

Disentangling consistency and reward in repeated moral decisions with mouse tracking and fMRI

Xinyi Julia Xu (yc27301@um.edu.mo)

Centre for Cognitive and Brain Sciences and Department of Psychology, University of Macau
Macau SAR, China

Haiyan Wu (haiyanwu@um.edu.mo)

Centre for Cognitive and Brain Sciences and Department of Psychology, University of Macau
Macau SAR, China

Abstract

Tracking response history and current rewards is critical for making moral decisions. By integrating fMRI and mouse tracking (MT) with a value-based moral decision task, we quantify the level of choice conflict with the MT metric, and examine how individuals incorporate information from the response history to make repeated moral decisions. Our study uses response entropy and cumulative responses (CR) to define choice consistency on both the subject-level and trial-level. We find that a stronger correlation between choice conflict and response entropy is mediated by the weight of reward in decisions. On the neural level, the brain adapts to conflict over the experiment sessions, and the adaption in reward-related brain regions is linked to response entropy. Meanwhile, multivariate representations in cognitive control and self-referential brain regions encode the weight of relative reward and CR. Through understanding choice conflict and response history, our research sheds light on its significance in multi-trial moral decision-making from the consistency perspective. These findings lay the groundwork for studying the underlying mechanisms in repeated decision processes.

Keywords: choice consistency, dishonesty, fMRI, moral decision

Introduction

Repeated decisions are common in daily life, where we often adapt to different contexts while also maintaining choice consistency—the tendency to repeat past choices. Choice consistency refers to repeating former choices, which is quantified by **cumulative responses (CR)** - the time of choosing the same option across all repetitions (Alós-Ferrer & Garagnani, 2021; Sen, 1993). In moral decision-making, maintaining consistency is crucial to ensure that our decisions align with our values, beliefs, or self-image (Jefferson, 2020). Although previous studies presented repeated moral scenarios (Garrett, Lazzaro, Ariely, & Sharot, 2016), the role of self-consistency remains unclear.

Moral decisions engage the reward, self-referential, and cognitive control networks (Speer, Smidts, & Boksem, 2022). Dishonest behavior involves weighing external rewards against moral costs (Allingham & Sandmo, 1972; Becker, 1968), with cognitive control required to override moral defaults (Speer, Smidts, & Boksem, 2020). Brain regions

like the ACC, IFG, and NAcc are more active in deceivers, while honest individuals show greater activation in the self-referential network (PCC, TPJ, MPFC).

However, the neurocomputational basis of consistency in moral choices is still unclear. Prior work focused on trial-by-trial effects (e.g., switching costs) (Luu & Stocker, 2018; Weber et al., 2023) and rarely quantified choice history explicitly (Luu & Stocker, 2018). In this study, we address these gaps by using a computer mouse in an fMRI setting, where participants repeatedly made moral decisions.

Methods

Task Procedure

The experiment included nine self-paced runs, each with 20 randomized questions. In each trial (Fig. 1a), participants first saw a question and a start button. Upon clicking start, two answers (correct/incorrect) appeared in the screen corners, along with information: correctness (black circle), monetary reward, and past choice frequency (red triangles, visible from run 2 onward). Incorrect answers offered higher rewards in over 50% of trials to create a moral-reward conflict. Participants had 4 seconds to respond, and feedback was shown for 1 second. Trials with no response were excluded. Cumulative rewards were shown after each session, and mouse position reset at the bottom center each trial.

Computation of response entropy

Entropy (Shannon, 1948) is one way to quantify the randomness of a system, which is adopted to quantify the choice consistency.

$$H(p) = -p \log_2 p - (1 - p) \log_2 (1 - p); \quad (1)$$

Mouse Trajectory Measurement

Mouse trajectories underwent standard spatial and temporal normalization (Freeman & Ambady, 2010; Xu, Liu, Hu, & Wu, 2021). We computed the area under the curve (AUC)—the geometric area between the actual and ideal paths—to quantify response conflict in decision-making (Stillman, Krajbich, & Ferguson, 2020).

fMRI general linear model (GLM) analysis

For each run and participant, a GLM was constructed with stimulus onset as the event time and trial-wise AUC as a para-

metric modulator of BOLD signal. Parametric modulator beta maps were used for inter-subject representational similarity analysis (IS-RSA; see below).

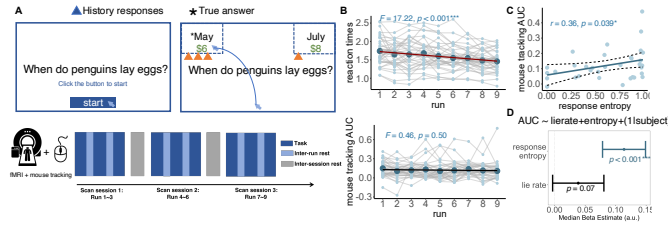


Figure 1: Illustration of the experimental paradigm and behavioral results.

Results

Choice conflict is correlated with response entropy

Reaction times (RTs) decreased over runs ($F_{(1,304)} = 17.22$, $p < 0.001$), while the mouse-tracking area under the curve (AUC)—a measure of response conflict—remained stable ($F_{(1,304)} = 0.46$, $p = 0.50$; Figure 1B). AUC reflects hesitation and decisional conflict.

Higher entropy indicated more inconsistent choices, regardless of honesty. AUC was positively correlated with response entropy ($r = 0.36$, $p = 0.039$; Figure 1C), suggesting that greater decisional conflict accompanied more random behavior. A linear mixed model further confirmed that response entropy—but not lie rate—significantly predicted AUC (entropy: $\beta = 0.11$, $p < 0.001$; lie rate: $\beta = 0.04$, $p < 0.001$; Figure 1D).

Relative reward weight mediates the link between choice conflict and response entropy

Higher relative reward and more prior dishonest choices increased lie rates, shifting both threshold and overall dishonesty (Figure A). Dishonest responses were primarily driven by relative reward and relative CR, but not by the total number of prior dishonest responses (Figure B).

Next, we applied a Bayesian hierarchical drift-diffusion model (HDDM) (Wiecki, Sofer, & Frank, 2013) to examine the effect of relative reward and CR. In the best model, both relative reward and CR significantly influenced the drift rate (posterior probability = 1), with CR having a stronger weight (posterior probability = 0.9965; Figure C). Notably, the weight on CR negatively correlated with response entropy, while the weight on reward positively correlated with both response entropy and lie rate (Figure D–E). Furthermore, the weight on relative reward significantly mediated the relationship between choice conflict (AUC) and response entropy (Figure F).

Patterns of choice conflict in cognitive control and self-referential ROIs represented both consistency and reward

We extracted the beta maps of AUC from cognitive-control related ROIs (see *Methods*). We observed a linear decrease in

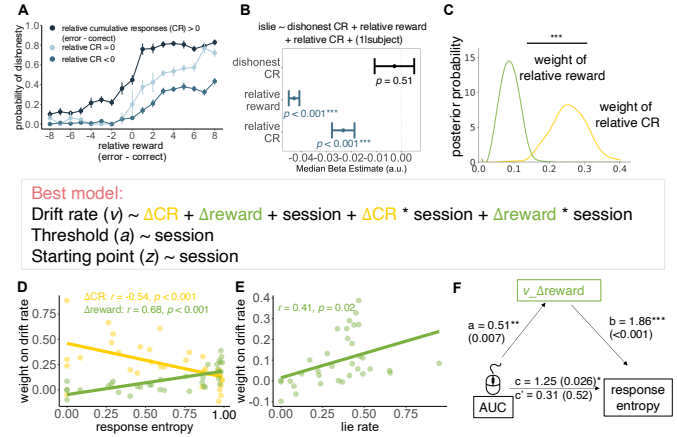


Figure 2: (A and B) The effect of reward and CR. (C) Results of HDDM. (D–F) Correlation of HDDM parameters and behavioral results.

IFG, preSMA, and vIPFC (Figure 3A), indicating the conflict adaption. Specifically, the mean activity in IFG was correlated with response entropy (Figure 3B).

We further conducted an inter-subject representational similarity analysis (ISRSA) to seek the brain regions encoding the evaluation of both the weight of relative CR and reward. We first created a geometric representational space of HDDM parameter space (weights of relative CR and reward). For brain data, we averaged the beta maps obtained from AUC modulation effect of all runs. We calculated multi-voxel activity pattern distances (1-correlation) between each pair of participants (Figure 3C). After correlating the behavioral distance matrix with the brain dissimilarity matrix, we observed significant inter-subject representational similarity effects in IFG ($r = 0.12$, $p = 0.0017$), ACC ($r = 0.09$, $p = 0.01$) and ITPJ ($r = 0.09$, $p = 0.02$; Figure 3D). These results indicated that conflict-related activity patterns in these regions were more similar than in other brain regions across participants who shared similar weights when making repeated moral choices.

Conclusion

We combined mouse-tracking and fMRI to study how people weigh current rewards and past choices during moral decision-making. Choice conflict, indexed by AUC, was associated with response entropy—not lie rate—and this relationship was mediated by sensitivity to reward. Participants relied more on cumulative response history than on immediate payoff. At the neural level, dishonesty-related regions encoded both consistency and reward signals, suggesting that moral decisions arise from dynamic integration of past behavior and current incentives across key brain networks.

Acknowledgments

This work was mainly supported by the Natural Science Foundation of China (U1736125), Science and Technology Development Fund (FDCT) of Macau [0127/2020/A3,

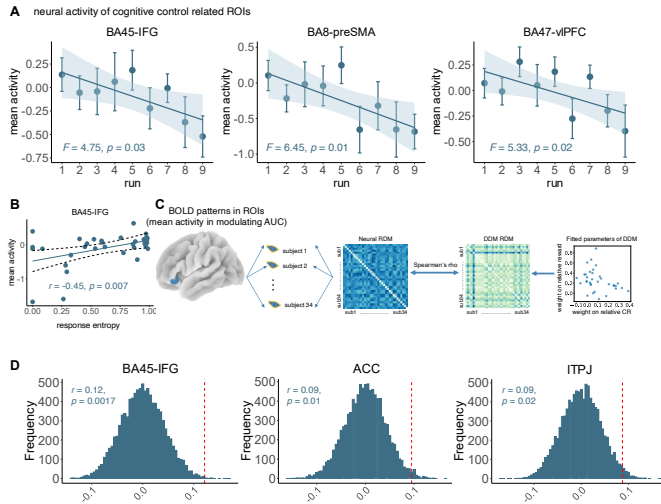


Figure 3: (A) Conflict adaption. (B) The mean modulation effect of AUC in IFG was correlated with response entropy (C) Illustration of inter-subject representational analysis (ISRSA). (D) The results of ISRSA

0041/2022/A], the Natural Science Foundation of Guangdong Province (2021A1515012509), MYRG of University of Macau (MYRG2022-00188-ICI), Shenzhen-Hong Kong-Macao Science and Technology Innovation Project (Category C) (SGDX2020110309280100), and the SRG of University of Macau (SRG2020-00027-ICI). We thank Prof. Xiaolin Zhou and Prof. Kang Lee for their helpful comments on the manuscript.

References

- Allingham, M. G., & Sandmo, A. (1972). Income tax evasion: A theoretical analysis. *Taxation: critical perspectives on the world economy*, 3, 323–338.
- Alós-Ferrer, C., & Garagnani, M. (2021). Choice consistency and strength of preference. *Economics Letters*, 198, 109672.
- Becker, G. S. (1968). Crime and punishment: An economic approach. In *The economic dimensions of crime* (pp. 13–68). Springer. doi: 10.1007/978-1-349-62853-7₂
- Freeman, J. B., & Ambady, N. (2010). Mousetracker: Software for studying real-time mental processing using a computer mouse-tracking method. *Behavior research methods*, 42(1), 226–241. doi: 10.3758/brm.42.1.226
- Garrett, N., Lazzaro, S. C., Ariely, D., & Sharot, T. (2016, October). The brain adapts to dishonesty. *Nature Neuroscience*, 19(12), 1727–1732. Retrieved from <https://doi.org/10.1038/nn.4426> doi: 10.1038/nn.4426
- Jefferson, A. (2020). Confabulation, rationalisation and morality. *Topoi*, 39(1), 219–227.
- Luu, L., & Stocker, A. A. (2018). Post-decision biases reveal a self-consistency principle in perceptual inference. *Elife*, 7, e33334. doi: 10.7554/elife.33334

- Sen, A. (1993). Internal consistency of choice. *Econometrica: Journal of the Econometric Society*, 495–521. doi: 10.2307/2951715
- Shannon, C. E. (1948). A mathematical theory of communication. *The Bell system technical journal*, 27(3), 379–423.
- Speer, S. P., Smidts, A., & Boksem, M. A. (2020). Cognitive control increases honesty in cheaters but cheating in those who are honest. *Proceedings of the National Academy of Sciences*, 117(32), 19080–19091. doi: 10.1073/pnas.2003480117
- Speer, S. P., Smidts, A., & Boksem, M. A. (2022). Cognitive control and dishonesty. *Trends in Cognitive Sciences*. doi: 10.1016/j.tics.2022.06.005
- Stillman, P. E., Krajbich, I., & Ferguson, M. J. (2020, November). Using dynamic monitoring of choices to predict and understand risk preferences. *Proceedings of the National Academy of Sciences*, 117(50), 31738–31747. Retrieved from <https://doi.org/10.1073/pnas.2010056117> doi: 10.1073/pnas.2010056117
- Weber, J., Iwama, G., Solbakk, A.-K., Blenkmann, A. O., Larsson, P. G., Ivanovic, J., ... Helfrich, R. (2023). Subspace partitioning in the human prefrontal cortex resolves cognitive interference. *Proceedings of the National Academy of Sciences*, 120(28), e2220523120.
- Wiecki, T. V., Sofer, I., & Frank, M. J. (2013). Hddm: Hierarchical bayesian estimation of the drift-diffusion model in python. *Frontiers in neuroinformatics*, 7, 55610.
- Xu, X. J., Liu, X., Hu, X., & Wu, H. (2021). Mt-aiat: Integrating mouse tracking into memory-detection aiat.