Neuroscience Insights



overlapping neural representations or psychological constructs.⁹ An alternative approach is to develop a multivariate brain-based signature (or bio-marker) of each construct. The idea here is that if the bio-marker of physical pain does not respond to social "pain" and vice versa, then these 2 constructs do not share the same neural representation.¹⁰

Developing a guilt-related brain signature

Inspired by this approach, we recently identified a multivariate brain-based signature of guilt based on a paradigmatic case of guilt—causing harm to an innocent person.¹¹ We trained and validated the signature on 2 fMRI datasets. In the training dataset (N = 24, Chinese population), participants and an anonymous co-player performed a perceptual task, where failure would cause pain to the co-player.¹² We induced guilt by manipulating the responsibility of the participants in causing the pain. Specifically, if a participant performed poorly and the co-player performed well, then the performance failure, and the resulting co-player's pain, was caused by the participant. In comparison, if both the participants and the co-player performed poorly, then both of them were responsible for the coplayer's pain. Behaviorally, both self-reported guilt and reparation were positively correlated with participants' responsibility. We trained a multivariate Support-Vector-Machine (SVM) classifier to dissociate the sole-responsible and the both-responsible conditions. This classifier, or guilt-related brain signature (GRBS), was not only able to discriminate the 2 conditions on which it was trained, it was also able to discriminate the sole-responsible condition with other closely matched control conditions in the training dataset. Moreover, the predictive power of GRBS was generalizable to an

independent test dataset (N=19; Swiss population) that adopted a similar interpersonal action-monitoring task.¹³ We further demonstrated that GRBS did not discriminate different levels of painful thermal stimulation or different degree of vicarious pain, nor did it differentiate recalled guilt from recalled sadness or shame induced by person-specific episodes. Together, these results demonstrate that GRBS satisfies the 3 criteria proposed for bio-markers: sensitivity, specificity, and generalizability.¹⁴ Specifically, GRBS: (1) detects the presence of the "cognitive antecedents" of guilt in social interactions, here operationalized as responsibility; (2) does not discriminate other types of negative experiences, such as physical pain and emotion memory; and (3) detects the presence of the cognitive antecedent on which the signature is trained are present in studies and samples other than the training sample.

Using GRBS as an indicator of guilt-related neurocognitive processes in social-moral decision-making

Guilt-related neurocognitive processes are involved in many social-moral decision-making contexts. However, agents in those contexts are not always aware of or have biases in reporting guilt and guilt-related processes. In these situations, GRBS has the potential to provide an implicit, brain-based measure of guilt-related neurocognitive processes that are not easily forged by the agents. Returning to the self-restraint failure example, one theoretically important question is when people claim that they feel guilty for eating too much or for not working hard enough, are the neurocognitive processes underlying this affective phenomenon the same one as when they feel guilty for hurting their partner? If GRBS could discriminate self-restraint failure from self-restraint success, then we would be more confident that we are talking about the same *kind* of emotion in interpersonal and intrapersonal scenarios.⁴

Guilt is also relevant to moral evaluations of actions and character. When evaluating the moral status of an action or the moral character of an agent, the agent's inner states, such as attention and emotion accompanying their action, usually play an integral role. ¹⁵ Take hypocrisy as an example. A commonly held conceptualization characterizes a hypocrite as someone whose behaviors fall short of their expressed attitudes regarding some moral standards, namely "saying one thing, doing another." ¹⁶ Note, however, that this conceptualization speaks only to observable behaviors, irrespective of the mental states of the agent who behaves this way. Some philosophers, however, make the distinction between deceptive and *akratic* hypocrite. ¹⁷

Deceptive hypocrites "appears moral while, if possible, avoiding the cost of actually being moral." These hypocrites do not genuinely care about the moral standards that they publicly preach or cite to blame others, and therefore deserve the moral objections that laypeople assign to hypocrites. Akratic hypocrites, on the other hand, do genuinely care about the moral standards that they preach, but occasionally give in to temptations at the time of decision-making, perhaps due to weakness-of-the-will. A hallmark of

Yu et al 3

akratic hypocrites, therefore, is their feelings of conflict and guilt when they realize that what they do violates the moral standards they genuinely believe to be relevant and valuable. ¹⁷ Judging and treating deceptive and akratic hypocrites differently according to their mental states (ie, moral conflict, guilt) seems fairer and leaves room for moral education and self-improvement. ¹⁹ Behavioral measures alone are difficult, if not impossible, to distinguish these 2 types of hypocrites, because self-reported conflicted feelings and guilt can be easily faked. Applying the GRBS to neural response patterns associated with moral decision-making may offer a way to gauge the guilt-related neurocognitive processes involved and therefore provides a way to characterize the extent of deceptive versus akratic hypocrisy. ²⁰

Understanding the diversity and complexity of guilt via the brain-based signature approach

There are some limitations to GRBS that are worth noting. First, GRBS was trained on the datasets where the experimen-