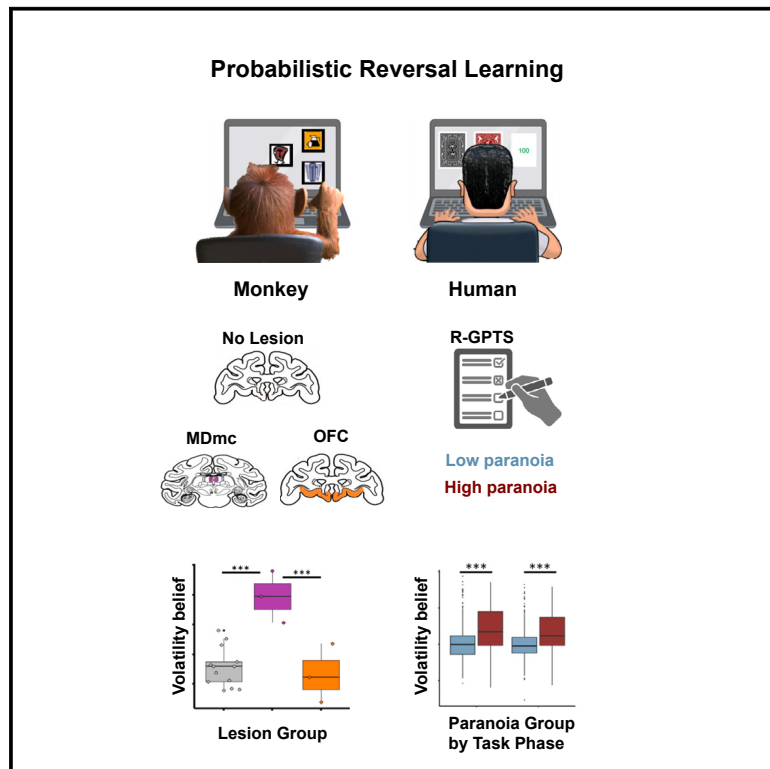


# Lesions to the mediodorsal thalamus, but not orbitofrontal cortex, enhance volatility beliefs linked to paranoia

## Graphical abstract



## Authors

Praveen Suthaharan,  
Summer L. Thompson,  
Rosa A. Rossi-Goldthorpe, ...,  
Jane R. Taylor, Philip R. Corlett,  
Steve W.C. Chang

## Correspondence

philip.corlett@yale.edu (P.R.C.),  
steve.chang@yale.edu (S.W.C.C.)

## In brief

Suthaharan et al. demonstrate a causal role of the primate mediodorsal thalamus (MD) in beliefs about environmental volatility. Applying a behavioral paradigm and a computational model, they establish that belief volatility increases in paranoid people and monkeys with lesions to the MD. This suggests a role for the MD in paranoia.

## Highlights

- Lesions to the mediodorsal thalamus (MD) in monkeys heighten belief volatility
- Orbitofrontal cortex (OFC) lesions impair value learning
- This is a double dissociation of belief learning between brain regions
- Shared computational modeling across species implicates MD in paranoia



## Report

# Lesions to the mediodorsal thalamus, but not orbitofrontal cortex, enhance volatility beliefs linked to paranoia

Praveen Suthaharan,<sup>1,2,3,18</sup> Summer L. Thompson,<sup>3,18</sup> Rosa A. Rossi-Goldthorpe,<sup>1,3</sup> Peter H. Rudebeck,<sup>4</sup> Mark E. Walton,<sup>5</sup> Subhojit Chakraborty,<sup>5,6</sup> Maryann P. Noonan,<sup>5,7</sup> Vincent D. Costa,<sup>8</sup> Elisabeth A. Murray,<sup>9</sup> Christoph D. Mathys,<sup>10,11</sup> Stephanie M. Groman,<sup>3,12,13</sup> Anna S. Mitchell,<sup>5,14</sup> Jane R. Taylor,<sup>1,3,15,16,17</sup> Philip R. Corlett,<sup>1,2,3,15,16,19,20,\*</sup> and Steve W.C. Chang<sup>1,2,15,16,17,19,20,21,\*</sup>

<sup>1</sup>Interdepartmental Neuroscience Program, Yale School of Medicine, New Haven, CT, USA

<sup>2</sup>Kavli Institute for Neuroscience, Yale University, New Haven, CT, USA

<sup>3</sup>Department of Psychiatry, Yale University, New Haven, CT, USA

<sup>4</sup>Icahn School of Medicine at Mount Sinai, New York, NY, USA

<sup>5</sup>Department of Experimental Psychology, Oxford University, Oxford, UK

<sup>6</sup>NIHR Biomedical Research Centre, Moorfields Eye Hospital NHS Foundation Trust and UCL Institute of Ophthalmology, London EC1V 9EL, UK

<sup>7</sup>Department of Psychology, University of York, York, UK

<sup>8</sup>Division of Neuroscience, Oregon National Primate Research Center, Oregon Health and Science University, Beaverton, OR, USA

<sup>9</sup>Laboratory of Neuropsychology, NIMH, Bethesda, MD, USA

<sup>10</sup>Interacting Minds Centre, Aarhus University, Aarhus, Denmark

<sup>11</sup>Translational Neuromodeling Unit (TNU), Institute for Biomedical Engineering, University of Zurich and ETH Zurich, Zurich, Switzerland

<sup>12</sup>Department of Neuroscience, University of Minnesota Medical School, Minneapolis, MN, USA

<sup>13</sup>Department of Anesthesia and Critical Care, University of Chicago, Chicago, IL, USA

<sup>14</sup>School of Psychology, Speech, and Hearing, University of Canterbury, Christchurch, New Zealand

<sup>15</sup>Department of Psychology, Yale University, New Haven, CT, USA

<sup>16</sup>Wu Tsai Institute, Yale University, New Haven, CT, USA

<sup>17</sup>Department of Neuroscience, Yale University, New Haven, CT, USA

<sup>18</sup>These authors contributed equally

<sup>19</sup>These authors contributed equally

<sup>20</sup>Senior author

<sup>21</sup>Lead contact

\*Correspondence: [philip.corlett@yale.edu](mailto:philip.corlett@yale.edu) (P.R.C.), [steve.chang@yale.edu](mailto:steve.chang@yale.edu) (S.W.C.C.)

<https://doi.org/10.1016/j.celrep.2024.114355>

## SUMMARY

**Beliefs—attitudes toward some state of the environment—guide action selection and should be robust to variability but sensitive to meaningful change. Beliefs about volatility (expectation of change) are associated with paranoia in humans, but the brain regions responsible for volatility beliefs remain unknown. The orbitofrontal cortex (OFC) is central to adaptive behavior, whereas the magnocellular mediodorsal thalamus (MDmc) is essential for arbitrating between perceptions and action policies. We assessed belief updating in a three-choice probabilistic reversal learning task following excitotoxic lesions of the MDmc ( $n = 3$ ) or OFC ( $n = 3$ ) and compared performance with that of unoperated monkeys ( $n = 14$ ). Computational analyses indicated a double dissociation: MDmc, but not OFC, lesions were associated with erratic switching behavior and heightened volatility belief (as in paranoia in humans), whereas OFC, but not MDmc, lesions were associated with increased lose-stay behavior and reward learning rates. Given the consilience across species and models, these results have implications for understanding paranoia.**

## INTRODUCTION

The ability to form and update beliefs about our actions and their consequences, especially in volatile environments, is at the core of advanced cognition.<sup>1</sup> Indeed, disruptions in volatility beliefs are associated with maladaptive cognitive and behavioral outcomes, such as paranoia, the belief that others intend to exert

harm.<sup>2</sup> Volatility beliefs and manifestations of their disruption, like paranoia, can be captured in decision-making paradigms that incorporate environmental volatility, such as probabilistic reversal learning (PRL) tasks.<sup>3,4</sup> In PRL tasks, contingencies between choices and outcomes reverse so that previously richly reinforced options become lean and vice versa. These tasks incentivize individuals to form a set of beliefs that are robust to



probabilistic noise but flexible enough to track true change in action-outcome contingencies to maximize reward (Figure 1A). In prior work with human participants, key behavioral metrics of belief connote flexible<sup>3</sup> and perseverative behavior.<sup>5</sup> “Win switching” occurs when the subject selects a different option on the subsequent trial following a reward, or “win,” which, when frequent, marks excessive flexibility. On the other hand, “lose staying” happens when the same choice is repeated following an unrewarded trial, or “lose,” signifying persistence, which becomes perseverative if sustained. Therefore, paranoia is associated with an increase in win switching, a dearth of lose staying, and an elevated volatility belief.<sup>3,4</sup>

We can discern such volatility beliefs quantitatively by fitting Bayesian inference (BI) models to participants’ behavior and estimating the parameter values to account for individual patterns of choices.<sup>3</sup> This computational modeling approach to study brain and behavior has been employed to better understand various aspects of human psychopathology, including paranoia. While traditional reinforcement learning models have provided insight into choice behavior in PRL tasks, there is evidence to suggest that BI models may be better equipped to capture rapid behavioral change in response to volatility (e.g., a reversal event).<sup>6,7</sup> Indeed, BI models seem necessary to capture the elevated win-switching behavior observed with paranoia.<sup>3</sup> It has been proposed that BI models might detect these nuances of behavior in PRL tasks by tracking latent (i.e., hidden) states, the unobservable features that dictate an environment’s underlying dynamics, which may allow more rapid recognition of state changes and more efficient behavioral adaptation.<sup>8,9</sup> Two components are critical in such a model: (1) participants’ assumptions about how the task works and (2) how those assumptions lead to their decisions. One such BI model that is designed to capture volatility beliefs is the hierarchical Gaussian filter (HGF; Figure 1B).<sup>6,10</sup> There are two model parameters that are particularly central to volatility processing: (1) an equilibrium value that attracts and stabilizes the agent’s estimate of volatility ( $m_3$ , volatility belief; higher values indicate greater expectation of change) and (2) the rate of adjustment of beliefs about the values of each option ( $\omega_2$ , value learning rate; higher values indicate more rapid learning). In recent work,  $m_3$  has been shown to relate to risk for psychosis.<sup>6</sup>

It is of course possible that other model structures and types may yield different insights to paranoia in humans and reversal behavior in animals. The exercise here is to use a computational model whose parameters relate to paranoia in humans to model behavior in nonhuman primates to connect those literatures. Other connections may become apparent in future work.

Several brain regions and circuits have been implicated in the formation and updating of such volatility processing<sup>11</sup> and/or paranoia in PRL tasks (reviewed in Soltani and Izquierdo<sup>12</sup>). fMRI, particularly in patients with paranoid persecutory delusional beliefs, has suggested that this volatility belief parameter ( $m_3$ ) correlates with the engagement of the dorsolateral prefrontal cortex (DLPFC) during PRL tasks.<sup>6</sup> Data in nonhuman primates suggest that lesions of the magnocellular mediodorsal thalamus (MDmc) may also recapitulate the associated win-switching pattern.<sup>13</sup> We note that MDmc receives inputs from the DLPFC, in addition to inputs from the ventral PFC,<sup>14,15</sup> which

may support the increased win-switching behavior after MDmc lesions. Animal work has shown that the thalamus, in partnership with the cortex, is a key nexus of perception and action<sup>16,17</sup> and a mediator of behavioral change in the face of evolving contingencies.<sup>18</sup> Yet, the MDmc has not been extensively implicated in human paranoia and thus warrants further investigation (but see Crail-Melendez et al.<sup>19</sup>). Aspiration lesions of the orbitofrontal cortex (OFC) in nonhuman primates have been shown to disrupt PRL performance,<sup>20,21</sup> though these effects may be more attributable to parts of the ventrolateral prefrontal cortex (VLPFC).<sup>22</sup> Resting-state functional connectivity studies suggest that the OFC also plays a role in human paranoia.<sup>23</sup> However, the HGF model has yet to be applied to choice behavior in PRL tasks performed by nonhuman primates, which could elucidate the neural underpinnings of volatility beliefs relevant to human paranoia.

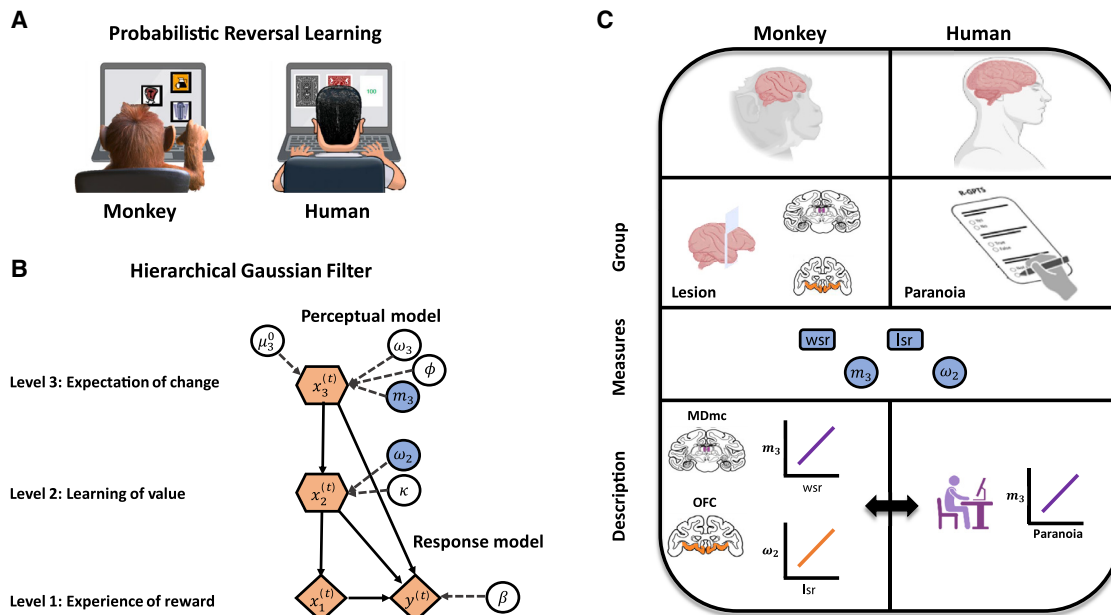
In the present work, we sought to infer the potential role of the MDmc and OFC in paranoia by estimating volatility belief in a PRL task with the HGF model in monkeys with lesions to these regions. We hypothesized that damage to brain regions involved in action and perception would lead to a pattern of changes in PRL behavior and volatility beliefs in monkeys that are associated with paranoia in human participants. By comparing the effects of MDmc and OFC lesions on actions that can be attributed to volatility belief, we can begin to elucidate the functional contributions of different brain regions, and possibly neural circuits, to belief updating and, transitively, to paranoia in human participants. Ultimately, the identification of human participants with similar patterns of behavior as lesioned monkeys, and examination of how those patterns relate to self-reported paranoia symptoms (Figure 1C), can afford us the ability to draw inferences about the specificity of differential relationships between task behavior, computational model parameters, brain lesions, and psychopathology.

## RESULTS

To understand the contributions of the MDmc and OFC to beliefs about volatility in the primate brain, we analyzed existing datasets from multiple labs that administered the same three-choice PRL task to rhesus macaques with selective, excitotoxic lesions of the MDmc or the OFC (STAR Methods). We focused particularly on win-switch and lose-stay behavior and beliefs about volatility using a BI model.

### Flexibility and perseveration following brain lesions

We examined win-switch and lose-stay rates in monkeys with excitotoxic lesions to either the MDmc ( $n = 3$ ) or OFC ( $n = 3$ ) compared with non-lesioned controls ( $n = 14$ ) (Figure 2A). We observed a significant two-way interaction between the lesion group (OFC, MDmc, or control) and reversal phase (before or after reversal) on win-switch behavior (Figure 2B, left;  $\chi^2 = 89.67$ ,  $p < 0.001$ ; generalized linear mixed model [GLMM]). To resolve this interaction, we first assessed group effects within each reversal phase. We identified an effect of lesion group on win switching before the reversal ( $\chi^2 = 61.17$ ,  $p < 0.001$ ) and after the reversal ( $\chi^2 = 580.6$ ,  $p < 0.001$ ). Pairwise comparisons within the pre-reversal phase revealed that, whereas MDmc-lesioned monkeys exhibited elevated win-switching behavior compared



**Figure 1. A translational mapping of belief updating in monkeys and human participants**

(A) A monkey and a human participant playing a three-choice PRL task, selecting from a set of three options and learning which is the *best* option, through trial and error, in an environment where the underlying reinforcement schedule changes periodically (i.e., reversals).

(B) A model used to investigate how beliefs about changes in the environment influence decision-making. The graphical notation is adopted from a prior study,<sup>6</sup> parameters of interest for this study are shaded in blue.

(C) Integrating data between monkeys with targeted brain lesions and human participants with self-reported paranoia through computational models to identify links between neural circuits and psychopathology. MDmc, magnocellular mediodorsal thalamus; OFC, orbitofrontal cortex; wsr, win-switch rate; lsr, lose-stay rate;  $m_3$ , volatility belief;  $\omega_2$ , value learning.

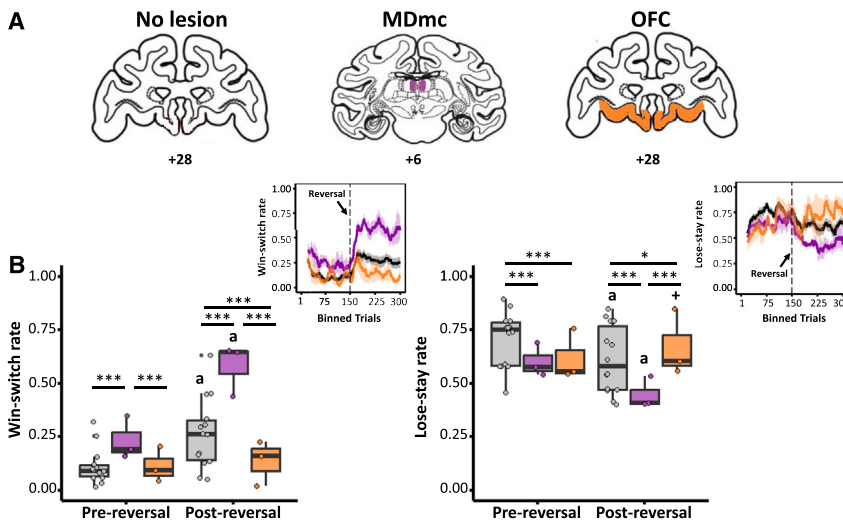
to non-lesioned controls ( $\chi^2 = 57.6, p < 0.001$ ) and OFC-lesioned monkeys ( $\chi^2 = 32.8, p < 0.001$ ), OFC-lesioned monkeys did not differ from non-lesioned controls ( $\chi^2 = 0.14, p = 0.712$ ). Pairwise comparisons within the post-reversal phase revealed differences between all groups; MDmc-lesioned monkeys again showed greater win-switching behavior compared to non-lesioned controls ( $\chi^2 = 410.22, p < 0.001$ ) and OFC-lesioned monkeys ( $\chi^2 = 486.63, p < 0.001$ ), and here we observed lower win switching in OFC-lesioned monkeys versus non-lesioned controls ( $\chi^2 = 96.06, p < 0.001$ ). We observed an effect of reversal phase on win switching in MDmc-lesioned monkeys ( $\chi^2 = 272.47, p < 0.001$ ) but not OFC-lesioned monkeys ( $\chi^2 = 1.58, p = 0.209$ ), which was characterized by a more dramatic increase in win switching (mean difference = 0.169,  $p < 0.001$ ) from pre to post reversal than in non-lesioned controls ( $\chi^2 = 324.19, p < 0.001$ ).

We also observed a similar two-way interaction between lesion group and reversal phase on lose-stay behavior (Figure 2B, right;  $\chi^2 = 41.45, p < 0.001$ ). An effect of lesion group on lose-stay behavior was observed both before reversal ( $\chi^2 = 78.13, p < 0.001$ ) and after reversal ( $\chi^2 = 96.78, p < 0.001$ ). Compared to non-lesioned controls, lose-stay rates were altered both pre and post reversal in MDmc-lesioned (pre reversal:  $\chi^2 = 56.9, p < 0.001$ ; post reversal:  $\chi^2 = 79.92, p < 0.001$ ) and OFC-lesioned monkeys (pre reversal:  $\chi^2 = 37.19, p < 0.001$ ; post reversal:  $\chi^2 = 6.17, p = 0.013$ ). Lose-stay behavior did not differ between the two lesion groups before the reversal ( $\chi^2 = 0.85, p = 0.356$ ), while

we did see differences in lose-stay behavior after the reversal ( $\chi^2 = 68.08, p < 0.001$ ). The effect of reversal within each group revealed a marked *decrease* in lose-stay behavior in both MDmc-lesioned monkeys ( $\chi^2 = 49.76, \beta = -0.58, p < 0.001$ ) and non-lesioned controls ( $\chi^2 = 115.75, \beta = -0.48, p < 0.001$ ), whereas a trend-level *increase* in lose-stay behavior was observed in OFC-lesioned monkeys ( $\chi^2 = 3.16, \beta = 0.18, p = 0.076$ ). These data, collectively, indicate that lesions to the MDmc were associated with erratic switching behaviors regardless of reversal phase. Lesions of the OFC, however, were associated with inflexible choice behavior that is reversal phase dependent.

### A double dissociation of volatility and value

To investigate the impact of brain lesions on how beliefs are updated under uncertainty, we applied the HGF model to trial-by-trial choice data from the MDmc lesion, OFC lesion, and non-lesion control monkeys (Figure S1). Computational modeling of behavior revealed distinct effects of lesion group on volatility belief and value learning (Figure 3). First, reversal enhanced volatility beliefs in MDmc-lesioned ( $m_3; \chi^2 = 7736.09, p < 0.001$ ; Figure 3A) and non-lesioned control monkeys ( $\chi^2 = 37.9, p < 0.001$ ; Figure 3A), whereas lesions to the OFC blocked the increase in volatility beliefs associated with reversal ( $\chi^2 = 0.38, p = 0.535$ ; Figure 3A). Furthermore, we observed a marginally greater increase in volatility beliefs in MDmc-lesioned compared to non-lesioned control monkeys (mean difference = 1.33,  $p = 0.081$ ).



**Figure 2. Behavior of lesioned monkeys**

(A) Illustrations of the anatomical brain lesion locations in monkeys, where the excitotoxic MDmc and OFC lesion locations are illustrated on standard coronal sections from a monkey brain atlas (Laboratory of Neuropsychology, National Institute of Mental Health). Prior studies show the actual lesion locations of the MDmc and OFC in individual monkeys.<sup>13,22</sup>

(B) Differences in win-switch behavior (left) and lose-stay behavior (right) for the lesion groups between pre-reversal and post-reversal phases. Top right: PRL behaviors over time (moving average across 20 lagged trials) of win-switch (left) and lose-stay (right) choices across lesion groups, with a reversal occurring after 150 trials.

Each data point overlaid indicates an individual monkey's data. Asterisks indicate significant effects of lesion group within each reversal phase (\*\* $p < 0.001$ , \* $p < 0.05$ ; GLMM). Other symbols indicate significant effects of reversal phase within each lesion group (a,  $p < 0.001$ ; +,  $p < 0.10$ ).

With regard to learning about value, monkeys with OFC lesions and non-lesioned controls exhibited increased value learning post reversal relative to pre reversal ( $\omega_2$ ; OFC:  $\chi^2 = 5.70$ ,  $p = 0.017$ ; non-lesion:  $\chi^2 = 21.50$ ,  $p < 0.001$ ; Figure 3B), whereas lesions in MDmc blocked the increase in value learning rate that accompanied reversal in the other groups ( $\chi^2 = 1.09$ ,  $p = 0.296$ ; Figure 3B).

Pairwise comparisons of effects of lesion group on volatility beliefs (Figure 3A) within the pre-reversal phase revealed that, whereas MDmc-lesioned monkeys exhibited elevated volatility belief compared to non-lesioned controls ( $\chi^2 = 3.99$ ,  $\beta = 0.99$ ,  $p = 0.046$ ), OFC-lesioned monkeys did not differ from non-lesioned controls ( $\chi^2 = 0.40$ ,  $p = 0.526$ ) or from MDmc-lesioned monkeys ( $\chi^2 = 0.47$ ,  $p = 0.495$ ). Pairwise comparisons within the post-reversal phase revealed differences between all groups; MDmc-lesioned monkeys showed enhanced volatility beliefs compared to non-lesioned controls ( $\chi^2 = 15.43$ ,  $\beta = 2.33$ ,  $p < 0.001$ ) and OFC-lesioned monkeys ( $\chi^2 = 61.42$ ,  $\beta = 3.72$ ,  $p < 0.001$ ), and here we observed a decrease in volatility beliefs in OFC-lesioned monkeys compared to non-lesioned controls ( $\chi^2 = 5.29$ ,  $\beta = -1.40$ ,  $p = 0.021$ ).

Pairwise comparisons of effects of lesion group on value learning (Figure 3B) within the pre-reversal phase revealed that, whereas OFC-lesioned monkeys exhibited an increase in value learning compared to non-lesioned controls ( $\chi^2 = 4.68$ ,  $\beta = 1.26$ ,  $p = 0.030$ ) and MDmc-lesioned monkeys ( $\chi^2 = 8.98$ ,  $\beta = 1.72$ ,  $p = 0.003$ ), MDmc-lesioned monkeys did not differ from non-lesioned controls ( $\chi^2 = 0.75$ ,  $p = 0.386$ ). Within the post-reversal phase, MDmc-lesioned monkeys exhibited reduced value learning relative to non-lesioned controls ( $\chi^2 = 18.74$ ,  $p < 0.001$ ) and OFC-lesioned monkeys ( $\chi^2 = 19.86$ ,  $p < 0.001$ ), whereas OFC-lesioned monkeys did not differ in value learning from non-lesioned controls ( $\chi^2 = 1.54$ ,  $p = 0.214$ ). This dissociation implies distinct neural and computational mechanisms underlying actions attributable to volatility and value beliefs; MDmc lesions disinhibited volatility beliefs but blunted sensitivity to the reversal for updating value learning, whereas OFC lesions disinhibited value learning but blunted

sensitivity to the reversal for updating volatility beliefs. Given that lesions in these areas lead to opposing effects on these belief-updating parameters, this established a double dissociation between the function of MDmc and OFC regions.

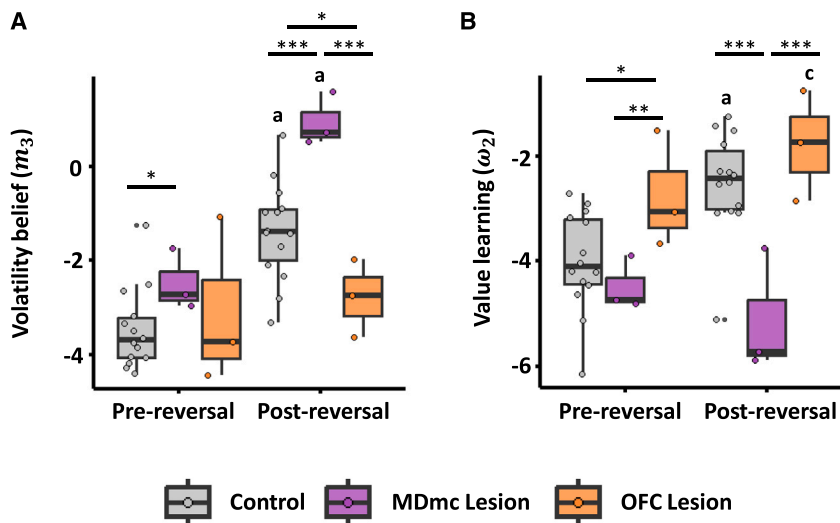
In summary, MDmc-lesioned monkeys demonstrated increased win-switching and reduced lose-stay behavior after reversal; OFC-lesioned monkeys showed the opposite behavior. In terms of belief updating, MDmc lesions increased volatility beliefs, notably post reversal, while OFC lesions blocked changes in volatility beliefs; however, OFC lesions elevated value learning, while MDmc lesions blocked changes in value learning. These findings confirm a double dissociation. See Figure S2 for more fine-grained, within-monkey analyses.

### Translational insights from monkey behavioral neuroscience to human psychopathology

In order to gain clinically relevant insights from the changes in beliefs about volatility and value learning parameters due to MDmc and OFC lesions, we compared PRL behavior of lesioned monkeys to that of humans with paranoia (Figure 4), as human paranoia has been associated with erratic switching behavior and an elevated sense of uncertainty about the environment.<sup>3</sup> We analyzed existing datasets in human participants who completed either of two types of PRL tasks: a single-reversal task whose structure matches the monkey version<sup>24</sup> and a multi-reversal task.<sup>3</sup> The shift-induced increase in unexpected volatility, in the multi-reversal task, has been shown to increase sensitivity in detecting paranoia group differences in task performance.<sup>3</sup>

Paranoid participants evinced higher win switching pre and post reversal, a pattern also observed in the monkeys with MDmc lesions, in both versions of the task (Figures 4A and 4B). In the single-reversal task, we found a significant interaction between paranoia group and reversal phase for win-switch behavior (Figure 4A;  $\chi^2 = 18.11$ ,  $p < 0.001$ , GLMM). In particular, we observed significantly more post-reversal win switching in paranoid compared to non-paranoid human participants ( $\chi^2 = 10.43$ ,  $p = 0.001$ ) but no pre-reversal, win-switching





**Figure 3. MDmc and OFC lesions differentially impact beliefs about volatility and value in monkeys**

(A) Differences in volatility beliefs ( $m_3$  parameter from the HGF model; higher values indicate greater volatility beliefs) for the lesion groups between the pre-reversal and post-reversal phases.

(B) Differences in reward value learning ( $\omega_2$  parameter; higher values indicate more rapid learning about value) for the lesion groups between pre-reversal and post-reversal phases.

Each data point overlaid indicates an individual monkey's data. Asterisks indicate significant lesion group differences (\*\* $p < 0.01$ , \* $p < 0.05$ ; GLMM). Other symbols indicate significant reversal phase differences in each lesion group (a,  $p < 0.001$ ; c,  $p < 0.05$ ; GLMM).

difference between paranoia groups ( $\chi^2 = 0.05$ ,  $p = 0.818$ ). On the other hand, in the multi-reversal task, we found a greater separation than for the single-reversal task in win-switching behavior between paranoid and non-paranoid human participants (Figure 4B; pre-shift:  $\chi^2 = 112.23$ ,  $p < 0.001$ ; post-shift:  $\chi^2 = 74.74$ ,  $p < 0.001$ ; Figure S3). Presumably, this difference occurs because participants are experiencing the first contingency reversal just like the monkeys completing the single-reversal task. However, in the multi-reversal version, reversals continued prior to and following the contingency shift, which we believe would increase expected and unexpected volatility.

Similarly, we estimated volatility beliefs and value learning in paranoid and non-paranoid human participants completing both versions of the PRL task. The pattern of elevated volatility beliefs and reduced value learning, observed in the monkeys, was reflected in both versions of the PRL task in human participants but more clearly in the multi-reversal version. In the single-reversal task, no effects of paranoia were identified for volatility belief (Figure 4C, top). In contrast, in the multi-reversal task, high paranoia substantially augmented overall volatility belief (Figure 4D, top;  $m_3$ :  $\chi^2 = 61.16$ ,  $p < 0.001$ ), while shift marginally decreased it ( $m_3$ :  $\chi^2 = 3.02$ ,  $\beta = -0.09$ ,  $p = 0.08$ ). Thus, the multi-reversal task was more sensitive to differences in volatility belief between paranoia groups, regardless of shift, which mirrored the effects seen in MDmc-lesioned monkeys.

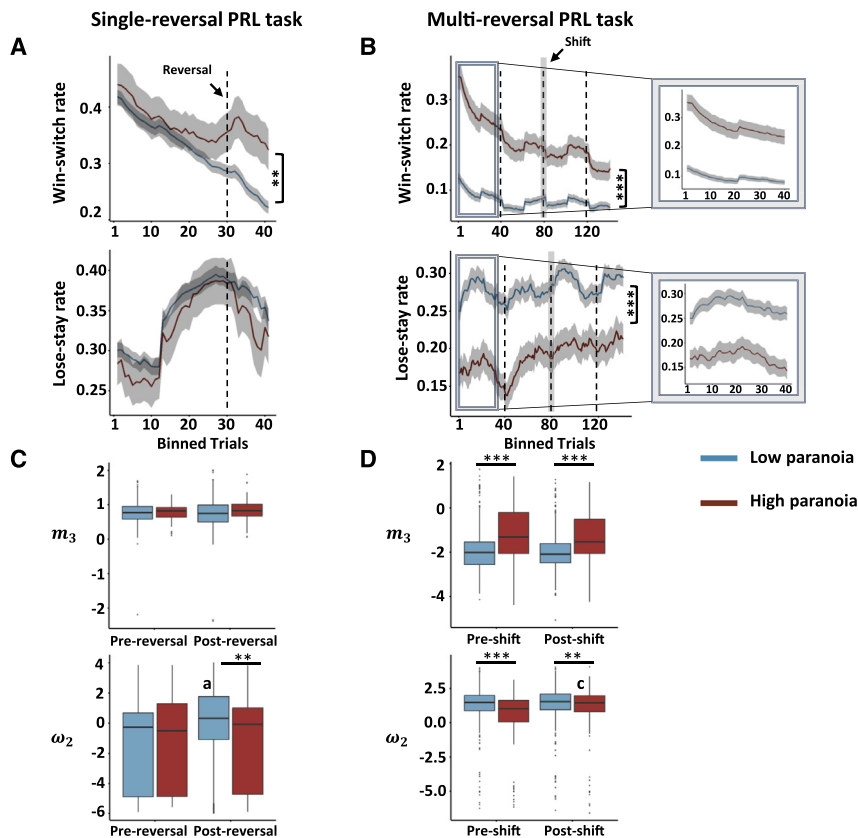
For value learning, an interaction between paranoia group and reversal phase was found in the single-reversal task (Figure 4C, bottom;  $\omega_2$ :  $\chi^2 = 4.41$ ,  $p = 0.036$ ; GLMM). Lower value learning was observed in high-paranoia individuals within post but not pre reversal ( $\omega_2$ :  $\chi^2 = 7.02$ ,  $p = 0.008$ ), whereas the reversal increased value learning within low- but not high-paranoia individuals ( $\omega_2$ :  $\chi^2 = 46.05$ ,  $p < 0.001$ ). In the multi-reversal task, an interaction between paranoia group and reversal phase was also observed (Figure 4D, bottom;  $\omega_2$ :  $\chi^2 = 6.43$ ,  $p = 0.011$ ; GLMM). However, unlike the single-reversal version, high-paranoia individuals had lower value learning than their low-paranoia counterparts within both pre ( $\omega_2$ :  $\chi^2 = 27.34$ ,  $p < 0.001$ ) and post shift ( $\omega_2$ :  $\chi^2 = 7.22$ ,  $p = 0.007$ ). The shift increased value learning

only within the high-paranoia group ( $\omega_2$ :  $\chi^2 = 4.62$ ,  $p = 0.032$ ). Overall, lower value learning was observed in the high-paranoia groups, much like MDmc-lesioned monkeys.

## DISCUSSION

It is controversial whether it is feasible to model psychiatric symptoms like paranoia in nonhuman animals. Presently, we met that challenge and conclude that nonhuman primates do indeed display choice responses akin to participants with high paranoia during a PRL task, and, thus, they may serve as models for the exploration of psychiatric symptoms hitherto considered outside the realm of translational neuroscience. We describe a series of analyses of PRL datasets from monkeys and human participants. We found that, in a single-reversal PRL task, monkeys with MDmc lesions but not OFC lesions exhibited increased win-switch behavior, increased volatility beliefs, and decreased value learning rates, especially after the reversal in reward contingencies. By contrast, OFC but not MDmc lesions increased lose-stay behavior and led to increased value learning rates while blunting updating of volatility belief. This pattern of behavioral responses and altered neurocomputations observed in monkeys with MDmc lesions (i.e., elevated volatility belief) was similar to that observed in human participants with high levels of paranoia. Indeed, the pattern of win-switching and volatility belief change was consistent across two independent datasets that employed different reinforcement contingency schemes. In cognitive neuropsychology, double dissociations are essential for the identification of independent functions. We report a double dissociation in the behavioral effects and computational consequences of MDmc versus OFC lesions on belief updating.

The challenge with PRL is to harbor a set of beliefs that is robust to noise but sensitive to real underlying contingency change. It appears that lesions of the MDmc and OFC have doubly dissociable effects on these processes. Specifically, in a single-reversal task, win-switching behavior and volatility beliefs increased in response to the contingency reversal in control animals, whereas lose-stay rates decreased. This represents



**Figure 4. Impact of the PRL task design on win-switch and lose-stay behavior and on beliefs about volatility and value in humans**

(A) Win-switch and lose-stay behaviors in non-paranoid and paranoid human participants during a single-reversal (after 30 trials) three-choice PRL task.

(B) Win-switch and lose-stay behaviors in non-paranoid and paranoid human participants during a multi-reversal (with a mid-way contingency shift after 80 trials) three-choice PRL task. Insets illustrate the effect of performance-based reversal (i.e., experience reversal upon 9 of 10 consecutive selections of the highest reinforcement probability deck) on trial-by-trial behavior in the multi-reversal task to show the similarity in the characteristic of task behavior in the single-reversal task versions.

(C and D) Differences in beliefs about volatility ( $m_3$  parameter from the HGF) and reward value learning ( $\omega_2$  parameter) in non-paranoid and paranoid human participants during a single-reversal three-choice PRL task (C) and during a multi-reversal three-choice PRL task (D).

Asterisks indicate significant paranoia group differences ( $***p < 0.001$ ,  $**p < 0.01$ ; GLMM). Other symbols indicate significant reversal phase differences in each lesion group (a,  $p < 0.001$ ; c,  $p < 0.05$ ; GLMM).

sensitivity to contingency change and the concomitant increase in volatility of the environment. Monkeys with OFC lesions fail to show this increase. This manifests as greater post-reversal lose-stay behavior. Conversely, monkeys with MDmc lesions make choices that are suggestive of elevated volatility beliefs both prior to reversal and a greater increase post reversal compared to control and OFC-lesioned animals, which is paralleled by the higher rates of win switching in the MDmc-lesioned animals. Furthermore, value learning rates increased post reversal in control animals, an effect mirrored in monkeys with OFC-lesions too. In contrast, monkeys with MDmc lesions fail to show this effect. It is possible that a higher-level task structure (level 3, HGF; Figure 1B), rather than value (level 2, HGF; Figure 1B), is more salient in MDmc-lesioned animals. In addition, in the absence of appropriate value learning, optimal choice responses may become noisier in MDmc-lesioned animals.<sup>25,26</sup>

Taken together, these data highlight the differential roles of MDmc and OFC nodes in belief processing. MDmc-lesioned animals appear to confuse stochasticity for real change (win switching more pre reversal, and even more post reversal), whereas OFC-lesioned animals do not update their beliefs in response to the volatility occasioned by reversal. Such maladaptive behavior in MDmc-lesioned monkeys manifests functionally as a failure to optimally learn complex reward associations and, thus, exploit the most rewarding option; distorted volatility beliefs could disrupt that ability to more frequently choose the highest probability reinforced option. Alternatively, maladaptive behavior manifests in OFC-lesioned monkeys as perseverative

responding. Such perseveration is not driven by sticking to outdated value (indeed, changes in value learning rates due to reversal are consistent with non-lesioned controls) but, rather, by a lack of adaptive change in volatility belief in the face of real change. Finally, these data further highlight the overlap between MDmc-lesioned monkey behavior and that of people who are paranoid, who also confuse stochasticity for volatility prior to reversal and have higher  $m_3$  (i.e., greater expectation of change) and lower  $\omega_2$  (i.e., slower value learning).

In addition to similarities, there were also some differences across datasets. While the single-reversal task in humans has superficial structural similarities to the monkey task, the fact that the monkeys underwent multiple reversal sessions (Figure S4) across days (and weeks) suggests that the beliefs that they build and update are more similar to the beliefs that human participants formed and updated during the multi-reversal task. This suggestion is supported by the computational analyses, which showed greater parallels between the MDmc-lesioned monkeys and humans with paranoia on the multi-reversal task; volatility beliefs were increased, and value learning rates were decreased. Hence, we demonstrate the utility of administering similar tasks and models across species and argue that computation may be a lingua franca<sup>27</sup> that could connect disparate fields of research and model organisms. However, close inspection of the pre- and post-reversal data shows that monkey win switching increases post reversal. Human win-switching does not increase. We do not know precisely why the monkey and human behavioral data differ in this way. We can speculate that monkeys are more sensitive to rewards than humans are to abstract points. After reversal, this sensitivity could drive monkeys

to switch more often to seek more rewarding options. Alternatively, monkeys and humans may be building differentially complex models of the task. In humans, this added complexity would accommodate reversals and shifts and contingencies, whereas monkeys may be switching erratically in search of an explanatory model.<sup>28,29</sup> This would actually be consistent with something we observed in humans. Paranoid participants tend to behave more randomly than non-paranoid participants,<sup>3</sup> although not differentially pre and post reversal.

Our data implicate MDmc contributing a role in volatility belief updating. Previous work has linked volatility belief updating to paranoia and persecutory delusions.<sup>3,4,6</sup> What is the role (if any) of MD in paranoia and psychosis? There are isolated cases where circumscribed lesions cause new-onset paranoia in humans<sup>19,30</sup> and frightening distortions of social stimuli,<sup>31</sup> but there are also more diffuse (often mnemonic) problems associated with MD damage, which can, as a result of impaired interpersonal behavior, lead indirectly to paranoia. Alternatively, genetic changes and the atrophy in frontotemporal dementia (FTD) are often quite specific to MD,<sup>32</sup> and psychotic symptoms (specifically paranoia) are common in FTD.<sup>25,33</sup> We suggest that our data are more aligned with a direct impact of thalamic damage on belief updating rather than a secondary impairment in relational cognition, which, although possible in monkeys, would not explain the specific pattern of findings we report presently because relational cognition and social interaction were not part of the present task. Furthermore, our human data suggest that the differences in high vs. low paranoia are not more readily attributable to confounding differences in mood or general intellectual function (Figure S5).

In prior work we argued that, because our PRL task was relatively non-social, and rats are a relatively asocial species (but not completely asocial),<sup>34–36</sup> the similarity in behavior between paranoid humans and rats treated with methamphetamine was consistent with a non-social explanation of paranoia. Our non-social result did not differ when we made the PRL more social,<sup>4,37</sup> although participants did imbue the non-social stimuli with harmful intentions, so social cognition may have been in play. Nonhuman primates are more social than rodents. Sociality is often related to theory of mind (ToM)—the cognitive capacity to attribute mental states to others.<sup>38</sup> This mentalizing process has been demonstrated in macaques<sup>39</sup> and apes.<sup>40</sup> However, ToM in humans and nonhuman primates may differ substantially; macaques do not attribute mental states to the Heider and Simmel animations.<sup>41</sup> Furthermore, putative ToM-like capacities in primates seem to fall short of human abilities.<sup>42</sup> ToM is associated with dorsomedial prefrontal cortex (DMPFC) activity in humans<sup>43</sup> and macaques.<sup>39</sup> The MDmc projects directly to the DMPFC.<sup>44</sup> Disrupting the MDmc may perturb social valuation and inference in the DMPFC. However, the DMPFC may also be less social specific than previously believed.<sup>45–47</sup> Indeed, a recent nonhuman primate study suggested that the DMPFC is involved in tracking and weighting the reliability of social and non-social information,<sup>48</sup> and in humans, the DMPFC may be tracking uncertainty rather than ToM.<sup>49</sup> Furthermore, human fMRI studies of the PRL suggest DMPFC responses during non-social reversal learning.<sup>5,50</sup> It is possible, on one hand, that mediodorsal thalamus (MD) disruption could lead to abnormal

ToM processing in the DMPFC via spurious inputs to the DMPFC. On the other hand, social belief volatility may reside in an encapsulated module independent from MD-mediated computations. Non-social environmental volatility is often attributed to social agents. Domain specificity may manifest at the level of algorithm (beliefs about environmental volatility vs. beliefs about agents' intentions) or at the level of implementation (domain-general vs. domain-specific nodes and circuits).<sup>51</sup> The involvement of the MD in the animism bias bears further investigation, although the bias seems absent in macaques.<sup>41</sup> Even if the underlying mechanisms of paranoia are not inherently social, the ramifications of their perturbation often bear social implications.<sup>51</sup> When people perceive or expect excess volatility, they typically find that anxiogenic and seek to explain it away. Typically, volatility is ascribed to other agents<sup>52</sup> and their intentions. When the world feels unpredictable, paranoid beliefs arise.<sup>4</sup> Consequently, our findings on the MDmc open an avenue to further explore the potential intersections and divergences between belief volatility and paranoia. Future work ought to leverage these tasks and models as well as other computational approaches to further elucidate the domain generality or specificity of paranoia.

Our findings have important implications for understanding the underlying mechanisms of paranoia and, relatedly, psychosis and schizophrenia. While prior work has implicated the OFC<sup>23,53</sup> and the DLPFC<sup>6</sup> in paranoia as well as the link between paranoia and reversal learning, our lesion findings highlight that the MDmc is a necessary node for controlling volatility belief. Future work ought to explore the human MDmc in the context of paranoia and its treatment. Furthermore, the widespread projection targets of MDmc within the prefrontal cortex may well prove relevant to belief updating; for example, the VLPFC. In the case of thalamo-prefrontal dynamics, it has been posited that damage to the thalamus often recapitulates the effects of prefrontal damage.<sup>54</sup> This was not what we observed with OFC damage. Even though the OFC projects directly to the MDmc, we found a double dissociation of the effects of OFC and MDmc damage on behavior and computational model parameters, suggesting that they underwrite independent, dissociable functions. We further explored the effects of VLPFC lesions on PRL task behavior (Figure S6). VLPFC lesions had effects redolent of MDmc lesions on behavior, elevating win switching and volatility beliefs. This is consistent with the hypothesized relay function between the thalamus and PFC in the thalamo-prefrontal-cortical circuits. What, then, are we to conclude about OFC and MDmc function? Rodent work suggests that the direction of coupling between brain regions is relevant to the circuit processing of reversal learning—damage to projections from the OFC has different effects than damage to projections to the OFC.<sup>55</sup> Similar manipulations of the projections between the MDmc and OFC will be revelatory in this regard. Furthermore, it is possible that the effects of circumscribed lesions have differential effects on some other brain region(s) with which both the OFC and MDmc interact. One possibility is the nucleus accumbens.<sup>55</sup>

More broadly, the current results suggest that not all PRL tasks are equally sensitive to paranoia in humans. While the overall patterns of behavioral response and model parameter change are consistent across tasks and species, the multi-reversal version seemed to separate paranoid from non-paranoid participants



more clearly and to cleave more closely to the patterns observed in monkeys. This task type specificity implies that future studies intended to predict paranoia or track its prognosis might result in most predictive data when applying the multi-reversal version of the task. Based on the evidence provided here from MDmc-lesioned monkeys, we predict that human neuroimaging studies would detect an MDmc involvement in the more complex, multi-reversal version, which may be more difficult to detect using the single-reversal task. Furthermore, human participants with MDmc damage are predicted to be particularly challenged by a multi-reversal version of the PRL task.

Previous work has established an association between belief updating and paranoia through computational modeling.<sup>3,4,56</sup> Here, we map computational markers of belief updating to specific brain regions in monkeys, where the MDmc and OFC play dissociable roles. Parallels between MDmc-lesioned monkeys and paranoid human participants suggest an influence of the MDmc in paranoia that warrants further investigation. Our study represents an initial cross-species effort to use computational modeling and common behavioral tasks to translate the neural mechanisms of belief updating from monkeys to humans and, ultimately, into the clinic.

### Limitations

While we observed similarities between lesioned monkeys and humans with paranoia performing PRL, there were differences in behavioral patterns and in the tasks to which participants were exposed. Future work might present participants of both species with a multi-reversal task featuring a contingency shift. Furthermore, though our aim was to extend existing human computational psychiatry into the nonhuman primate realm, the possibility exists that a better-fitting model exists that spans these datasets. That will be a target for future investigation. Finally, the human imaging and neuropsychology literatures implicate the mediodorsal thalamus in paranoia and schizophreniform psychosis; however, this association bears investigation using PRL tasks, functional imaging, and computational modeling in human participants and patients in particular.

### STAR★METHODS

Detailed methods are provided in the online version of this paper and include the following:

- **KEY RESOURCES TABLE**
- **RESOURCE AVAILABILITY**
  - Lead contact
  - Materials availability
  - Data and code availability
- **EXPERIMENTAL MODEL AND SUBJECT DETAILS**
  - Monkeys and humans
- **METHOD DETAILS**
  - Surgical and lesion procedures for monkeys
  - Questionnaire for human participants
  - Three-choice PRL task
  - Quantification and statistical analysis
  - Computational modeling
- **STATISTICAL ANALYSIS**
  - General
  - GLMMs
  - Permutation tests

### SUPPLEMENTAL INFORMATION

Supplemental information can be found online at <https://doi.org/10.1016/j.celrep.2024.114355>.

### ACKNOWLEDGMENTS

This work was supported by the Yale University Department of Psychiatry, the Connecticut Mental Health Center (CMHC) and Connecticut State Department of Mental Health and Addiction Services (DMHAS), and the Kavli Foundation (to P.S.). It was funded by NIMH R01MH12887 (to P.R.C.), NIMH R21MH120799-01 (to P.R.C.), NIMH R01MH128190 (to S.W.C.C.), and NIMH R01MH120081 (to S.W.C.C.). The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

### AUTHOR CONTRIBUTIONS

P.S., P.R.C., and S.W.C.C. conceived the study. P.H.R., M.E.W., S.C., A.S.M., M.P.N., and E.A.M. collected these data. C.D.M. provided analytical tools. V.D.C., S.M.G., and J.R.T. helped design the study. P.S., S.L.T., and R.A.R.-G. analyzed these data. All authors discussed the overarching approach and wrote and edited the manuscript.

### DECLARATION OF INTERESTS

P.R.C. is a cofounder of Tetricus Labs Inc., which did not fund this work.

Received: November 2, 2023

Revised: April 13, 2024

Accepted: May 29, 2024

Published: June 13, 2024

### REFERENCES

1. Walker, E.Y., Pohl, S., Denison, R.N., Barack, D.L., Lee, J., Block, N., Ma, W.J., and Meyniel, F. (2023). Studying the neural representations of uncertainty. *Nat. Neurosci.* 26, 1857–1867. <https://doi.org/10.1038/s41593-023-01444-y>.
2. Feeney, E.J., Groman, S.M., Taylor, J.R., and Corlett, P.R. (2017). Explaining Delusions: Reducing Uncertainty Through Basic and Computational Neuroscience. *Schizophr. Bull.* 43, 263–272. <https://doi.org/10.1093/schbul/sbw194>.
3. Reed, E.J., Uddenberg, S., Suthaharan, P., Mathys, C.D., Taylor, J.R., Groman, S.M., and Corlett, P.R. (2020). Paranoia as a deficit in non-social belief updating. *Elife* 9, e56345. <https://doi.org/10.7554/eLife.56345>.
4. Suthaharan, P., Reed, E.J., Leptourgos, P., Kenney, J.G., Uddenberg, S., Mathys, C.D., Litman, L., Robinson, J., Moss, A.J., Taylor, J.R., et al. (2021). Paranoia and belief updating during the COVID-19 crisis. *Nat. Hum. Behav.* 5, 1190–1202. <https://doi.org/10.1038/s41562-021-01176-8>.
5. Albein-Urios, N., Chase, H., Clark, L., Kirkovski, M., Davies, C., and Enticott, P.G. (2019). Increased perseverative errors following high-definition transcranial direct current stimulation over the ventrolateral cortex during probabilistic reversal learning. *Brain Stimul.* 12, 959–966. <https://doi.org/10.1016/j.brs.2019.02.013>.
6. Cole, D.M., Diaconescu, A.O., Pfeiffer, U.J., Brodersen, K.H., Mathys, C.D., Julkowsky, D., Ruhrmann, S., Schilbach, L., Tittgemeyer, M., and Vogeley, K. (2020). Atypical processing of uncertainty in individuals at risk for psychosis. *Neuroimage. Clin.* 26, 102239. <https://doi.org/10.1016/j.nicl.2020.102239>.
7. Costa, V.D., Tran, V.L., Turchi, J., and Averbach, B.B. (2015). Reversal Learning and Dopamine: A Bayesian Perspective. *J. Neurosci.* 35, 2407–2416. <https://doi.org/10.1523/JNEUROSCI.1989-14.2015>.
8. Gershman, S.J., and Uchida, N. (2019). Believing in dopamine. *Nat. Rev. Neurosci.* 20, 703–714. <https://doi.org/10.1038/s41583-019-0220-7>.

9. Eckstein, M.K., Master, S.L., Dahl, R.E., Wilbrecht, L., and Collins, A.G.E. (2022). Reinforcement learning and Bayesian inference provide complementary models for the unique advantage of adolescents in stochastic reversal. *Dev. Cogn. Neurosci.* 55, 101106. <https://doi.org/10.1016/j.dcn.2022.101106>.
10. Mathys, C.D., Lomakina, E.I., Daunizeau, J., Iglesias, S., Brodersen, K.H., Friston, K.J., and Stephan, K.E. (2014). Uncertainty in perception and the Hierarchical Gaussian Filter. *Front. Hum. Neurosci.* 8, 825. <https://doi.org/10.3389/fnhum.2014.00825>.
11. Massi, B., Donahue, C.H., and Lee, D. (2018). Volatility Facilitates Value Updating in the Prefrontal Cortex. *Neuron* 99, 598–608.e4. <https://doi.org/10.1016/j.neuron.2018.06.033>.
12. Soltani, A., and Izquierdo, A. (2019). Adaptive learning under expected and unexpected uncertainty. *Nat. Rev. Neurosci.* 20, 635–644. <https://doi.org/10.1038/s41583-019-0180-y>.
13. Chakraborty, S., Kolling, N., Walton, M.E., and Mitchell, A.S. (2016). Critical role for the mediodorsal thalamus in permitting rapid reward-guided updating in stochastic reward environments. *Elife* 5, e13588. <https://doi.org/10.7554/eLife.13588>.
14. Mitchell, A.S. (2015). The mediodorsal thalamus as a higher order thalamic relay nucleus important for learning and decision-making. *Neurosci. Biobehav. Rev.* 54, 76–88. <https://doi.org/10.1016/j.neubiorev.2015.03.001>.
15. Xiao, D., Zikopoulos, B., and Barbas, H. (2009). Laminar and modular organization of prefrontal projections to multiple thalamic nuclei. *Neuroscience* 161, 1067–1081. <https://doi.org/10.1016/j.neuroscience.2009.04.034>.
16. Phillips, J.M., Kambi, N.A., Redinbaugh, M.J., Mohanta, S., and Saalman, Y.B. (2021). Disentangling the influences of multiple thalamic nuclei on prefrontal cortex and cognitive control. *Neurosci. Biobehav. Rev.* 128, 487–510. <https://doi.org/10.1016/j.neubiorev.2021.06.042>.
17. Wolff, M., and Halassa, M.M. (2024). The mediodorsal thalamus in executive control. *Neuron* 112, 893–908. <https://doi.org/10.1016/j.neuron.2024.01.002>.
18. Mukherjee, A., Lam, N.H., Wimmer, R.D., and Halassa, M.M. (2021). Thalamic circuits for independent control of prefrontal signal and noise. *Nature* 600, 100–104. <https://doi.org/10.1038/s41586-021-04056-3>.
19. Crail-Melendez, D., Atriano-Mendieta, C., Carrillo-Meza, R., and Ramirez-Bermudez, J. (2013). Schizophrenia-like psychosis associated with right lacunar thalamic infarct. *Neurocase* 19, 22–26. <https://doi.org/10.1080/13554794.2011.654211>.
20. Rudebeck, P.H., Behrens, T.E., Kennerley, S.W., Baxter, M.G., Buckley, M.J., Walton, M.E., and Rushworth, M.F.S. (2008). Frontal Cortex Subregions Play Distinct Roles in Choices between Actions and Stimuli. *J. Neurosci.* 28, 13775–13785. <https://doi.org/10.1523/JNEUROSCI.3541-08.2008>.
21. Walton, M.E., Behrens, T.E.J., Buckley, M.J., Rudebeck, P.H., and Rushworth, M.F.S. (2010). Separable learning systems in the macaque brain and the role of orbitofrontal cortex in contingent learning. *Neuron* 65, 927–939. <https://doi.org/10.1016/j.neuron.2010.02.027>.
22. Rudebeck, P.H., Saunders, R.C., Lundgren, D.A., and Murray, E.A. (2017). Specialized Representations of Value in the Orbital and Ventrolateral Prefrontal Cortex: Desirability versus Availability of Outcomes. *Neuron* 95, 1208–1220.e5. <https://doi.org/10.1016/j.neuron.2017.07.042>.
23. Walther, S., Lefebvre, S., Conring, F., Gangl, N., Nadesalingam, N., Alexaki, D., Wüthrich, F., Rüter, M., Viher, P.V., Federspiel, A., et al. (2022). Limbic links to paranoia: increased resting-state functional connectivity between amygdala, hippocampus and orbitofrontal cortex in schizophrenia patients with paranoia. *Eur. Arch. Psychiatry Clin. Neurosci.* 272, 1021–1032. <https://doi.org/10.1007/s00406-021-01337-w>.
24. Barnby, J.M., Mehta, M.A., and Moutoussis, M. (2022). The computational relationship between reinforcement learning, social inference, and paranoia. *PLoS Comput. Biol.* 18, e1010326. <https://doi.org/10.1371/journal.pcbi.1010326>.
25. Perry, B.A.L., Lomi, E., and Mitchell, A.S. (2021). Thalamocortical interactions in cognition and disease: The mediodorsal and anterior thalamic nuclei. *Neurosci. Biobehav. Rev.* 130, 162–177. <https://doi.org/10.1016/j.neubiorev.2021.05.032>.
26. Ouhaz, Z., Fleming, H., and Mitchell, A.S. (2018). Cognitive Functions and Neurodevelopmental Disorders Involving the Prefrontal Cortex and Mediodorsal Thalamus. *Front. Neurosci.* 12, 33. <https://doi.org/10.3389/fnins.2018.00033>.
27. Corlett, P.R., and Schoenbaum, G. (2021). Leveraging Basic Science for the Clinic-From Bench to Bedside. *JAMA Psychiatr.* 78, 331–334. <https://doi.org/10.1001/jamapsychiatry.2020.3656>.
28. Tervo, D.G.R., Kuleshova, E., Manakov, M., Proskurin, M., Karlsson, M., Lustig, A., Behnam, R., and Karpova, A.Y. (2021). The anterior cingulate cortex directs exploration of alternative strategies. *Neuron* 109, 1876–1887.e6. <https://doi.org/10.1016/j.neuron.2021.03.028>.
29. Mitchell, A.S., Baxter, M.G., and Gaffan, D. (2007). Dissociable Performance on Scene Learning and Strategy Implementation after Lesions to Magnocellular Mediodorsal Thalamic Nucleus. *J. Neurosci.* 27, 11888–11895. <https://doi.org/10.1523/JNEUROSCI.1835-07.2007>.
30. Loh, A., Germann, J., Qazi, S., Husain, R., Boutet, A., Lozano, A.M., and Mansouri, A. (2022). Lesion network mapping of ectopic craniopharyngioma identifies potential cause of psychosis: a case report. *Acta Neurochir.* 164, 3285–3289. <https://doi.org/10.1007/s00701-022-05355-y>.
31. Delgado, M.G., and Bogousslavsky, J. (2013). “Distortoidolias” - fantastic perceptive distortion. A new, pure dorsomedial thalamic syndrome. *Eur. Neurol.* 70, 6–9. <https://doi.org/10.1159/000348361>.
32. Bocchetta, M., Iglesias, J.E., Neason, M., Cash, D.M., Warren, J.D., and Rohrer, J.D. (2020). Thalamic nuclei in frontotemporal dementia: Mediodorsal nucleus involvement is universal but pulvinar atrophy is unique to C9orf72. *Hum. Brain Mapp.* 41, 1006–1016. <https://doi.org/10.1002/hbm.24856>.
33. Landqvist Waldö, M., Gustafson, L., Passant, U., and Englund, E. (2015). Psychotic symptoms in frontotemporal dementia: a diagnostic dilemma? *Int. Psychogeriatr.* 27, 531–539. <https://doi.org/10.1017/S1041610214002580>.
34. Donaldson, T.N., Barto, D., Bird, C.W., Magcalas, C.M., Rodriguez, C.I., Fink, B.C., and Hamilton, D.A. (2018). Social order: Using the sequential structure of social interaction to discriminate abnormal social behavior in the rat. *Learn. Motiv.* 61, 41–51. <https://doi.org/10.1016/j.lmot.2017.03.003>.
35. Kondrakiewicz, K., Kostecki, M., Szadzińska, W., and Knapska, E. (2019). Ecological validity of social interaction tests in rats and mice. *Genes Brain Behav.* 18, e12525. <https://doi.org/10.1111/gbb.12525>.
36. Urbach, Y.K., Bode, F.J., Nguyen, H.P., Riess, O., and von Hörsten, S. (2010). Neurobehavioral Tests in Rat Models of Degenerative Brain Diseases. In *Rat Genomics: Methods and Protocols Methods in Molecular Biology*, I. Anegon, ed. (Humana Press), pp. 333–356. [https://doi.org/10.1007/978-1-60327-389-3\\_24](https://doi.org/10.1007/978-1-60327-389-3_24).
37. Suthaharan, P., and Corlett, P.R. (2023). Assumed shared belief about conspiracy theories in social networks protects paranoid individuals against distress. *Sci. Rep.* 13, 6084. <https://doi.org/10.1038/s41598-023-33305-w>.
38. Barnby, J.M., Dayan, P., and Bell, V. (2023). Formalising social representation to explain psychiatric symptoms. *Trends Cogn. Sci.* 27, 317–332. <https://doi.org/10.1016/j.tics.2022.12.004>.
39. Hayashi, T., Akikawa, R., Kawasaki, K., Egawa, J., Minamimoto, T., Kobayashi, K., Kato, S., Hori, Y., Nagai, Y., Iijima, A., et al. (2020). Macaques Exhibit Implicit Gaze Bias Anticipating Others’ False-Belief-Driven Actions via Medial Prefrontal Cortex. *Cell Rep.* 30, 4433–4444.e5. <https://doi.org/10.1016/j.celrep.2020.03.013>.
40. Krupenye, C., Kano, F., Hirata, S., Call, J., and Tomasello, M. (2016). Great apes anticipate that other individuals will act according to false beliefs. *Science* 354, 110–114. <https://doi.org/10.1126/science.aaf8110>.

41. Schafroth, J.L., Basile, B.M., Martin, A., and Murray, E.A. (2021). No evidence that monkeys attribute mental states to animated shapes in the Heider–Simmel videos. *Sci. Rep.* *11*, 3050. <https://doi.org/10.1038/s41598-021-82702-6>.
42. Devaine, M., San-Galli, A., Trapanese, C., Bardino, G., Hano, C., Saint Jalme, M., Bouret, S., Masi, S., and Daunizeau, J. (2017). Reading wild minds: A computational assay of Theory of Mind sophistication across seven primate species. *PLoS Comput. Biol.* *13*, e1005833. <https://doi.org/10.1371/journal.pcbi.1005833>.
43. Schurz, M., Radua, J., Aichhorn, M., Richlan, F., and Perner, J. (2014). Fractionating theory of mind: A meta-analysis of functional brain imaging studies. *Neurosci. Biobehav. Rev.* *42*, 9–34. <https://doi.org/10.1016/j.neurobiorev.2014.01.009>.
44. Pergola, G., Danet, L., Pitel, A.-L., Carlesimo, G.A., Segobin, S., Pariente, J., Suchan, B., Mitchell, A.S., and Barbeau, E.J. (2018). The Regulatory Role of the Human Mediodorsal Thalamus. *Trends Cogn. Sci.* *22*, 1011–1025. <https://doi.org/10.1016/j.tics.2018.08.006>.
45. Chen, L.L., and Wise, S.P. (1996). Evolution of Directional Preferences in the Supplementary Eye Field during Acquisition of Conditional Oculomotor Associations. *J. Neurosci.* *16*, 3067–3081. <https://doi.org/10.1523/JNEUROSCI.16-09-03067.1996>.
46. Jamali, M., Grannan, B.L., Fedorenko, E., Saxe, R., Báez-Mendoza, R., and Williams, Z.M. (2021). Single-neuronal predictions of others' beliefs in humans. *Nature* *591*, 610–614. <https://doi.org/10.1038/s41586-021-03184-0>.
47. Corlett, P.R., Mollick, J.A., and Kober, H. (2022). Meta-analysis of human prediction error for incentives, perception, cognition, and action. *Neuropsychopharmacology* *47*, 1339–1349. <https://doi.org/10.1038/s41386-021-01264-3>.
48. Mahmoodi, A., Harbison, C., Bongioanni, A., Emberton, A., Roumazeilles, L., Sallet, J., Khalighinejad, N., and Rushworth, M.F.S. (2024). A frontopolar-temporal circuit determines the impact of social information in macaque decision making. *Neuron* *112*, 84–92.e6. <https://doi.org/10.1016/j.neuron.2023.09.035>.
49. Berkay, D., and Jenkins, A.C. (2023). A Role for Uncertainty in the Neural Distinction Between Social and Nonsocial Thought. *Perspect. Psychol. Sci.* *18*, 491–502. <https://doi.org/10.1177/17456916221112077>.
50. Cools, R., Clark, L., Owen, A.M., and Robbins, T.W. (2002). Defining the Neural Mechanisms of Probabilistic Reversal Learning Using Event-Related Functional Magnetic Resonance Imaging. *J. Neurosci.* *22*, 4563–4567. <https://doi.org/10.1523/JNEUROSCI.22-11-04563.2002>.
51. Lockwood, P.L., Apps, M.A.J., and Chang, S.W.C. (2020). Is There a 'Social' Brain? Implementations and Algorithms. *Trends Cogn. Sci.* *24*, 802–813. <https://doi.org/10.1016/j.tics.2020.06.011>.
52. Sullivan, D., Landau, M.J., and Rothschild, Z.K. (2010). An existential function of enemyship: evidence that people attribute influence to personal and political enemies to compensate for threats to control. *J. Pers. Soc. Psychol.* *98*, 434–449. <https://doi.org/10.1037/a0017457>.
53. Premkumar, P., Fannon, D., Sapara, A., Peters, E.R., Anilkumar, A.P., Simmons, A., Kuipers, E., and Kumari, V. (2015). Orbitofrontal cortex, emotional decision-making and response to cognitive behavioural therapy for psychosis. *Psychiatry Res.* *231*, 298–307. <https://doi.org/10.1016/j.psychres.2015.01.013>.
54. Parnaudeau, S., Bolkan, S.S., and Kellendonk, C. (2018). The Mediodorsal Thalamus: An Essential Partner of the Prefrontal Cortex for Cognition. *Biol. Psychiatry* *83*, 648–656. <https://doi.org/10.1016/j.biopsych.2017.11.008>.
55. Groman, S.M., Keistler, C., Keip, A.J., Hammarlund, E., DiLeone, R.J., Pittenger, C., Lee, D., and Taylor, J.R. (2019). Orbitofrontal Circuits Control Multiple Reinforcement-Learning Processes. *Neuron* *103*, 734–746.e3. <https://doi.org/10.1016/j.neuron.2019.05.042>.
56. Sheffield, J.M., Suthaharan, P., Leptourgos, P., and Corlett, P.R. (2022). Belief Updating and Paranoia in Individuals With Schizophrenia. *Biol. Psychiatry. Cogn. Neurosci. Neuroimaging* *7*, 1149–1157. <https://doi.org/10.1016/j.bpsc.2022.03.013>.
57. Frässle, S., Aponte, E.A., Bollmann, S., Brodersen, K.H., Do, C.T., Harrison, O.K., Harrison, S.J., Heinzle, J., Iglesias, S., Kasper, L., et al. (2021). TAPAS: An Open-Source Software Package for Translational Neuromodeling and Computational Psychiatry. *Front. Psychiatry* *12*, 680811.
58. Ilinsky, I.A., and Kultas-Ilinsky, K. (1987). Sagittal cytoarchitectonic maps of the Macaca mulatta thalamus with a revised nomenclature of the motor-related nuclei validated by observations on their connectivity. *J. Comp. Neurol.* *262*, 331–364. <https://doi.org/10.1002/cne.902620303>.
59. Browning, P.G.F., Chakraborty, S., and Mitchell, A.S. (2015). Evidence for Mediodorsal Thalamus and Prefrontal Cortex Interactions during Cognition in Macaques. *Cereb. Cortex* *25*, 4519–4534. <https://doi.org/10.1093/cercor/bhv093>.
60. Freeman, D., Loe, B.S., Kingdon, D., Startup, H., Molodynski, A., Rosebrock, L., Brown, P., Sheaves, B., Waite, F., and Bird, J.C. (2021). The revised Green et al., Paranoid Thoughts Scale (R-GPTS): psychometric properties, severity ranges, and clinical cut-offs. *Psychol. Med.* *51*, 244–253. <https://doi.org/10.1017/S0033291719003155>.
61. Beck, A.T., Steer, R.A., and Brown, G. (1996). Beck Depression Inventory–II. <https://doi.org/10.1037/t00742-000>.
62. Noonan, M.P., Walton, M.E., Behrens, T.E.J., Sallet, J., Buckley, M.J., and Rushworth, M.F.S. (2010). Separate value comparison and learning mechanisms in macaque medial and lateral orbitofrontal cortex. *Proc. Natl. Acad. Sci.* *107*, 20547–20552. <https://doi.org/10.1073/pnas.1012246107>.

## STAR★METHODS

### KEY RESOURCES TABLE

REAGENT or RESOURCE	SOURCE	IDENTIFIER
<b>Other</b>		
Control lesion data	Walton et al., 2010 <sup>21</sup>	<a href="https://doi.org/10.1016/j.neuron.2010.02.027">https://doi.org/10.1016/j.neuron.2010.02.027</a>
MDmc lesion and control data	Chakraborty et al., 2016 <sup>13</sup>	<a href="https://doi.org/10.7554/eLife.13588">https://doi.org/10.7554/eLife.13588</a>
OFC lesion and control data	Rudebeck et al., 2017 <sup>22</sup>	<a href="https://doi.org/10.1016/j.neuron.2017.07.042">https://doi.org/10.1016/j.neuron.2017.07.042</a>
Multi-reversal PRL human data	Suthaharan et al., 2021 <sup>4</sup>	<a href="https://doi.org/10.1038/s41562-021-01176-8">https://doi.org/10.1038/s41562-021-01176-8</a>
Single-reversal PRL human data	Barnby et al., 2022 <sup>24</sup>	<a href="https://doi.org/10.1371/journal.pcbi.1010326">https://doi.org/10.1371/journal.pcbi.1010326</a>
<b>Software and algorithms</b>		
MATLAB (Data modeling)	Mathworks	2019a
RStudio (Data visualization)	RStudio	4.3.0
HGF toolbox v5.3.1	Frässle et al., 2021 <sup>57</sup>	<a href="https://translationalneuromodeling.github.io/tapas">https://translationalneuromodeling.github.io/tapas</a>
Codes	This paper	<a href="https://github.com/psuthaharan/belief-update-monkeys">https://github.com/psuthaharan/belief-update-monkeys</a>

### RESOURCE AVAILABILITY

#### Lead contact

Further information and requests for resources and reagents should be directed to and will be fulfilled by either of the co-senior authors, Philip R. Corlett ([philip.corlett@yale.edu](mailto:philip.corlett@yale.edu)) and Steve W. C. Chang ([steve.chang@yale.edu](mailto:steve.chang@yale.edu)).

#### Materials availability

This study did not generate new unique reagents.

#### Data and code availability

Monkey behavioral data presented in this paper will be available upon request. Human behavioral data from the single-reversal PRL task<sup>24</sup> and the multi-reversal with contingency shift PRL task data<sup>4</sup> are publically available. The analysis scripts for this paper can be found at <https://github.com/psuthaharan/belief-update-monkeys>.

### EXPERIMENTAL MODEL AND SUBJECT DETAILS

All monkey and human experiments in this study were in accordance with local and national acts and committees for the use of animals in scientific research and behavioral research in humans. Experiments from Study 1<sup>21</sup> and Study 2<sup>13</sup> were performed in compliance with the United Kingdom Animals (Scientific Procedures) Act of 1986. A Home Office (UK) Project License (PPL 30/2678) obtained after review by the University of Oxford Animal Care and Ethical Review Committee licensed all procedures. The housing and husbandry were in compliance with the guidelines of the European Directive (2010/63/EU) for the care and use of laboratory animals.

All procedures from Study 3<sup>22</sup> were reviewed and approved by the National Institute of Mental Health (NIMH) Animal Care and Use Committee.

For one dataset,<sup>4</sup> all human behavioral experiments were conducted at the Connecticut Mental Health Center in strict accordance with Yale University's Human Investigation Committee who provided ethical review and exemption approval (no. 2000026290). Written informed consent was provided by all research participants. For the other dataset,<sup>24</sup> all human behavioral experiments were internally reviewed and approved by the Research Ethics Committee at King's College London, UK (ref: RESCM-19/20-0603). Human participants gave consent by ticking checkboxes online following the information sheet, and prior to the administration of questionnaires or tasks.

#### Monkeys and humans

A total of twenty rhesus macaque monkeys (*Macaca mulatta*, all males) across three laboratories and a total of 1,225 online human participants from two laboratories were included in the study. Each animal was individually-, pair-, or group-housed, and was kept on a 12-h light dark cycle and had access to water 24 h a day. All experiments were conducted during the light phase.

### Monkeys

Fourteen monkeys – three from Study 1, three from Study 2, and eight from Study 3 – served as unoperated controls. Six monkeys – three from Study 2 and three from Study 3 – received excitotoxic lesions of different parts of the brain.

### Humans

We had 692 online human participants (collected via *Prolific*), of whom 84 had high paranoia, who completed the single-reversal task of the task most similar to the monkey task.<sup>24</sup> We had 533 online human participants (collected via *CloudResearch*), of whom 140 had high paranoia, who completed a multi-reversal version of the task.<sup>4</sup>

## METHOD DETAILS

### Surgical and lesion procedures for monkeys

Surgical procedures for the monkeys have been previously described in detail.<sup>13,22</sup> Aseptic neurosurgery was conducted under general anesthesia (isoflurane, 1-2% to effect), in a dedicated operating theater and with pre- and post-operative analgesia and antibiotics. In each animal that received a lesion, during the surgery the skin and fascia were opened, the muscles retracted, and a bilateral bone flap was taken in the cranium over the target region. The dura over the posterior part of the hemisphere was cut and retracted to the midline. The splenium of the corpus callosum was then sectioned and the choroid plexus cut to provide access to and enable visualization of the posterior medial thalamus. Injections were then made as described below. A semi-circular dural flap was reflected toward the orbit to allow access to the ventral surface of the frontal lobe in each hemisphere. Injections were then made into the OFC on the basis of sulcal landmarks as described below.

#### Magnocellular mediodorsal thalamus (MDmc)

In brief, the procedure focused on the dorsal thalamic nuclei, specifically the magnocellular mediodorsal thalamus (MDmc) in the monkey brain. The coordinates of the intended lesion site (see [Figure S7](#)) within the mediodorsal thalamus were taken from a comprehensive monkey brain atlas.<sup>58</sup> Moreover, evidence from prior studies<sup>13,59</sup> on the use of ibotenic acid and NMDA for successful mediodorsal thalamic lesions in monkeys supported this surgical procedure. Ten x 1.0ul injections (5 per hemisphere) of ibotenic acid/NMDA, spaced at least 1 mm apart were made into the thalamus after visualizing the anterior-posterior extent of the MD in monkeys' MD1, MD2, and MD3. MD lesions were confirmed histologically.<sup>13</sup>

#### Orbitofrontal cortex (OFC)

The procedure targeted Walker's areas 11, 13, and 14 in each hemisphere of the granular orbitofrontal cortex of the monkey brain. Injections of ibotenic acid were made into each hemisphere in separate surgeries with at least two weeks recovery time between surgeries. In each surgery, multiple 0.1ul ibotenic acid injections were made into the cortex on the orbital surface, specifically between the fundus of the lateral orbital sulcus and the rostral sulcus on the medial surface of the hemisphere. The anterior-posterior boundary of the lesion was the anterior and posterior ends of the medial and lateral orbital sulci. In each hemisphere, between 71 and 119 injections were made to cover the entirety of the OFC. The extent of the bilateral excitotoxic lesions in each monkey (see [Figure S7](#)) were verified by observing white hypersignal associated with edema in T2-weighted MRI scans taken within 5 days of the surgery and T1-weighted scans taken >1 year after the initial surgery.<sup>22</sup>

In sum, the monkeys we chose to study had bilateral excitotoxic lesions to their magnocellular mediodorsal thalamus (MDmc) and orbitofrontal cortex (OFC). Because we focused on win-switch and lose-stay behavior for estimating belief volatility, one animal from Study 3 was excluded from the present analysis as he was a formal outlier with regards to his win-switch behavior (see [Figure S8](#)).

### Questionnaire for human participants

Human participants completed a revised Green et al. paranoid thought scale (R-GPTS) as a means to collect self-report paranoia levels.<sup>60</sup> The paranoia group labels were based on subscale B of the R-GPTS; that is, *high* paranoia labels were based on a clinical cut-off of  $\geq 11$  (this threshold discriminates persecutory delusions from non-clinical paranoia, see [Figure S9](#) for link between paranoia and PRL task behavior). Self-report depression levels were measured using the Beck Depression Inventory II (BDI-II) scale,<sup>61</sup> excluding 'suicidal thoughts' item as per protocol; a clinical cut-off of  $\geq 17$  for clinical depression (i.e., *high* depression).

### Three-choice PRL task

#### Single-reversal

A PRL task – with three options to choose from, and a single reversal of the reward probabilities of the options – was used in this study for the behavior of both monkeys and human participants.

#### Monkeys

This task has been described extensively elsewhere.<sup>13,20–22,62</sup> In brief, while inside a wheeled transport cage, monkeys were positioned in front of a touch sensitive monitor. On each trial they were presented with three stimuli that were novel at the beginning of each 300-trial session. When monkeys touched one of the stimuli on the screen, a reward (190 mg food pellets, Noyes) was delivered based on a predetermined probabilistic reward schedule. Initially, one stimulus was associated with the highest probability of reward delivery. At around trial 150, there was a reversal in reward contingencies such that the stimulus associated with highest probability of reward became the lowest and the stimulus associated with the lowest changed to being associated with the highest probability of reward. For this report we analyzed one of four reward schedules that was similar to the human task: 'stable'. This corresponded to



Schedule 1 from a previous study<sup>22</sup>; pre-reversal reward contingencies followed a 0.61-0.2-0 pattern and changed to a post-reversal reward contingency of 0.19-0.44-0.76 pattern such that the most rewarding stimuli became the least rewarding stimuli and vice versa. Each monkey completed the stable version of the task for a total of 300 trials per session for five sessions (monkeys in the MDmc study completed five sessions so we selected the first five sessions from the OFC study). After each trial there was either a 2 s (Studies 1 and 2) or a 5 s (Study 3) intertrial interval before the next trial began. For Studies 1 and 2, a lunch box containing the monkey's food for the day was opened at the conclusion of the 300-trial session whereas in Study 3, monkeys were given their daily food amount in their home cage.

### Humans

The performance of the monkeys was compared to that of human participants who completed a similar version of the PRL task.<sup>24</sup> This non-social PRL task involved presenting participants with a series of symbols across 60 trials. Each symbol yielded either a positive outcome (+10 points) or a negative outcome (−5 points). At the beginning of the task, participants were informed that among the three symbols, one had a high likelihood (80%) of providing +10 points, another had an even chance (50%), and the remaining one had a low probability (20%) of yielding +10 points. Crucially, participants were made aware that the probabilities associated with each symbol could change unpredictably during the course of the task. At trial 30, participants were explicitly asked to indicate which symbol they believed offered the highest probability of delivering points. Following this judgment phase, the contingencies of the symbols were altered for the final 30 trials – the symbol previously associated with the lowest probability of gaining points became the one with the highest probability, the symbol with the highest probability became the even probability symbol, and the symbol with an even probability became the one with the lowest probability. Upon task completion, participants were once again prompted to identify the symbol they thought had provided the most points throughout the task.

### Multi-reversal

This version of the task is similar in theory but different in structure.<sup>3</sup> Human participants completed this PRL task in two different scenarios – one with a non-social aspect using a deck of cards another with a social aspect using avatars representing partners. In the card deck scenario, participants were asked to choose from three decks – a positive (+100) or negative (−50) outcome – with the goal of earning the most points. It was also mentioned that the deck with the highest reinforcement probability could change. In the partner scenario, participants were asked to choose from three avatars and to imagine working on a group project with them. The avatars could represent helpful (+100) or hurtful (−50) partners and the partner that would give them the best chance to succeed on a project could also change. Data between these two versions were collapsed for the present analysis based on prior findings of no task performance differences.<sup>4</sup> However, where things differ in this version from the single-reversal PRL, is 2-fold. The first difference is this idea of *multiple* reversals – every 40 trials participants were provided a break, following which probabilities automatically reassigned. In addition, there were performance-based reversals – after every 9 out of 10 consecutive selections of the highest reinforcement probability deck, unbeknownst to the participants, the underlying probabilities of receiving rewards were altered. The second difference is this idea of a *shift* (that occurs halfway through the experiment) – the contingencies of rewards initially follow a 90-50-10 pattern for the first 80 trials but switch to a 80-40-20 for the last 80 trials, making it more difficult to distinguish whether a loss was due to chance or a change in the best deck (or partner). Thus, at this halfway mark, individuals experience both a reversal and a shift.

## Quantification and statistical analysis

### Behavioral choice analysis

We quantified strategies of behavior on the task similar to prior work.<sup>3</sup> The likelihood that participants would choose alternative options after positive feedback (*win-switch*) and select the same option after negative feedback (*lose-stay*) was measured. The rate of win-switch behavior was calculated as the total number of trials in which participants switched after positive feedback divided by the number of trials in which they received positive feedback. Similarly, the rate of lose-stay behavior was calculated as the total number of trials in which a participant persisted after negative feedback divided by the total negative feedback trials. The rates of win-switch and lose-stay behavior on each trial was calculated using a lagged moving average approach with a 20-trial window. This means that for each binned trial, we computed the average rate based on the previous 20 trials. This lagged moving average is depicted in [Figure 2B](#) (and [4A](#), [4B](#)), where *reversal* (and *shift* in the multi-reversal version) is marked at the mid-way period of the experiment.

### Computational modeling

We used a mean-reverting HGF model.<sup>6</sup> In this model ([Figure 1B](#)), the component that incorporates task beliefs is called 'the perceptual model'. The component that governs how beliefs are converted into choices is called 'the response model'. The perceptual model has three hierarchical layers of belief about the task. The layers interact and influence one another through learning rate parameters. At the highest level (i.e., level 3), the model captures beliefs about changes in the task environment (i.e., how are values of the choices changing over time?). Level 2 characterizes beliefs on reward probabilities (i.e., the tendency of a choice to be rewarding). Level 1 characterizes task reward feedback (i.e., win or loss). These three levels of belief are then integrated and fed through a sigmoid response function to produce a decision (i.e., whether to either stay with the same option or switch to a different one). To calculate subject-specific estimates of belief about volatility ( $m_3$ ) and belief about reward value learning ( $\omega_2$ ), we fit the Hierarchical Gaussian Filter (HGF) to the three-choice PRL data.<sup>57</sup> We estimated these perceptual parameters for the first (monkeys: trials 1–150; single-reversal humans: trials 1–30, multi-reversal humans: trials 1–80) and second (monkeys: trials 151–300; single-reversal

humans: trials 31–60, multi-reversal humans: trials 81–160) halves of the task. It is important to note that we did not model these two halves separately; we used the estimated priors from the first half of the task as priors for the second half, thereby accounting for the mid-way contingency shift in a unified model approach (see Figure S10). Each of the agent's choice (i.e., option 1, 2, or 3) and outcomes (reward or no reward) were entered as separate column vectors with rows corresponding to trials. Reward was encoded as '1', no reward as '0', and choices as '1', '2', or '3'. We modified the perceptual model configuration file (*tapas\_hgf\_ar1\_binary\_mab\_config.m*) to reflect the 'winning' model which uses a mean-reverting HGF perceptual model (refer to Table S1a in<sup>6</sup>) and a softmax-mu03 decision model (*tapas\_softmax\_mu3.m*). **Parameter recovery.** We fit our HGF model to the observed PRL task data to estimate parameters that describe how an agent's (monkey or human) beliefs update throughout the experiment. Our approach to demonstrating parameter recovery involved attempting to recover the parameter estimates that were originally determined from the true choice data using simulated task data. For the simulation, task performance for each agent was generated using the true model parameter estimates, the perceptual model, and the decision model. The values applied for the perceptual and decision configuration files were drawn from a prior study (specifically, the M6 – winning model; refer to their Table S1a).<sup>6</sup> A dummy 'zerth' trial served as our starting point to simulate the initial response and outcome. The agent's responses (choice 1, 2, or 3) were simulated based on the softmax function, considering the true model parameters and the experienced, simulated outcomes. Meanwhile, the outcomes (whether a reward of 1 or no reward of 0) were generated through a Bernoulli distribution, reflective of the agent's reward schedules (for instance, the stable schedule for monkeys). For every agent, we created a set of simulated responses and outcomes. This entire process was repeated for  $i = 10$  iterations, yielding 10 unique sets of simulated response-outcome data per agent. Subsequently, we refitted the HGF to the simulated data, which produced a set of belief parameters for each agent during every iteration. These simulated (or recovered) model parameters were correlated with the observed (or estimated) model parameters (see Figure S11) for the equilibrium value of belief volatility ( $m_3$ ) and beliefs about value learning ( $\omega_2$ ) parameters, to assess performance of parameter recovery across the different variations of the PRL task for both monkeys and humans. Furthermore, plotting the simulated win-switch and lose-stay rates, at the group level, recapitulated the observed differences between lesion groups (see Figure S12). The HGF toolbox v5.3.1 is freely available for download in the TAPAS package at <https://translationalneuromodeling.github.io/tapas>. We installed and ran the package in MATLAB and Statistics Toolbox Release 2016a (MathWorks\*, Natick, MA).

## STATISTICAL ANALYSIS

### General

Statistical analyses and effect size calculations were performed with an alpha of 0.05 and two-tailed  $p$ -values in RStudio: Integrated Development Environment for R, Version 4.3.0.

### GLMMs

To combat non-normality and random effects, we employed generalized linear mixed models (GLMMs) to identify significant group differences (lesion group or paranoia group) in behavior and beliefs. We used a *binomial GLMM* for behavior, incorporating the counts of win-switch and win-stay for win-switch rate and lose-stay and lose-switch for lose-stay rate, and a *Gaussian GLMM* for the HGF belief parameters. Significant interaction effects in omnibus models were resolved by generating lower-order models.

### Permutation tests

We shuffled the lesion groups to create random permutations of the data. For each permutation, we calculated the test statistic (i.e., difference of means) between the win-switching rates of the MDmc-lesioned monkeys and the non-lesioned control monkeys for post-reversal and pre-reversal trials; mean difference =  $(\mu_{post,mdmc} - \mu_{pre,mdmc}) - (\mu_{post,control} - \mu_{pre,control})$ . We repeated this process for  $n = 100$  permutation to build the empirical null distribution. We calculated the  $p$ -value as the proportion of permuted test statistics that are greater than or equal to the observed test statistic (in our case, this was 0.169). The  $p$ -value represents the probability of obtaining a test statistic as extreme as the observed one, assuming the null hypothesis ( $H_0$ : no difference between groups). If  $p$ -value  $< 0.05$ , we reject the null hypothesis and conclude that there is a significant difference in win-switching between the lesion groups after the reversal. Lastly, Interquartile Range (IQR) tests were used to label data as outliers.