

# Neurocomputational evidence of sustained Self-Other mergence after psychedelics

Pablo Mallaroni<sup>1</sup>, Natasha L. Mason<sup>1</sup>, Katrin H. Preller<sup>2</sup>, Adeel Razi<sup>3,4,5</sup>,  
Sam Ereira<sup>6\*</sup> & Johannes G. Ramaekers<sup>1\*</sup>

\*These authors contributed equally to this work

<sup>1</sup>Department of Neuropsychology and Psychopharmacology, Faculty of Psychology and Neuroscience, Maastricht, the Netherlands

<sup>2</sup>Department of Adult Psychiatry and Psychotherapy, Psychiatric University Clinic Zurich, University of Zurich, Zurich, Switzerland

<sup>3</sup>Turner Institute for Brain and Mental Health, School of Psychological Sciences, Monash University, Clayton, Victoria, Australia

<sup>4</sup> Queen Square Institute of Neurology, University College London, London, UK.

<sup>5</sup>CIFAR Azrieli Global Scholars Program, CIFAR, Toronto, Ontario, Canada

<sup>6</sup>Wolfson Institute of Population Health, Queen Mary University of London, London, United Kingdom

## Corresponding author email

[p.mallaroni@maastrichtuniversity.nl](mailto:p.mallaroni@maastrichtuniversity.nl); [j.ramaekers@maastrichtuniversity.nl](mailto:j.ramaekers@maastrichtuniversity.nl)

## Conflict of interest:

Johannes G. Ramaekers is a scientific consultant to GH Research who have no involvement in the preparation or conception of this manuscript or related data. K.P. is currently an employee of Boehringer Ingelheim GmbH & Co KG, and is currently Chief Scientist and on the Board of Directors for the Heffter Research Institute, and scientific advisor for the MIND foundation. All other authors have no relevant conflicts of interest to declare.

## Funding

Johannes G. Ramaekers acknowledges financial support from the Dutch Research Council (NWO, grant number 406.18. GO.019). NLM is financially supported by the Dutch Research Council (NWO, grant number VI.Veni.231G.011).

## Keywords

Psychedelics, social cognition, reinforcement learning, computational psychiatry, fMRI

## Abstract

Mental illness is often characterised by a maladaptive sense of *self*. The neurobiological basis of Self-Other distinction may provide targets for therapeutic interventions. Psychedelics alter the experience of selfhood, but the neurocomputational mechanism is unclear. We used a computationally-informed behavioural assay to investigate whether psychedelics disrupt Self-Other boundaries in belief formation. In a double-blind, crossover design, 22 participants received placebo, psilocybin or 2C-B (2,5-dimethoxy-4-bromophenethylamine). The next day, we fitted reinforcement learning models to probabilistic false-belief task behaviour, yielding objective Self-Other distinction measures. Compared to placebo, psychedelics induced a state of Self-Other mergence, associated with a multivariate signal of sustained psychosocial wellbeing. Effective-connectivity estimates from resting-state fMRI showed that Self-Other boundary disruption was associated with reduced inhibitory tone from right temporoparietal junction to dorsomedial prefrontal cortex. We show that psychedelics quantifiably act on the neural basis of Self-Other distinction, offering potential routes to precision therapeutics in psychedelic psychiatry.

## Introduction

Humans have evolved the ability to represent the contents of each others' minds. This ability, to *put myself into your shoes*, is often referred to as Theory of Mind<sup>1,2</sup> or mentalising<sup>3</sup>. It allows humans to make predictions about each other<sup>4,5</sup>, and may be a fundamental process behind the emergence of human social networks and societies<sup>6,7</sup>.

A critical component of Theory of Mind is the selective attribution of a mental state, such as a belief or a desire, to a specific agent, such as Self or Other. Neuroimaging and electrophysiological research show that humans can achieve this by simulating the neural computations of other people in segregated, agent-specific, neural circuits<sup>8-12</sup>. However, Self-Other distinction is a flexible process that adapts to different social environments<sup>11</sup>. On one hand, it would be unhelpful for an individual to persistently conflate the contents of their own mind with the contents of another person's mind<sup>13</sup>. On the other hand, empathic responding, where another person's mental states are experienced as one's own, can be adaptive to facilitate social bonding and mutual understanding<sup>14-19</sup>.

As well as enabling complex social behaviour, adaptive Self-Other distinction is necessary for furnishing us with a coherent sense of Self<sup>20,21</sup>. It is therefore unsurprising that aberrant mentalising is a transdiagnostic feature of several mental illnesses<sup>10,22</sup>, including autism spectrum disorder (ASD)<sup>23,24</sup>, schizophrenia<sup>25,26</sup>, borderline personality disorder (BPD)<sup>27,28</sup>, depression<sup>29</sup> and dementia<sup>30</sup>. People with these conditions experience difficulties in maintaining social connections and a coherent sense of Self, highlighting the potential value in identifying interventions that selectively modulate the neurobiological basis of Self-Other distinction.

Serotonergic psychedelics are one class of drugs that can transiently and profoundly alter a person's experience of selfhood. Compounds such as psilocybin, lysergic acid diethylamide (LSD), mescaline, 2,5-dimethoxy-4-bromophenethylamine (2C-B), and dimethyltryptamine (DMT) reliably induce a blurring of the boundary between Self and non-Self<sup>31-33</sup>, likely via agonism at 5-HT<sub>2A</sub> receptors<sup>34,35</sup>. This acute state, known as *ego dissolution*<sup>36-38</sup>, is thought to underpin the observed therapeutic effects of psychedelics in disorders marked by aberrant Self-referential processing<sup>39-42</sup>. In the days following the acute phase, there is often a subacute *afterglow* period of sustained psychosocial wellbeing, characterised by enhanced empathy, suggestibility, and a sense of social connectedness<sup>43-48</sup>.

The neurocomputational basis of ego dissolution and the ensuing afterglow period are not well understood. Experimental work on Self-Other processing under psychedelics has largely relied on subjective reports. Where behaviour has been assessed, studies consistently show that psychedelics induce a blurring of the Self-Other boundaries of affective and sensory states<sup>49-51</sup>. However, there has been little work to generate behavioural data amenable to computational modelling. Computationally-informed behavioural tasks, such as that used in the current study, elicit behavioural data that can test explicit models of underlying cognitive

processes. By moving from subjective reports to objective computational metrics, it may be possible to describe and explain psychedelic alterations in Self-Other processing within established neurocomputational frameworks.

Here within a double-blind, placebo-controlled, within-subject crossover design, we tested whether the psychedelic compounds, psilocybin and 2C-B, could induce enduring changes in Self-Other distinction with respect to prediction errors, fundamental computational units in learning and belief formation<sup>52,53</sup>. We employed a probabilistic false-belief task (pFBT) to record behavioural data and estimate subject-specific behavioural parameters. These parameters provide an objective measurement of a person's propensity to process their environment using agent-specific prediction errors or agent-independent prediction errors<sup>10,11,27</sup>. Agent-specific prediction errors are sensory surprise signals that contain information not only about the magnitude of the surprise, but also about the identity of the person who is experiencing the surprise (e.g. *me* or *you*), allowing for Self-Other distinction in belief states. Agent-independent prediction errors encode the surprise signal without any information about agent identity, allowing for a Self-Other mergence in belief formation. The pFBT takes inspiration from the field of computational psychiatry<sup>22,54-56</sup>, which seeks to identify objective biomarkers, by explaining cognition and mental illness using mathematical models. These models are constrained by quantitative parameters, estimated from a person's behavioural, neural or other physiological data.

To complement these behavioural measures, participants were scanned with ultra high field 7 Tesla resting-state fMRI (rsfMRI) during peak drug effects. By estimating effective connectivity using spectral dynamic causal modelling (DCM), we aimed to identify acute drug-induced connectivity changes, between brain regions known to be involved in Self-Other processing, that were associated with behavioural effects.

By employing tools from computational psychiatry, our aim was to assess whether and how psychedelic compounds alter an objective measure of Self-Other distinction in learning and belief formation. We predicted that psychedelics shift people into a cognitive state where they were more likely to conflate the mental states of Self and Other, and that this altered state is associated with the broad psychosocial wellbeing previously observed in people treated with psychedelics.

## Results

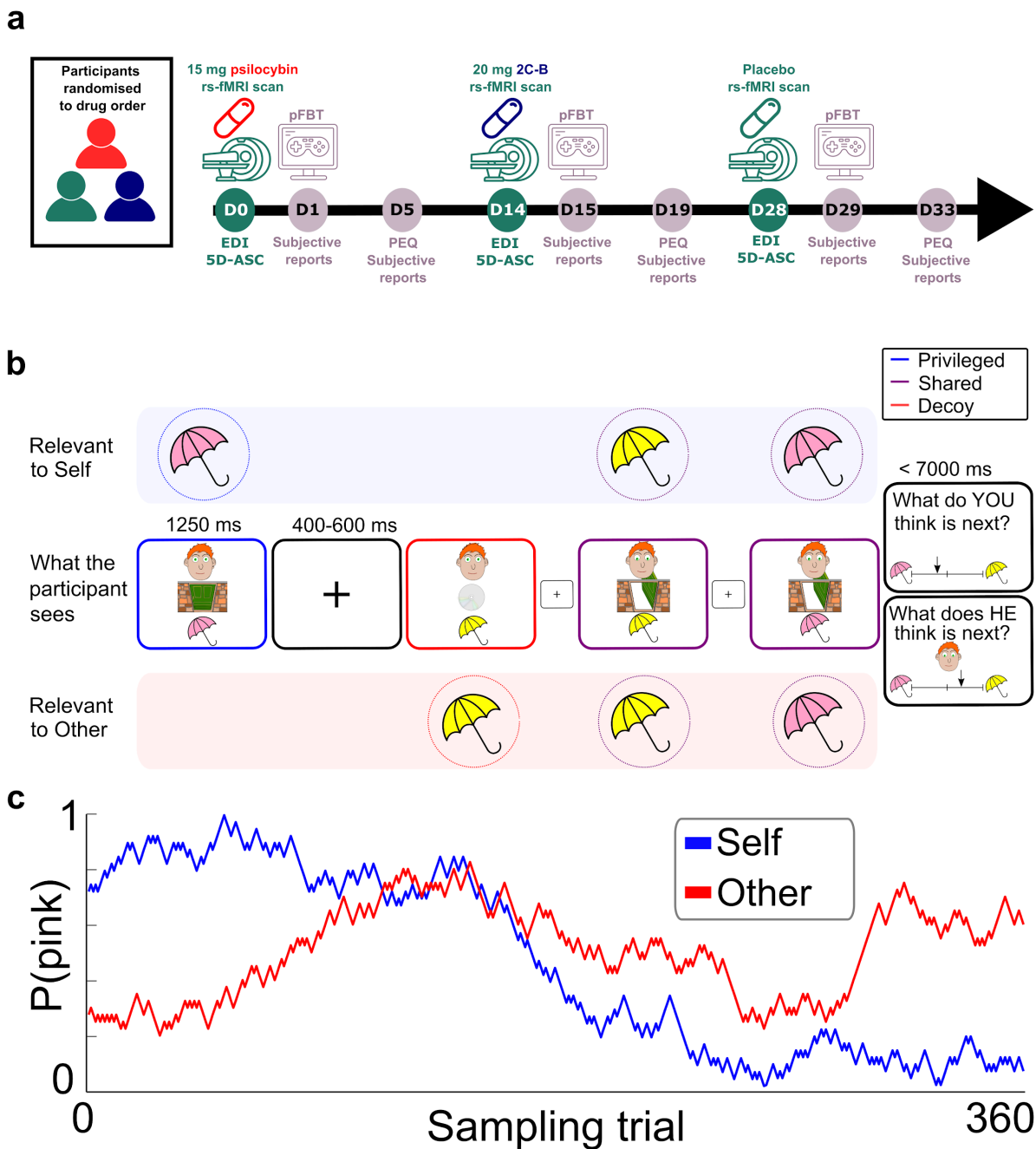
### ***Psychedelics do not impair performance on the probabilistic false belief task***

N = 22 participants were treated with matched doses of 15mg psilocybin, 20mg 2C-B, and placebo, in a randomised order, with two weeks between each dosing session (Fig. 1a). One day after each dosing session, we administered a probabilistic false belief task (pFBT)<sup>10</sup>, in order to measure components of Self-Other processing in the subacute phase following psychedelic treatment.

The pFBT (Fig. 1b) presented participants with a sequence of binary samples, in this case yellow sun-shades or pink umbrellas. Each sample was paired with a second image indicating whether or not a fictional other person could see the sample. On some trials, this image indicated that the information *could* be seen by the other person, but that it was misleading. The probability of a trial revealing a pink umbrella gradually drifted throughout the task along a random walk (Fig. 1c, blue line). Participants were instructed to infer and keep track of this drifting probability so that they could make accurate predictions about what the outcome might be on the next trial. By combining a randomised mixture of trials where the other agent could see some outcomes and not others, and also saw some misleading outcomes, the other person's belief about the probability also followed a random walk, decorrelated from the true underlying probability (Fig. 1c, red line). Participants were instructed to also infer and keep track of the other person's false belief about the probability. Behavioural data was collected on probe trials, in which participants reported their estimates of the probability, or their estimates of the other person's false belief about the probability.

Accuracy on the pFBT was calculated as the correlation between participants' reports on probe trials, and the true underlying probabilities generating outcomes relevant to Self and Other (see Methods). Most participants performed better than chance, in all three testing sessions (Fig. 2a), demonstrating a good understanding of the task. However, two participants performed below chance-level in multiple sessions, and they were excluded from all further analyses.

For the remaining N = 20 participants, we fitted a linear mixed-effects model, predicting accuracy from condition (placebo, psilocybin or 2C-B) with random intercepts estimated for each participant, to account for repeated measures. Overall, accuracy was significantly higher than chance ( $F_{1,57} = 60.97, p < 0.001$ ) and there was no effect of condition on accuracy ( $F_{2,57} = 1.69, p = 0.19$ ) These results suggest that during the subacute phase following psychedelic administration, participants' comprehension of and ability to perform the pFBT was not impaired.



**Figure 1 - Experimental design**

**a**, This was a randomised crossover study where participants received psilocybin, 2C-B and placebo, in a randomised order, with a 2 week washout between each dosing day. On each dosing day (days 0, 14 and 28) rsfMRI data was collected and participants were assessed using the Ego Dissolution Inventory (EDI) and 5-Dimensions of Altered States of Consciousness (5D-ASC) questionnaire. The day after each dosing day (days 1, 15 and 29), participants played a probabilistic false belief task (pFBT) and the following subjective self-report questionnaires

were administered: Inclusion of Other in the Self (IOS), Positive and Negative Affect Schedule (PANAS), Social Connectedness Scale–Revised (SCS-R). Five days after each dosing day (days 5, 19 and 33) the same battery of questionnaires was administered, in addition to the persisting effects questionnaire (PEQ). **b**, pFBT task design. Participants observed a series of *sampling trials* where each trial revealed either a pink umbrella (with probability  $P$ ) or a yellow sun-shade (with probability  $1-P$ ). Participants were required to infer and keep track of  $P$  as well as another agent's false belief about  $P$ . On *privileged* trials the information was hidden from the other agent. On *shared* trials, the information was available to the other agent. On *decoy* trials, the information was available to the other agent, but unbeknown to the other agent, the information was misleading. Participants were intermittently probed to report their estimate of  $P$  or the estimate of the other agent's belief about  $P$ . **c**, The pFBT was designed such that the drifting variable  $P$  (blue) and the other agent's belief about  $P$  (red) were not correlated.

### ***Psychedelics promote Self-Other mergence in prediction errors***

We fitted a family of reinforcement-learning style models to participants' behavioural data from probe trials (see Methods). These models shared the assumption that, on each trial, participants updated their own belief about the underlying probability, and also updated a simulation of the other agent's belief using a parallel and simultaneous belief-update. This family of models has been shown to approximate behaviour in the pFBT well<sup>10,11,27</sup> and also furnish computational variables that can explain brain activity during task engagement<sup>10,11</sup>.

The models were nested, and they differed in whether they allowed the parallel belief-updates for Self and Other to be parameterised independently or jointly, and whether or not they allowed for cross-agent leakage. In one version of the model, the parallel belief-updates were entirely segregated, reflecting a perfect implementation of Self-Other distinction. In another version of the model, the participant's prediction errors were leaky, leading them to erroneously update their simulation of the other agent's belief. This *idiocentric updating*<sup>27</sup>, reflects a situation where *my own beliefs influence my simulation of your beliefs*. In an alternative version of the model, the simulated prediction errors of the other agent were leaky, leading to erroneous updating of the participant's own beliefs. This *allocentric updating*<sup>27</sup>, reflects a situation where *my simulation of your beliefs influences my own beliefs*.

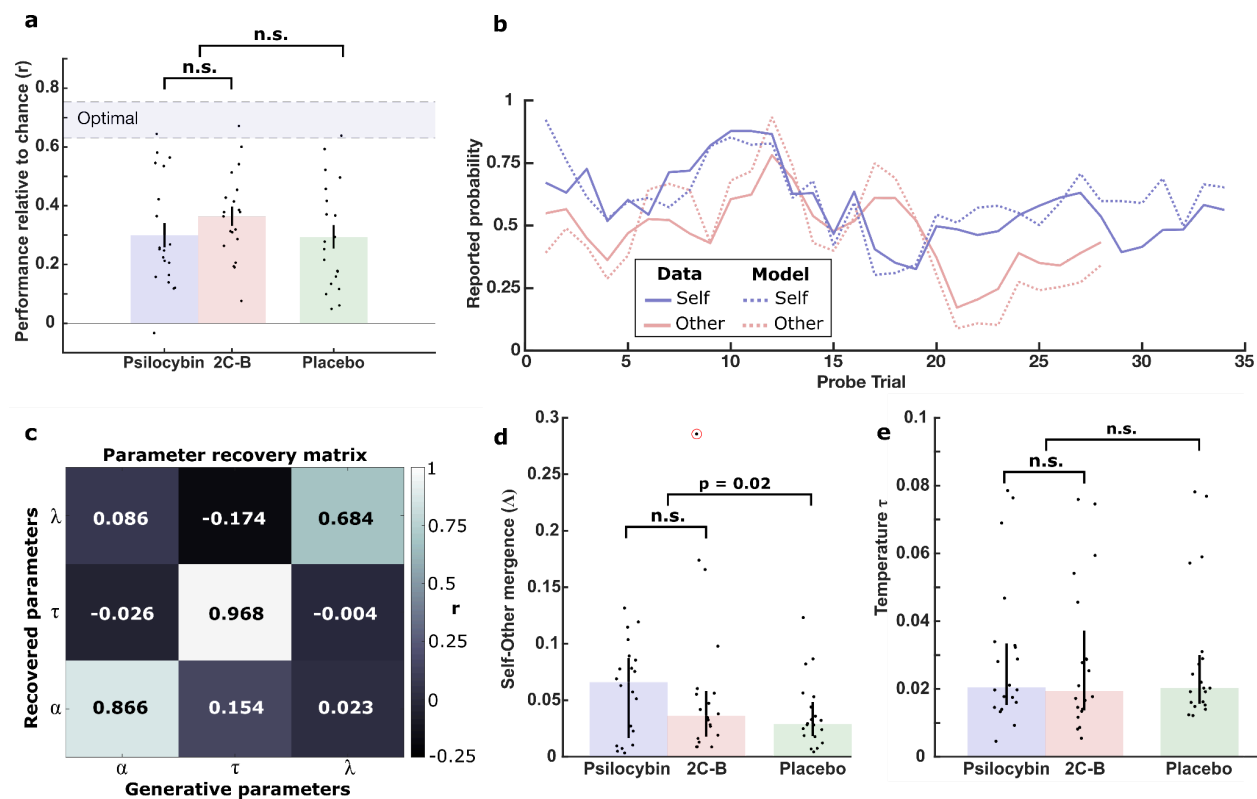
In a Bayesian model comparison of 36 models (see Methods) we found that the best model to explain behaviour across all three testing sessions was a three parameter model. The model specifies a learning rate ( $\alpha$ ), which governs the rate at which new information is incorporated into old beliefs (i.e. sensitivity to surprise), a decision temperature ( $\tau$ ), which governs behavioural stochasticity in decision-making, and a Self-Other leak parameter ( $\lambda$ ), governing bi-directional leakage from Self-to-Other and Other-to-Self, capturing both idiocentric and allocentric updating. This model explained true choice data well (Fig. 2b) with  $r^2 = 0.41 \pm$

0.24 (median  $\pm$  IQR). The three parameters were recoverable and separately identifiable (Fig. 2c).

We were specifically interested in whether psychedelic compounds alter the sensitivity to leaked prediction errors. As described in Methods, this is captured by the parameter  $\Lambda$ , which is the product of  $\alpha$  and  $|\lambda|$  and reflects the extent to which *my belief is influenced by my simulation of your prediction errors* and *my simulation of your belief is influenced by my own prediction errors*. In other words,  $\Lambda$  is a metric that quantifies Self-Other mergence in belief updating. We did not predict any differences between 2C-B and psilocybin. Consistent with this, in a two-tailed Wilcoxon sign rank test we found no difference in Self-Other mergence ( $\Lambda$ ) between the post-psilocybin session and the post-2C-B session ( $z = 0.22$ ,  $p = 0.82$ ).

To capture the overall effect of either psychedelic compound, we fitted a linear mixed-effects model predicting Self-Other mergence ( $\Lambda$ ) from whether the session involved drug or placebo (binary). The model included random intercepts and random condition-specific slopes for each participant, allowing the effect of each drug (placebo, 2C-B, psilocybin) to vary by subject. A likelihood ratio test favoured this model compared to a nested model with only random intercepts ( $\chi^2(5) = 16.80$ ,  $p = .0049$ ). We found that psychedelics significantly increased Self-Other mergence relative to placebo ( $F_{1,57} = 5.77$ ,  $p = 0.02$ ). In a sensitivity analysis, we refitted the model after excluding one participant's 2C-B session, which had an anomalously high Self-Other mergence metric (Fig. 2d); results were consistent ( $F_{1,57} = 4.65$ ,  $p = 0.035$ ).

Fitting another linear mixed effects model predicting decision temperature (Fig. 2e) instead of Self-Other mergence, we found no effect of psychedelics on decision temperature ( $F_{1,57} = 0.54$ ,  $p = 0.82$ ). There was also no drug effect on learning rate ( $F_{1,57} = 0.54$ ,  $p = 0.59$ ), suggesting that it was not simply enhanced sensitivity to surprise that was driving the effect on Self-Other mergence ( $\Lambda$ ). These results suggest that these psychedelic compounds induce a specific effect on belief formation, wherein Self-Other mergence is enhanced in the subacute phase following drug administration.



**Figure 2 - Analysis of pFBT behavioural task data**

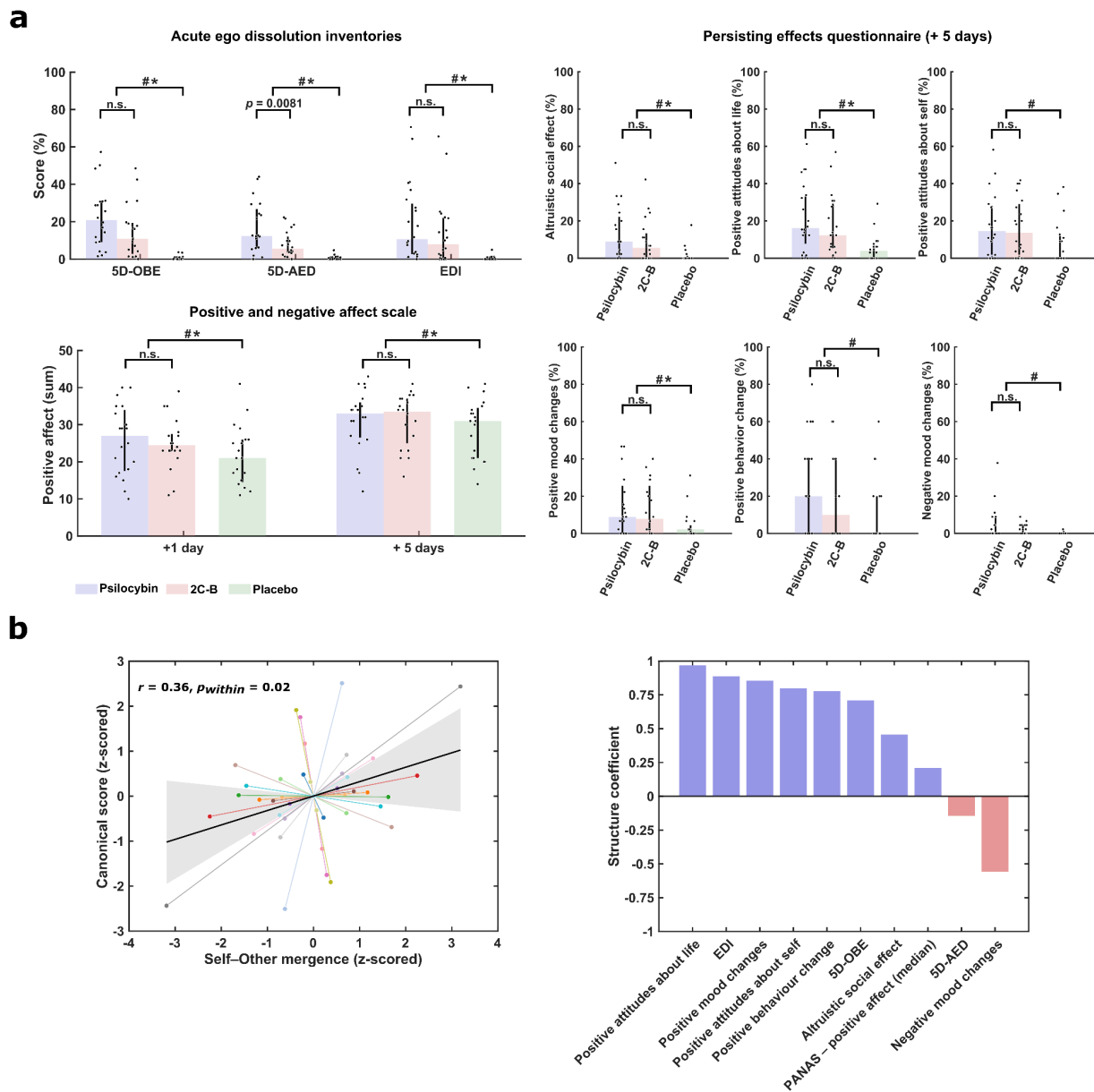
**a**, Task performance, calculated as the correlation between participant reports and true underlying probabilities generating *sampling trials*. As described in Methods, “optimal” performance was estimated by simulating observers with perfect memory and no behavioural stochasticity. There was no difference in performance between post-2C-B sessions and post-psilocybin sessions. There was also no difference in performance between post-placebo sessions and post-psychedelic sessions. The bar graph shows the median +/- interquartile range (IQR) and each dot represents a different participant.  $N = 20$ . **b**, Exemplar model fit for a single participant session. To avoid cherry-picking, the participant session with the median log-model evidence across all participant sessions was selected for visualisation. Model predictions are illustrated for *Self probe* and *Other probe* trials separately. **c**, Parameter recovery matrix. Pearson correlation coefficients between generative parameters of 1000 simulated participants and recovered parameters demonstrate that the three model parameters ( $\alpha$ /learning rate,  $\tau$ /temperature,  $\lambda$ /leak) are separately identifiable. **d**, Self-Other mergence ( $\Lambda$ , i.e. product of  $\alpha$  and  $|\lambda|$ ) was significantly lower following placebo administration than following psychedelic administration. The bar graph shows the median +/- IQR and each dot represents a different participant. One datapoint for the 2C-B session, circled in red, was identified as an outlier. In a sensitivity analysis where this participant was excluded, statistical results were consistent.  $N =$

20. e, There was no difference in  $\tau$ /temperature between post-placebo and post-psychedelic sessions. The bar graph shows the median +/- IQR and each dot represents a different participant. N = 20.

***Psychedelic-induced Self-Other mergence predicts (sub)acute subjective effects.***

To assess whether drug administration was accompanied by changes in subjective changes in psychosocial wellbeing, we characterised acute ego dissolution using the Ego Dissolution Inventory (EDI) and the 5-Dimensions of Altered States of Consciousness (5D-ASC), and follow-up measures using the Persisting Effects Questionnaire (PEQ), Positive and Negative Affect Scale (PANAS), the Inclusion of Other in the Self (IOS) scale, and the Social Connectedness Scale–Revised (SCS-R) (see Fig. 3a, Methods and Supplementary Table 1).

Acutely, all ego dissolution dimensions were significantly elevated under both 2C-B and psilocybin, with the strongest effect in 5D-ASC oceanic boundlessness ( $F_{1,39} = 33.08$ ,  $p < .0001$ ). Psilocybin additionally produced greater elevations in dysphoric aspects of ego dissolution (5D-ASC anxious ego dissolution) relative to 2C-B ( $p = .0134$ ,  $d = 0.94$ ). Subacutely, a robust main effect of drug was observed across all positive valence dimensions of the PEQ, with altruistic prosocial effects showing the clearest enhancement at +5 days ( $F_{1,39} = 18.37$ ,  $p < .0001$ ) with no difference between compounds ( $p = .168$ ). No effects emerged for antisocial negative effects ( $p = .691$ ). Within negative valence measures, only negative mood showed a significant drug effect ( $F_{1,39} = 4.44$ ,  $p = .0415$ ), driven by elevated scores under psilocybin relative to placebo ( $p = .0215$ ,  $d = 0.88$ ), but not 2C-B ( $p = .163$ ). PANAS positive affect paralleled the PEQ results, with a main effect of drug ( $F_{2,95} = 4.56$ ,  $p = .0128$ ) but no interaction with follow-up time ( $p = .713$ ), nor significant differences between drugs ( $p = .939$ ). By contrast, PANAS negative affect, IOS, and SCS-R showed no significant drug or drug  $\times$  time effects (minimum  $p = .26$ ). These results indicate that psychedelics consistently enhance psychosocial wellbeing subacutely, with psilocybin potentially producing greater dysphoric ego dissolution and transient increases in negative mood.



**Figure 3. (Sub)acute psychosocial effects and their multivariate association with Self-Other merger.**

**a**, Subjective effect outcomes. Acute ego dissolution was assessed using the retrospective 5-Dimensions of Altered States of Consciousness questionnaire (5D-ASC; oceanic boundlessness [OBE], anxious ego dissolution [AED]) and the Ego Dissolution Inventory (EDI). Positive and Negative Affect Scale (PANAS) scores were collected at 1 and 5 days post-session, and persisting effects were measured using the Persisting Effects Questionnaire (PEQ) subscales at +5 days. Both psychedelics were associated with higher scores than placebo on several outcomes (see

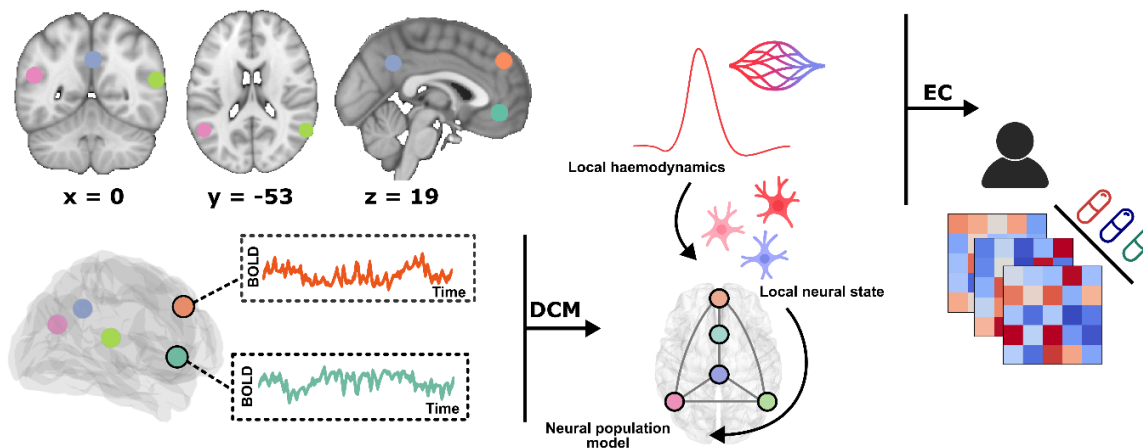
brackets), with significance assessed using Tukey-corrected pairwise comparisons. Bar graphs show the median +/- IQR; each dot represents one participant (N = 20). Significance ( $p < 0.05$ ) per contrast is represented as the following; # = psilocybin vs placebo, \* = 2C-B vs placebo, n.s. = not significant. **b**, Left: Multivariate association between Self-Other mergence and subjective outcomes. The scatter plot shows the repeated-measures association between Self-Other mergence and the canonical (sub)acute psychosocial response profile. Points are within-subject, drug-aware residuals (z-scored); segments connect distinct sessions from the same participant. The black line is the shared repeated-measures correlation fit and the shaded band shows its 95% CI (subject-level bootstrap). Because residualisation includes an intercept and drug term, placebo observations plot near the origin. Right: The bar chart depicts structure coefficients (correlation between each outcome and the canonical variate), indicating the relative contribution of each measure to the multivariate association.

We examined whether individual differences in psychedelic-induced Self-Other mergence were associated with a coherent (sub)acute psychological response profile. To test this, we conducted a drug-aware permutation MANOVA (10,000 permutations), residualising average drug effects so that only between-participant variance was assessed. We included all subjective outcomes that showed a significant drug effect. A significant overall association between Self-Other mergence and the outcome set (Pillai's trace = 0.64,  $p_{\text{perm}} = 0.0012$ ) was identified. Canonical analysis (Fig. 3b) showed that Self-Other mergence correlated with the first canonical variate ( $r_{\text{in-sample}} = 0.36$ ,  $p_{\text{perm}} = 0.0041$ ). Importantly, a repeated-measures correlation confirmed that within-participant fluctuations in Self-Other mergence tracked corresponding changes in the canonical variate ( $p_{\text{within}} = 0.02$ ). The effect also generalised under cross-validation ( $r_{\text{out-of-sample}} = 0.32$ ,  $p_{\text{perm}} = 0.044$ ; Fig. 3b), indicating that the canonical relationship was stable and predictive. Canonical weights and auxiliary descriptive univariate regression results (none surviving FDR correction) are reported in the Supplementary Materials. Structure coefficients indicated strong positive correlations for positive-valence and ego dissolution measures, and a strong negative correlation for negative mood and ego dissolution changes. Altogether, these findings suggest individuals with greater increases in Self-Other mergence under psychedelics also reported more positive psychosocial responses, independent of general drug effects.

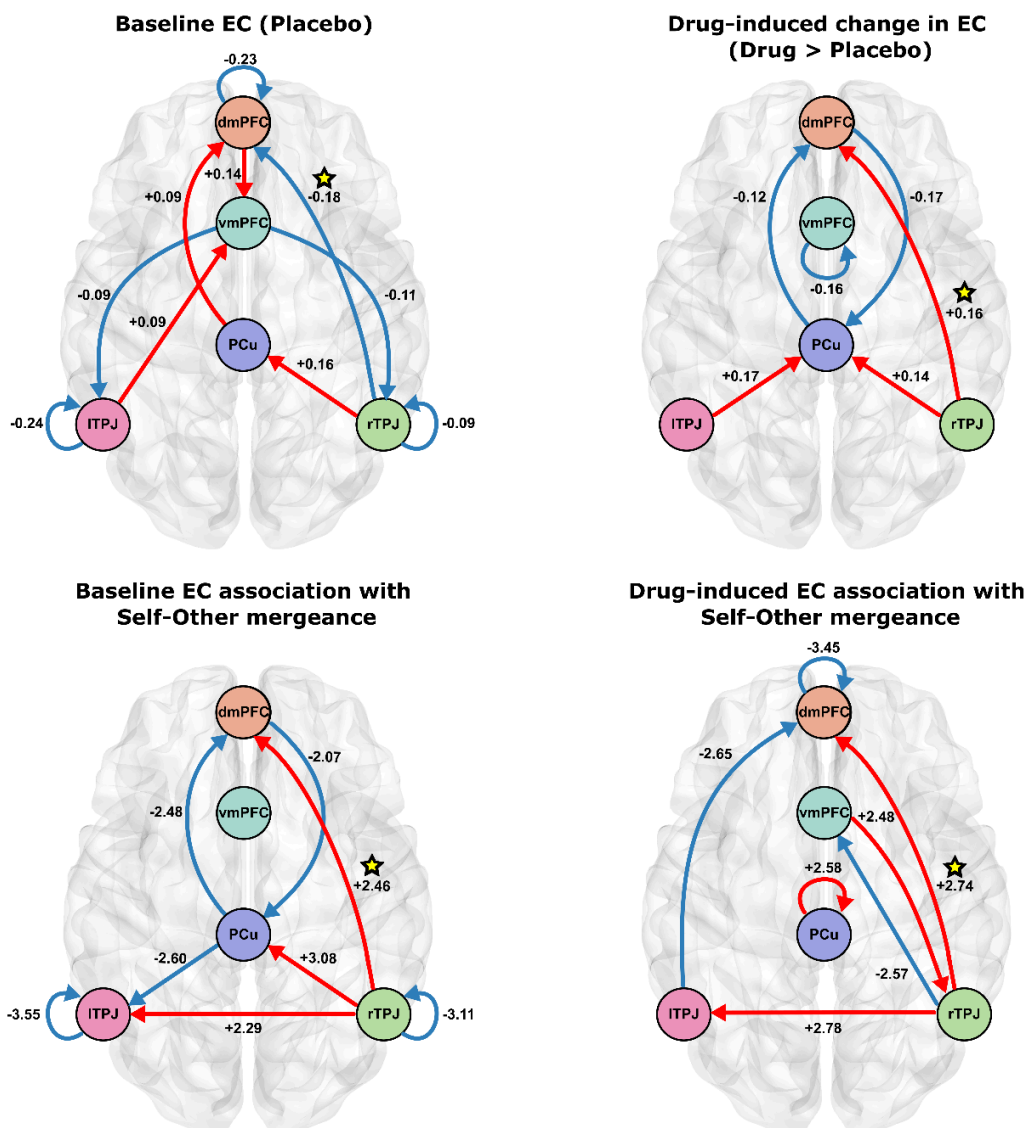
### ***Acute effective connectivity changes underlie persisting Self-Other mergence***

Psychedelics disrupt the functional organisation of association cortices involved in self-referential processing, to produce a transient loss of coherent self-awareness<sup>57</sup>. To examine how such disruptions might underlie persisting Self-Other mergence, we applied spectral dynamic causal modelling (DCM) to peak-effect 7T rsfMRI data. Focusing on a Theory of Mind (ToM) network (Fig. 4a, Methods), we tested whether psychedelics alter the organisation of effective connectivity (EC) between core ToM regions in a manner consistent with altered belief updating.

**a**



**b**



#### **Figure 4. Dynamic causal modelling of effective connectivity (EC) under psychedelics and its association with Self-Other mergence**

**a**, Theory-of-mind (ToM) network regions of interest (ROIs) are shown on canonical brain slices, with schematic BOLD time series for two ROIs. rsfMRI data acquired approximately 110 min post-administration were fitted using spectral dynamic causal modelling (DCM), which estimates directed effective connectivity (EC) between regions by modelling local haemodynamics and neural population dynamics. **b**, Top: first-level DCMs (placebo, psilocybin, 2C-B) were combined into participant-level models estimating baseline EC (placebo) and drug-induced change ([drug – placebo]) using parametric empirical Bayes (PEB). Bottom: a third-level PEB-of-PEBs tested group-average effects and associations with changes in Self-Other mergence ( $\Delta$ ). For our baseline model, blue arrows denote inhibitory, red arrows excitatory influences; arrow labels indicate direction. For drug–placebo contrasts, arrows indicate increases (red) or decreases (blue) relative to baseline, not absolute valence. Effect sizes are reported posterior expectations in Hz (self-connections are log-scaled, inherently inhibitory). Connections are only shown if the posterior probability of being non-zero is  $>0.99$ ; ★ marks effects convergent across contrasts. Results at posterior probability and EC difference matrices at  $>0.50$  are provided in Supplementary Materials. Abbreviations: dmPFC = dorsomedial prefrontal cortex, vmPFC = ventromedial prefrontal cortex, PCu = precuneus, lTPJ = left temporoparietal junction, rTPJ = right temporoparietal junction.

Bayesian model reduction and averaging revealed very strong evidence (posterior probability  $> 0.99$ ) for six connectivity parameters differing between drug and placebo conditions (Fig. 4b). At baseline (placebo), the network exhibited a mixture of excitatory and inhibitory influences. We will illustrate how to interpret these maps by taking one example connection, rTPJ→dmPFC. At baseline (placebo), group-level DCM estimates indicated an inhibitory connection from rTPJ to dmPFC. Under psychedelics, this connection exhibited an attenuation of inhibitory influence, consistent with relative disinhibition. Across participants, those with a more excitatory baseline connection (i.e. less inhibitory) exhibited larger drug-induced increases in Self-Other mergence ( $\Delta$ ). Furthermore, participants who showed the strongest excitatory shift in this connection under psychedelics also displayed the strongest behavioural evidence of increased Self-Other mergence 24 hours later. Exploratory analyses suggested that this attenuation was positively associated with the multivariate psychosocial response profile described above (see Supplementary Materials). Together, these findings suggest that the rTPJ→dmPFC pathway is normally inhibitory, that psychedelics attenuate this inhibitory tone, and that the extent of this attenuation predicts the magnitude of behavioural change in Self-Other processing.

## Discussion

We show that 2C-B and psilocybin induce a state of Self-Other mergence. More specifically, there is an increase in sensitivity to agent-independent prediction errors, which encode a generic surprise signal, divorced from agent-identity i.e. (*me* or *you*). This shift, from Self-Other distinction to Self-Other mergence, persists beyond the acute drug phase, into the subacute *afterglow* period.

Human neuroimaging studies<sup>8,10,11,58,59</sup> and electrophysiological studies in humans<sup>9</sup> and non-human primates<sup>60-64</sup> have found converging evidence that computational units of learning and belief formation can be encoded in anatomically segregated, agent-specific neuronal populations. This facilitates Self-Other distinction and provides a neurocomputational architecture for mentalising. It has previously been shown that this agent-specificity in prediction error encoding is malleable, and sensitive to social context<sup>11</sup>.

Our findings add to this literature by showing that pharmacological manipulation with psychedelics also reduces the agent-specificity of prediction errors. This is consistent with previous work showing that ketamine and psilocybin impair one's ability to discriminate Self-generated from Other-generated touch<sup>50</sup> and modulate bodily Self-experience<sup>65</sup>, and that psilocybin attenuates the distinctiveness of neural responses to hearing the voice of oneself versus another person<sup>51</sup>. Where these previous studies speak to Self-Other distinction with respect to observable physical phenomena (e.g. *Whose voice is that?* or *Whose touch is that?*), the current work extends this research to unobservable mental phenomena (e.g. *Whose belief is that?*).

A recent meta-analysis found that classical psychedelics enhance emotional empathy, the propensity to experience other people's emotional states<sup>49</sup>. This finding has also been replicated in patients with major depression<sup>43</sup>. However, these previous studies found no effect of psychedelics on cognitive empathy, the ability to accurately infer and report the mental states of Others. Our findings show that the picture is more nuanced. Not only do psychedelics increase one's propensity to attribute another person's affective states to themselves, but also to attribute another person's non-affective belief updates to themselves (*allocentric* belief updating). This is consistent with previous findings that psychedelics acutely make people more suggestible<sup>66</sup> and more likely to adopt the opinions of Others<sup>67</sup>. In our data we found that behaviour was best explained by a model that included not only *allocentric* updating, but also *idiocentric* updating (propensity to attribute one's own mental states to Other), and that these two vectors of belief updating were modified in equal measure. This drug-induced Self-Other mergence does not manifest as enhanced cognitive empathy, which would result in improved accuracy at identifying the beliefs of the other person. However, our findings indicate that blurring of Self-Other boundaries with respect to affective mental states extends to

non-affective mental states. We therefore propose that the theory of psychedelics influencing affective empathy, and not cognitive empathy, is an oversimplification.

Our demonstration that 2C-B and psilocybin act on the processes underlying cognitive empathy has important implications for understanding the therapeutic effects of psychedelics in mental illness<sup>68,69</sup>. It is thought that the (sub)acute subjective effects of psychedelics, such as ego dissolution and perceived connectedness, are predictive of therapeutic efficacy<sup>39,40</sup>. Converging evidence from preclinical studies and human neuroimaging suggests that the enduring prosocial state during the subacute *afterglow* period may reflect a reopened critical period for social learning and cognitive restructuring<sup>70-74</sup>, necessary for the *integration* component of psychedelic-assisted psychotherapy<sup>75,76</sup>. This window of neuroplasticity appears to be induced by 5-HT<sub>2A</sub>-mediated increases in prefrontal excitatory signalling<sup>77</sup>. In our study we found that interindividual variability in psychedelic-induced Self-Other mergence, during this subacute period, was predictive of a canonical variate capturing both acute ego dissolution and subacute psychosocial wellbeing. This is consistent with previous work showing that people who exhibit less of a neural Self-Other distinction in reward processing, are more likely to engage in prosocial behaviour<sup>78</sup>.

We found that drug-induced changes to Self-Other mergence were positively associated with self-reported psychosocial wellbeing. The direction of this association is unlikely to be universal. It has previously been shown in healthy adults that a baseline propensity for Self-Other mergence is associated with transdiagnostic subclinical traits of mental illness<sup>10</sup>. Furthermore people with diagnosed borderline personality disorder (BPD) exhibit enhanced idiocentric updating, relative to the general population<sup>27</sup>. In healthy physiology, the brain is likely to balance Self-Other distinction with Self-Other mergence in a way that suits the current social circumstances<sup>11</sup>. We speculate that this ability to adaptively generalise between Self and Other, without over-relying on Self-Other mergence, is a marker of good mental health. The period of heightened plasticity following psychedelic administration may enable a re-learning of Self-Other boundaries, which could be facilitated by cognitive training or psychotherapy. Whilst, in our data, we observed heightened Self-Other mergence one day after drug administration, it will be useful for future work to establish trajectories of Self-Other boundary transformation over a longer period of time and to directly assess whether psychedelics are simply promoting Self-Other mergence, or promoting sensitivity to Self-Other boundary re-learning.

Psychedelics profoundly alter how people form beliefs about themselves and about the world, and this is not without risks<sup>79</sup>. In order to develop safe, individualised treatments in the emerging field of psychedelic psychiatry, it will be useful to have access to objective biomarkers that can predict and track treatment responses<sup>80</sup>. Quantifying Self-Other mergence in belief updating may provide a tool to predict how individuals will respond to psychedelic treatment. For instance, increased monitoring may be warranted during the subacute phase if giving

psychedelics to individuals with maladaptive cognitions and behaviours related to an over-reliance on Self-Other generalisation, such as those exhibited by people with personality disorder, or individuals showing high Self-Other mergence in pre-treatment behavioural screening. Individual variability in Self-Other processing may also explain why some of the general population experience adverse effects or particularly positive psychosocial effects when using psychedelics recreationally.

Our findings also speak to a possible computational mechanism behind acute psychedelic-induced ego dissolution. In previous work applying reinforcement learning models to investigate the impact of 5-HT<sub>2A</sub> agonists on learning and belief formation, it was shown that lysergic acid diethylamide (LSD) enhances learning rates in humans<sup>81</sup> whilst 5-HT<sub>2A</sub> antagonism decreases learning rates in rodents<sup>82</sup>. Recent work has also shown that LSD-induced learning rate enhancement is correlated with functional connectivity changes in the default-mode network (DMN) that are thought to reflect ego dissolution<sup>83</sup>, consistent with the REBUS (RELaxed Beliefs Under pSychedelics) model<sup>84</sup>. The REBUS model proposes that psychedelics reduce the precision of beliefs, leading to an upweighting of prediction errors and increased sensitivity to new information. In allowing the brain access to new hypotheses about reality, previously rigid beliefs underpinning one's notion of Self may also be loosened, leading to ego dissolution<sup>85</sup>. Interestingly, the effect observed in the current study was not driven by enhanced sensitivity to surprise, but specifically enhanced sensitivity to cross-agent leaked surprise. By explicitly operationalising selfhood within our data and model, we found that psychedelics directly relax the constraints defining Self-Other boundaries in the belief-updating process. It is possible to conceptualise this within the REBUS framework by considering the neural instantiation of prediction errors as itself susceptible to a meta-learning process wherein metarepresentational learning signals<sup>6</sup> are upweighted in the face of reduced precision of Self-Other boundaries. In any case, the dimension of altered belief-updating introduced in the current study may permit a direct and objective quantification of ego dissolution, beyond indirect physiological correlates.

Default-mode network (DMN) functional disintegration is a reported correlate of acute ego dissolution and enduring psychedelic effects<sup>71,73,86</sup>, and this is complemented by our neuroimaging findings. Using spectral dynamical causal modelling (DCM) to estimate effective connectivity between hubs of the so-called *social brain*, we found that 2C-B and psilocybin induced acute connectivity changes associated with subsequent Self-Other mergence. Most striking, was a drug-induced attenuation of an inhibitory connection from right temporoparietal junction (rTPJ) to dorsomedial prefrontal cortex (dmPFC), both components of the DMN. Both the rTPJ and dmPFC are involved in inferring the mental states of Others<sup>87,88</sup>. Previous work supports a role for the dmPFC in distinguishing Self from Other<sup>89,90</sup>; dmPFC neurons encode specific information about other agents' behaviours<sup>91</sup> and mental states<sup>12</sup>. At the same time, the rTPJ is thought to play a role in suppressing one's current perspective, to allow for alternative

perspectives to be considered<sup>92–94</sup>. It has recently been shown that psilocybin-induced disembodiment is associated with altered rTPJ effective connectivity<sup>95</sup>. Taken together, the existing literature suggests that these two brain regions act in concert to allow for perspective-taking with selective attribution of a mental state to a specific agent.

Our results show that when an inhibitory connection from rTPJ to dmPFC is weakened, mental states shift to become less agent-specific and more agent-independent. Moreover, individuals with weaker baseline rTPJ–dmPFC inhibition were more susceptible to psychedelic-induced Self-Other mergence, suggesting a potential biomarker for predicting treatment responsiveness. Together, these findings inform our understanding of how hubs of the *social brain* interact to modulate Self-Other distinction and how connectivity changes within this network relate to the subjective effects of psychedelics. These findings can be used to guide future work investigating adjunctive interventions (e.g., transcranial magnetic stimulation) to probe causal mechanisms or modulate maladaptive belief updating<sup>96,97</sup>.

There are several limitations to the current work. First, as the study was conducted in a small sample of healthy volunteers, we cannot determine whether the drug-induced changes we observed are clinically significant, nor can we exclude the possibility that effects would differ in patient or psychedelic-naive populations. While both 2C-B and psilocybin produced comparable effects on Self-Other processing and wellbeing at moderate doses, suggesting a generalisable mechanism across psychedelics, the absence of significant differences in this sample should not be taken as evidence that the two compounds act equivalently. Pharmacodynamic distinctions between psychedelics have been shown to scale up to whole-brain functional connectivity<sup>98,99</sup>, and may modulate the haemodynamic response function that underpins the core neurovascular coupling assumptions of biophysical models such as DCM<sup>100</sup>. Larger comparative studies and pharmacological challenge designs (e.g., selective 5-HT<sub>2A</sub> antagonists, serotonin reuptake inhibitors) will be needed to specify the neuropharmacological basis of the effects reported here.

One caveat of the probabilistic False Belief Task (pFBT) is that it is not strictly social per se. Whilst this task can be adapted to make use of real-life social partners<sup>10,11</sup>, this is not essential for the dynamics of the task since there is no social interaction involved. The task simply requires the subject to keep track of two fluctuating random variables, presented with a social framing. This social framing has been shown to be effective and induces changes in behaviour on the task and the neural representations of the computational variables underpinning the learning model<sup>10</sup>. Nevertheless, future work should aim to assess how the metrics elicited by the pFBT translate to real social interactions.

In summary, our findings show that classical and novel psychedelics induce Self-Other mergence with respect to prediction errors, the fundamental computational units behind learning and belief formation. By integrating principles of computational psychiatry with

psychedelic psychiatry, we have shown that there are objective, quantifiable brain-based and behaviour-based biomarkers of psychedelic-induced Self-Other mergence. These markers capture biological, psychological and social dimensions of an individual's state. We propose they may be useful for holistically tracking the effects of a person's psychedelic treatment. More broadly, our findings support the view that the neurocomputational basis of selfhood lies in flexible and manipulable Self-Other boundaries. We hope these insights encourage a shift away from treating selfhood as an ineffable quality and toward viewing it as a set of quantifiable, measurable processes.

## Methods

### Participants

The data were collected as part of a larger trial (trial register: NL8813, approved by the Academic Hospital and the University's Medical Ethics Committee of Maastricht University (NL73539.068.20)<sup>99,101</sup>. Twenty-two healthy participants (11 female) aged 19-35 years (mean  $\pm$  SD: 25  $\pm$  4 years) were recruited by word of mouth and advertisement shared via Maastricht University social media. Participants were required to have had at least one previous lifetime psychedelic exposure (e.g., psilocybin, mescaline, LSD), but no exposure within the past 3 months. Those with psychiatric, major medical, endocrine, or neurological conditions, pregnant or lactating women, those not using reliable contraception, concomitant drug use, with a history of adverse reactions to psychedelics, with uncorrected or abnormal vision, or MRI contraindications were excluded (see Supplementary Materials for more details).

### Study design

We employed a double-blind, placebo-controlled, randomised crossover design with three acute rsfMRI sessions and two follow-ups per cycle (in-person at +1 day [FU1] and remotely at +5 days [FU2]; see Fig. 1). During dosing sessions, participants received 15 mg psilocybin, 20 mg 2C-B, or placebo (bittering agent). Doses were selected based on prior evidence of psychotropic equivalence, as detailed in companion manuscripts<sup>99,101</sup>. Acute eyes-open rsfMRI was acquired approximately 110 min after administration. In the present sample, psychotropic equivalence was reconfirmed immediately before and after rsfMRI acquisition across multiple measures, with no significant differences observed between compounds (see Supplementary Table 1). Further details are provided in the Supplementary Materials.

### The probabilistic false beliefs task (p-FBT)

One day after each dosing session (FU1) and following negative drug and alcohol screenings, participants returned to the laboratory to complete a probabilistic false belief task (p-FBT)<sup>11,27</sup> in person (Fig. 1). In the p-FBT, participants observed a sequence of binary outcomes generated from a Bernoulli distribution with an underlying probability,  $P$ . The task was framed such that participants imagined themselves observing a sequence of sales in a tourist store on a tropical island, where each purchase was either a pink umbrella (with probability  $P$ ) or a yellow sun-shade (with probability  $1-P$ ).  $P$  could take any value between 0 and 1, and the task was designed such that  $P$  gradually drifted along a random walk. Through a series of successive *sampling trials*, participants were instructed to infer, and keep track of, the drifting hidden variable  $P$ .

In addition to this, participants were instructed to infer, and keep track of, another agent's false beliefs about the drifting hidden variable  $P$ . The task included a fictional store manager who sampled only limited information, and sometimes misleading information, about the observed outcomes. Each sampling trial was *shared*, *privileged*, or *decoy*. In *shared sampling trials*, an icon of the cartoon manager looking through an open door indicated that the manager was also observing the outcome. In *privileged sampling trials*, an icon of the cartoon manager behind a closed door indicated that the manager was in a different room, unable to observe the information from the *sampling trial*. Finally, in *decoy sampling trials*, participants were told that the manager was observing sales on a security camera; however, unbeknown to the manager, the footage was out of date, and thereby uninformative about what the customers were really buying. These *decoy sampling trials* were indicated by an image of the cartoon manager's face next to an icon of a computer disc.

Trial sequences were pre-generated such that the underlying probability trajectories for the participant ( $P_{Self}$ ) and the manager ( $P_{Other}$ ) were uncorrelated. The uncorrelated random walks were bound between 0 and 1 and adjusted in steps of 0.025 on each sampling trial. *Privileged sampling trials* drew outcomes from  $P_{Self}$ , *decoy sampling trials* drew outcomes from  $P_{Other}$ . In *shared sampling trials* the outcomes were randomly drawn from either probability distribution.

A total of 360 trial sequences were pre-generated, each containing 120 of each of the three *sampling trial* types in a randomised order. One of these trial sequences was randomly selected for each participant's testing session. Each *sampling trial* was presented for 1250 ms, followed by a variable intertrial interval with a fixation cross on screen for 400–600 ms. Every four to nine *sampling trials*, participants were given one of two possible *probe trials*. In a *Self-probe* trial, participants were prompted to report their own belief about the probability that the next trial will reveal a pink umbrella ( $P_{Self}$ ). Alternatively, in *Other-probe* trials, participants were prompted to report an estimate of the manager's false belief about this probability ( $P_{Other}$ ). These reports were made by positioning a cursor along a continuous probability scale. For *Other-probe* trials, participants were explicitly instructed to consider which sales the manager had seen or missed, including information viewed on the misleading security footage.

Prior to testing, participants received a detailed walkthrough of the task, including training runs and an introductory example. This process explained the different trial types and the use of the continuous rating scale. Participants who rated their understanding as *poor* or *very poor* were required to repeat the instructions and the practice phase. Participants were informed that the task was intentionally challenging and were encouraged to perform to the best of their ability. Each testing session began with a practice run of the task. Feedback was provided only during the introductory example; no feedback was given during the main task. To

minimise fatigue, a self-paced break was introduced halfway through the task. The task was programmed using Matlab (MathWorks, Provo), visualised with Cogent 2000 (v125) and Cogent Graphics (v1.29), and administered in a quiet testing room.

### Quantifying task performance

For each behavioural testing session, we evaluated task performance by computing the Spearman correlation coefficient between the participant's responses on *probe trials* and the true underlying Bernoulli parameters  $P_{Self}$  and  $P_{Other}$ . In addition, for each testing session we also computed a permutation-based null performance measurement, by randomly shuffling the subject's responses on *probe trials* 1000 times, and each time computing the Spearman coefficient between the shuffled responses and the underlying Bernoulli parameters. We then subtracted the average of these null correlation coefficients from the true correlation coefficient. The resulting value represents behavioural performance, relative to chance-level random responding.

It is important to note that participants only sampled limited information about the true underlying probabilities  $P_{Self}$  and  $P_{Other}$ . It was therefore not possible for participants to make perfect inferences about  $P$ . We therefore approximated an "optimal" performance for each testing session by simulating a subject who inferred  $P_{Self}$  and  $P_{Other}$  by using a recency-weighted average of recently observed outcomes. We used a simplified version of the model described below (see 'Behavioural model') where the simulated subject used parallel belief updates for Self and Other, with a perfect memory and no behavioural stochasticity. Learning rate  $\alpha$  was arbitrarily set to 0.15. Across all the testing sessions that were delivered to participants, these simulated "optimal" observers had a mean performance-relative-to-chance correlation of 0.57 (IQR: 0.51 - 0.63). We used the upper quartile of values to denote a window of "optimal" performance (0.63 - 0.75). Performance metrics should therefore be interpreted so that 0 represents completely random behaviour and scores in excess of 0.6 represent near-optimal performance.

Two participants were excluded from behavioural modelling as they consistently performed at or below chance-level in multiple sessions. This left 20 participants for inclusion in the modelling analyses. One participant only completed two sessions, with missing data for the placebo session.

### Behavioural model

We fitted each participant's responses on *probe trials* with a simple reinforcement-learning style model, wherein the belief,  $B \in [0, 1]$ , about  $P$  on any trial  $t$ , is updated proportionally to a *prediction error* (PE). On any given *sampling trial*, the PE is the difference between the observed

outcome (1 or 0) and the previous belief state. To model agent-specificity, we used parallel learning processes for Self and Other:

$$B_t^{Self} = B_{t-1}^{Self} + (\alpha^{Self} \cdot PE_t^{Self}) + (\Lambda^{Other} \cdot PE_t^{Other}) + \delta^{Self} (0.5 - B_{t-1}^{Self})$$

$$B_t^{Other} = B_{t-1}^{Other} + (\alpha^{Other} \cdot PE_t^{Other}) + (\Lambda^{Self} \cdot PE_t^{Self}) + \delta^{Other} (0.5 - B_{t-1}^{Other})$$

Here,  $B^{Self}$  denotes the participant's belief, while  $B^{Other}$  denotes the participant's simulation of the other agent's belief. Thus, two belief representations are separately updated, based on observations available to Self and Other.  $\alpha \in [0, 1]$  is a learning rate parameter, while  $\delta \in [0, 1]$  is a memory decay rate, allowing beliefs to drift towards chance level. Importantly, updates for Self and Other may not be fully segregated and there may be a degree of Self-Other mergence. This is captured using *leak* parameter,  $\Lambda \in [-1, 1]$ , which allows one agent's prediction error to erroneously inform the other agent's belief. Where  $\alpha$  describes sensitivity to surprise,  $\Lambda$  describes sensitivity to cross-agent leaked surprise. If  $\Lambda = 0$  then prediction errors are entirely agent-specific. As  $\Lambda$  deviates from 0, prediction errors become progressively more agent-independent, resulting in Self-Other mergence of belief-updating channels. This form of model has previously been shown to fit behaviour well on the p-FBT<sup>10,11,27</sup>.

In order to keep belief states naturally bound between 0 and 1, a factorised version of the model is fit to the data<sup>27</sup>. In this factorised form,  $\Lambda$  (sensitivity to cross-agent leaked surprise) is obtained by multiplying  $\alpha$  with  $\lambda \in [-1, 1]$ .

$$B_t^{Self} = B_{t-1}^{Self} + \alpha^{Self} (PE_t^{Self} + \lambda^{Other} \cdot PE_t^{Other}) + \delta^{Self} (0.5 - B_{t-1}^{Self})$$

$$B_t^{Other} = B_{t-1}^{Other} + \alpha^{Other} (PE_t^{Other} + \lambda^{Self} \cdot PE_t^{Self}) + \delta^{Other} (0.5 - B_{t-1}^{Other})$$

### Model fitting procedure

The model yields point estimates for  $B^{Self}$  and  $B^{Other}$  on each trial. These estimates were used to derive the likelihood of a participant's responses across all *probe trials*. The likelihood was approximated using two Beta distributions (bound between 0 and 1), with the distribution mode set to the model-derived point estimate ( $B^{Self}$  or  $B^{Other}$ ) and the Beta distribution

variance set to a participant-specific free parameter,  $\tau^{Self}$  or  $\tau^{Other}$ , which can be thought of as decision temperature parameters, governing response stochasticity. High values in these parameters result in random responding, whilst lower values result in precise responses that are tightly coupled to the underlying belief states  $B^{Self}$  and  $B^{Other}$ .

The model can be specified in a variety of ways, resulting in a family of nested models that can be fitted to the data. Learning rate  $\alpha$  can be constrained such that  $\alpha^{Self} = \alpha^{Other}$  or they can be modelled as two separate free parameters. Temperature  $\tau$  can be constrained such that  $\tau^{Self} = \tau^{Other}$  or modelled as two separate free parameters. Memory decay  $\delta$  can be constrained such that  $\delta^{Self} = \delta^{Other}$  or modelled as two separate free parameters, or  $\delta$  can be removed completely such that  $\delta^{Self} = \delta^{Other} = 0$ . Finally, leak  $\lambda$  can be constrained such that  $\lambda^{Self} = \lambda^{Other}$  or can be modelled as two separate free parameters, or removed completely such that  $\lambda^{Self} = \lambda^{Other} = 0$ . In summary,  $\alpha$  has one or two parameters;  $\tau$  has one or two parameters;  $\delta$  has zero, one, or two parameters;  $\lambda$  has zero, one, or two parameters. This results in 36 nested model variants.

We fitted all 36 model variants to each individual testing session (three testing sessions per participant). Parameters were bounded through sigmoid transformation and fitted in inverse sigmoid space. Both  $\alpha$  and  $\delta$  parameters were bounded between 0 and 1,  $\lambda$  between -1 and 1, and  $\tau$  between 0.001 and 0.08 (this upper bound on  $\tau$  produces a near-uniform Beta distribution, corresponding to random responding). For each participant we sought model parameters with maximum posterior probability, given the data. We used fixed Gaussian priors over parameters and assumed a uniform prior over  $B$  models. For each model variant, we calculated the Bayesian model evidence  $P(x|M)$  (probability of data  $x$ , given the model  $M$ ), which favours models that provide a closer fit to the data, whilst penalising model complexity. Bayesian model evidence was calculated by marginalising the joint probability  $P(x, \theta|M)$ , with respect to the parameters  $\theta$ . The exact calculation for  $P(x|M)$  is given by:

$$P(x|M) = \int P(\theta)P(x|\theta, M) d\theta$$

As per Story et al.<sup>27</sup>, we approximated this integral numerically, by randomly sampling parameters from the prior distribution, such that:

$$P(x|M) \approx \frac{1}{K} \sum_{k=1}^K P(x|\theta_k, M)$$

Where  $\theta_k$  is a randomly sampled set of parameter values for a given model, and K is the total number of samples, set to 2000. Raw fitted parameter values in unconstrained native space are shown in supplementary figure S1.

To measure parameter recovery for the best fitting model, we fitted the model to data generated for 1000 simulated participants. Simulated participants' parameters were randomly sampled from the observed maximum a posteriori parameter estimates. Each parameter was sampled independently, with replacement. We computed a Pearson correlation coefficient between the generative parameters and fitted parameter estimates.

### **Subjective effects**

On each dosing day, psychotropic equivalency was reassessed using 0-100 visual analogue scales (VAS) to measure subjective high and effect intensity. The phenomenological content of each dosing visit was retrospectively evaluated using the 5 Dimensions of Altered States Questionnaire (5D-ASC) and Ego Dissolution Inventory (EDI) around 360 minutes post-drug intake<sup>102,103</sup>. Participants were contacted to assess “afterglow” effects at follow-up 1 (FU1, in person) and follow-up 2 (FU2, remotely via Qualtrics XM12). At both FU1 and FU2, participants completed the Social Connectedness Scale–Revised (SCS-R) and the Inclusion of Other in the Self (IOS) to capture changes in social functioning<sup>104,105</sup>, as well as the Positive and Negative Affect Schedule (PANAS) to index ongoing mood states<sup>106</sup>. At FU2, they also completed the Persisting Effects Questionnaire (PEQ) to evaluate longer-term psychosocial changes attributable to treatment<sup>107</sup>. All outcomes were assessed by linear mixed effect models (LMEMs) using drug condition as a main fixed effect (and where applicable an interaction of time) and participant as a random intercept. Follow-up assessments were performed using Tukey's method.

To test whether drug-induced changes in Self-Other mergence covaried with (sub)acute psychosocial outcomes, we conducted a subject-level permutation MANOVA. Outcome measures were residualised for mean drug effects (active vs. placebo) and z-scored, ensuring that associations reflected within-person, drug-evoked covariation. Multivariate association was evaluated with Pillai's trace, with significance determined from 10 000 within-subject permutations. Canonical correlation analysis then identified the linear outcome combination most strongly associated with Self-Other mergence, with significance tested by a 10 000-permutation Pearson correlation. A repeated-measures correlation further quantified the within-subject association between Self-Other mergence and canonical scores across conditions<sup>108</sup>. Robustness was evaluated using five-fold cross-validation, testing whether canonical scores could be predicted in held-out folds, with significance assessed against a 10 000-permutation null. This within-subject framework ensures that associations reflect how individual differences

in Self-Other mergence track psychosocial responses under psychedelics, rather than being driven by overall drug–placebo differences.

### **Administration order effects**

To examine potential order effects, we estimated subject fixed-effects regressions for all significant task and questionnaire outcomes, including an interaction term between drug condition and session day (1–3). Cluster-robust standard errors were applied. No significant drug × day interactions were detected (all  $F_{2,19} < 1.7$ ,  $p > 0.21$ ).

### **Neuroimaging acquisition**

Eyes-open rsfMRI data were acquired using a 7T Siemens Magnetom scanner (Siemens Medical, Erlangen, Germany) equipped with a 32-channel Nova Medical head coil (Nova Medical Inc., Wilmington, MA). Functional images (516 volumes, approximately 12 minutes) were obtained using a gradient echo-planar imaging (EPI) sequence with 1.5 mm isotropic voxels (TR = 1400 ms). Preprocessing was performed using SPM12 employing a previously described effective connectivity pipeline<sup>42,109</sup> adapted for ultra-high-field imaging. Steps included non-steady-state volume discarding, susceptibility distortion correction, slice-timing correction, realignment, spatial normalisation to the Montreal Neurological Institute (MNI) EPI template, and spatial smoothing with a 4-mm full width at half maximum Gaussian kernel. No significant differences in mean framewise displacement (FD, < 0.5 mm) were identified between conditions. Additional details are provided in the Supplementary Materials.

### **Extraction of region coordinates**

To evaluate whether dynamical changes in Theory of Mind (ToM) brain systems reflected acute psychedelic effects and subsequent alterations in mentalising abilities, we defined five regions reported as central to self-other processing such as belief attribution, perspective-taking, and understanding social intentions<sup>110–112</sup>. For this purpose, we isolated the dorsomedial prefrontal cortex, ventromedial prefrontal cortex, precuneus and the bilateral temporal parietal junction as targets (dmPFC, vmPFC, Pcu, ITPJ, rTPJ respectively). Neurosynth, a large-scale automated synthesis of functional neuroimaging data, was used to identify corresponding peak activation coordinates<sup>113</sup>, from which ROIs were centred and masked by an 8-mm radius sphere (see Table 1). More details can be found in the Supplementary Materials.

MNI coordinates			
	X	Y	Z
<b>Region</b>			
Dorsomedial Prefrontal Cortex (dmPFC)	-2	42	-8
Ventromedial Prefrontal Cortex (vmPFC)	2	46	38
Precuneus (PCu)	2	-52	36
Left Temporoparietal Junction (ITPJ)	-50	-56	22
Right Temporoparietal Junction (rTPJ)	58	-56	18

**Table 1. Coordinates of regions of interest.** Abbreviations: dorsomedial prefrontal cortex (dmPFC), ventromedial prefrontal cortex (vmPFC), precuneus (Prcun), left temporoparietal junction (ITPJ), right temporoparietal junction (rTPJ).

A general linear model was used to regress head motion and physiological noise from each regional time series using I) 6 head motion parameters (3 translation and 3 rotational), II) cerebrospinal fluid (CSF) signals (extracted from left ventricle using a 4-mm sphere [0 -40 -5] ), and III) white matter signals (extracted from pons using a 4-mm sphere [0 -30 -25]). Low-frequency signal drifts were filtered using a 128-s high-pass filter. Global signal regression was not performed as there is evidence that it does not substantially impact results in small network analyses<sup>114</sup>.

### **Effective connectivity estimation using spectral DCM**

We fitted a spectral dynamic causal model (DCM) to each rsfMRI timeseries, using the DCM toolbox in SPM12. Spectral DCM fits a biophysical state-space model to the observed cross-spectra of BOLD signals, to estimate underlying neuronal states and the rate of change in neural activity in each region (in hertz) as a function of activity in other regions<sup>115,116</sup>. For each of the 20 participants with p-FBT data, and each experimental session, a fully connected DCM was constructed using the five ROIs outlined in Table 1, with no inclusion of exogenous inputs. The DCM fitted the data well, and the amount of explained variance, averaged 88.13% (min = 79.08%, max = 96.26%) nor significantly differed across conditions (See Supplementary Materials). In DCM, self-connections are log-scaled and inherently inhibitory, reflecting a region's sensitivity to incoming inputs: positive values indicate increased self-inhibition, while negative values reflect disinhibition and enhanced synaptic gain. Note that only self-connections are log-transformed (see Supplementary for details).

### **Group inference using Parametric Empirical Bayes (PEB)**

For each participant, three first-level DCMs (one per session) were summarised with a subject-specific second-level PEB model<sup>117</sup>. This 2nd-level model included an intercept and a drug regressor (2C-B and psilocybin vs. placebo). Group-level inference was then performed using a PEB-of-PEBs approach<sup>118</sup>, where subject-specific PEBs served as inputs to a third-level model. The third-level design matrix contained two regressors: (i) a constant modelling group-average connectivity, and (ii) a behavioural covariate reflecting drug-induced changes in Self-Other merge (mean psychedelic minus placebo). This framework therefore allowed us to estimate group-average placebo connectivity, group-average drug effects, and associations between connectivity (both baseline and drug-induced changes) and behaviour. Models were pruned using exploratory Bayesian model reduction and Bayesian model comparison, employing an automatic greedy search over reduced models that iteratively removed parameters not contributing to model evidence. A Bayesian model average was then computed across the 256 best-fitting reduced models (default) from the final iteration. We applied a conservative statistical threshold of a posterior probability >0.99 equivalent to a Bayes Factor of >150 and considered a “very strong” evidence<sup>119</sup>. In the case of behaviour, the posterior probabilities express the likeliness of association between effective connectivity estimates and behavioural scores. The details of the biophysical model in DCM, model inversion at first-level and higher levels, and Bayesian model reduction, have already been extensively documented<sup>118,120,121</sup>. Further specifications are provided in the Supplementary Materials.

## **Funding**

J.G.R acknowledges financial support from Dutch Research Council (NWO, grant number 406.18.GO.019). NLM is also financially supported by the Dutch Research Council (NWO, grant number VI.Veni.231G.011).

## **Conflict of Interest**

J.G.R is a scientific consultant to GH Research who have no involvement in the preparation or conception of this manuscript or related data. K.P. is currently an employee of Boehringer Ingelheim GmbH & Co KG, and is currently Chief Scientist and on the Board of Directors for the Heffter Research Institute, and scientific advisor for the MIND foundation. All other authors have no relevant conflicts of interest to declare.

## **Data availability**

Code generated for analysis is to be made available at the following link: ([https://github.com/PabloMallaroni/project\\_etc\\_afterglow](https://github.com/PabloMallaroni/project_etc_afterglow)). Data and accompanying covariates can be made available to qualified research institutions upon request to J.G.R and a data use agreement executed with Maastricht University.

## **Author contributions**

P.M. contributed to conceptualisation, methodology, formal analysis, visualisation, investigation, data curation, project administration, and writing (original draft, review and editing); N.L.M. contributed to investigation, project administration, writing (review and editing), and supervision; K.H.P. contributed resources and writing (original draft, review and editing); A.R. contributed resources, methodology, and writing (original draft, review and editing); S.E. contributed to conceptualisation, methodology, formal analysis, visualisation, writing (original draft, review and editing), and supervision; J.G.R. contributed to project administration, data curation, writing (review and editing), funding acquisition, and supervision.

## Acknowledgements

We thank Emma de Brabander, Vitea Incerti-Medici, Ajna Jansson and Chantal Delaquis for their assistance with data collection. We are grateful to Dr. Cees van Leeuwen and Dr. Riccardo Paci for conducting the medical screenings.

## References

1. Fonagy, P. & Target, M. Playing with reality: I. Theory of mind and the normal development of psychic reality. *Int J Psychoanal* **77 ( Pt 2)**, 217–233 (1996).
2. Rakoczy, H. Foundations of theory of mind and its development in early childhood. *Nat. Rev. Psychol.* **1**, 223–235 (2022).
3. Fonagy, P. & Target, M. Mentalization and the changing aims of child psychoanalysis. *Psychoanal. Dialogues* **8**, 87–114 (1998).
4. Tamir, D. I. & Thornton, M. A. Modeling the predictive social mind. *Trends Cogn. Sci.* **22**, 201–212 (2018).
5. Koster-Hale, J. & Saxe, R. Theory of mind: a neural prediction problem. *Neuron* **79**, 836–848 (2013).
6. Ereira, S. A neurocomputational account of self-other distinction: from cell to society. (University College London, 2019).
7. Burkart, J. M. & Southgate, V. An evolutionary perspective on altercentrism. *Neurosci. Biobehav. Rev.* **176**, 106280 (2025).
8. Suzuki, S. *et al.* Learning to simulate others' decisions. *Neuron* **74**, 1125–1137 (2012).
9. Hill, M. R., Boorman, E. D. & Fried, I. Observational learning computations in neurons of

- the human anterior cingulate cortex. *Nat. Commun.* **7**, 12722 (2016).
10. Ereira, S., Dolan, R. J. & Kurth-Nelson, Z. Agent-specific learning signals for self-other distinction during mentalising. *PLoS Biol* **16**, e2004752 (2018).
  11. Ereira, S. *et al.* Social training reconfigures prediction errors to shape Self-Other boundaries. *Nat. Commun.* **11**, 3030 (2020).
  12. Jamali, M. *et al.* Single-neuronal predictions of others' beliefs in humans. *Nature* **591**, 610–614 (2021).
  13. Baron-Cohen, S., Leslie, A. M. & Frith, U. Does the autistic child have a 'theory of mind'? *Cognition* **21**, 37–46 (1985).
  14. Singer, T. *et al.* Empathy for pain involves the affective but not sensory components of pain. *Science* **303**, 1157–1162 (2004).
  15. Lamm, C., Decety, J. & Singer, T. Meta-analytic evidence for common and distinct neural networks associated with directly experienced pain and empathy for pain. *Neuroimage* **54**, 2492–2502 (2011).
  16. Mischkowski, D., Crocker, J. & Way, B. M. From painkiller to empathy killer: acetaminophen (paracetamol) reduces empathy for pain. *Soc. Cogn. Affect. Neurosci.* **11**, 1345–1353 (2016).
  17. Decety, J., Bartal, I. B.-A., Uzefovsky, F. & Knafo-Noam, A. Empathy as a driver of prosocial behaviour: highly conserved neurobehavioural mechanisms across species. *Philos. Trans. R. Soc. Lond. B Biol. Sci.* **371**, 20150077 (2016).
  18. Lamm, C., Bukowski, H. & Silani, G. From shared to distinct self-other representations in

- empathy: evidence from neurotypical function and socio-cognitive disorders. *Philos. Trans. R. Soc. Lond. B Biol. Sci.* **371**, 20150083 (2016).
19. Decety, J. & Fotopoulou, A. Why empathy has a beneficial impact on others in medicine: unifying theories. *Front. Behav. Neurosci.* **8**, 457 (2014).
  20. Zapparoli, L. *et al.* Self-other distinction modulates the sense of self-agency during joint actions. *Sci. Rep.* **14**, 30055 (2024).
  21. Neustadter, E. S., Fotopoulou, A., Steinfeld, M. & Fineberg, S. K. Mentalization and embodied selfhood in Borderline Personality Disorder. *J. Conscious. Stud.* **28**, 126–157 (2021).
  22. Zavlis, O., Story, G., Friedrich, C., Fonagy, P. & Moutoussis, M. A systematic review of computational modeling of interpersonal dynamics in psychopathology. *Nat. Ment. Health* (2025) doi:10.1038/s44220-025-00465-9.
  23. Andreou, M. & Skrimpa, V. Theory of Mind deficits and neurophysiological operations in autism spectrum disorders: A review. *Brain Sci.* **10**, 393 (2020).
  24. Yoshida, W. *et al.* Cooperation and heterogeneity of the autistic mind. *J. Neurosci.* **30**, 8815–8818 (2010).
  25. Thibaudeau, É., Cellard, C., Turcotte, M. & Achim, A. M. Functional impairments and theory of mind deficits in schizophrenia: A meta-analysis of the associations. *Schizophr. Bull.* **47**, 695–711 (2021).
  26. Alon, N. *et al.* (Mal)adaptive mentalizing in the cognitive hierarchy, and its link to paranoia. *Comput. Psychiatr.* **8**, 159–177 (2024).

27. Story, G. W. *et al.* A computational signature of self-other mergence in Borderline Personality Disorder. *Transl. Psychiatry* **14**, 473 (2024).
28. De Meulemeester, C., Lowyck, B. & Luyten, P. The role of impairments in self-other distinction in borderline personality disorder: A narrative review of recent evidence. *Neurosci. Biobehav. Rev.* **127**, 242–254 (2021).
29. Rothschild-Yakar, L. *et al.* Mentalizing self and other and affect regulation patterns in anorexia and depression. *Front. Psychol.* **10**, 2223 (2019).
30. Le Bouc, R. *et al.* My belief or yours? Differential theory of mind deficits in frontotemporal dementia and Alzheimer’s disease. *Brain* **135**, 3026–3038 (2012).
31. Dittrich, A. The standardized psychometric assessment of altered states of consciousness (ASCs) in humans. *Pharmacopsychiatry* **31 Suppl 2**, 80–84 (1998).
32. Letheby, C. & Gerrans, P. Self unbound: ego dissolution in psychedelic experience. *Neurosci. Conscious.* **2017**, nix016 (2017).
33. Preller, K. H. & Vollenweider, F. X. Phenomenology, structure, and dynamic of psychedelic states. *Curr. Top. Behav. Neurosci.* **36**, 221–256 (2018).
34. Preller, K. H. *et al.* The effect of 5-HT<sub>2A/1a</sub> agonist treatment on social cognition, empathy, and social decision-making. *Eur. Psychiatry* **30**, 22 (2015).
35. Preller, K. H. *et al.* Role of the 5-HT<sub>2A</sub> receptor in self- and other-initiated social interaction in lysergic acid diethylamide-induced states: A pharmacological fMRI study. *J. Neurosci.* **38**, 3603–3611 (2018).
36. Millière, R., Carhart-Harris, R. L., Roseman, L., Trautwein, F.-M. & Berkovich-Ohana, A.

- Psychedelics, meditation, and self-consciousness. *Front. Psychol.* **9**, 1475 (2018).
37. Klee, G. D. Lysergic acid diethylamide (LSD-25) and ego functions. *Arch. Gen. Psychiatry* **8**, 461–474 (1963).
  38. Fischman, L. G. Dreams, hallucinogenic drug states, and schizophrenia: a psychological and biological comparison. *Schizophr. Bull.* **9**, 73–94 (1983).
  39. Kałużna, A., Schlosser, M., Gulliksen Craste, E., Stroud, J. & Cooke, J. Being no one, being One: The role of ego-dissolution and connectedness in the therapeutic effects of psychedelic experience. *J. Psychedelic Stud.* **6**, 111–136 (2022).
  40. Yaden, D. B. & Griffiths, R. R. The subjective effects of psychedelics are necessary for their enduring therapeutic effects. *ACS Pharmacol. Transl. Sci.* **4**, 568–572 (2021).
  41. Nutt, D. & Carhart-Harris, R. The current status of psychedelics in psychiatry. *JAMA Psychiatry* **78**, 121–122 (2021).
  42. Stoliker, D. *et al.* Effective connectivity of functionally anticorrelated networks under lysergic acid diethylamide. *Biol. Psychiatry* **93**, 224–232 (2023).
  43. Jungwirth, J., von Rotz, R., Dziobek, I., Vollenweider, F. X. & Preller, K. H. Psilocybin increases emotional empathy in patients with major depression. *Mol. Psychiatry* **30**, 2665–2672 (2025).
  44. Majić, T., Schmidt, T. T. & Gallinat, J. Peak experiences and the afterglow phenomenon: When and how do therapeutic effects of hallucinogens depend on psychedelic experiences? *J. Psychopharmacol.* **29**, 241–253 (2015).
  45. Majić, T. *et al.* The Afterglow Inventory (AGI): Validation of a new instrument for measuring

- subacute effects of classic serotonergic psychedelics. *J. Psychopharmacol.* **39**, 474–488 (2025).
46. Evens, R., Schmidt, M. E., Majić, T. & Schmidt, T. T. The psychedelic afterglow phenomenon: a systematic review of subacute effects of classic serotonergic psychedelics. *Ther. Adv. Psychopharmacol.* **13**, 20451253231172254 (2023).
47. Forstmann, M., Yudkin, D. A., Prosser, A. M. B., Heller, S. M. & Crockett, M. J. Transformative experience and social connectedness mediate the mood-enhancing effects of psychedelic use in naturalistic settings. *Proc. Natl. Acad. Sci. U. S. A.* **117**, 2338–2346 (2020).
48. Kettner, H. *et al.* Psychedelic communitas: Intersubjective experience during psychedelic group sessions predicts enduring changes in psychological wellbeing and social connectedness. *Front. Pharmacol.* **12**, 623985 (2021).
49. Olami, A. & Peled-Avron, L. Effects of classical psychedelics on implicit and explicit emotional empathy and cognitive empathy: a meta-analysis of MET task. *Sci. Rep.* **14**, 24480 (2024).
50. Kaldewaij, R. *et al.* Ketamine reduces the neural distinction between self- and other-produced affective touch: a randomized double-blind placebo-controlled study. *Neuropsychopharmacology* **49**, 1767–1774 (2024).
51. Smigielski, L. *et al.* P300-mediated modulations in self-other processing under psychedelic psilocybin are related to connectedness and changed meaning: A window into the self-other overlap. *Hum. Brain Mapp.* **41**, 4982–4996 (2020).

52. den Ouden, H. E. M., Kok, P. & de Lange, F. P. How prediction errors shape perception, attention, and motivation. *Front. Psychol.* **3**, 548 (2012).
53. Schultz, W., Dayan, P. & Montague, P. R. A neural substrate of prediction and reward. in *Foundations in Social Neuroscience* 541–554 (The MIT Press, 2002).
54. Adams, R. A., Huys, Q. J. M. & Roiser, J. P. Computational Psychiatry: towards a mathematically informed understanding of mental illness. *J. Neurol. Neurosurg. Psychiatry* **87**, 53–63 (2016).
55. Huys, Q. J. M. Advancing clinical improvements for patients using the theory-driven and data-driven branches of computational psychiatry. *JAMA Psychiatry* **75**, 225 (2018).
56. Rouault, M., Seow, T., Gillan, C. M. & Fleming, S. M. Psychiatric symptom dimensions are associated with dissociable shifts in metacognition but not task performance. *Biol. Psychiatry* **84**, 443–451 (2018).
57. Stoliker, D., Egan, G. F., Friston, K. J. & Razi, A. Neural mechanisms and psychology of psychedelic ego dissolution. *Pharmacol. Rev.* **74**, 876–917 (2022).
58. Burke, C. J., Tobler, P. N., Baddeley, M. & Schultz, W. Neural mechanisms of observational learning. *Proc. Natl. Acad. Sci. U. S. A.* **107**, 14431–14436 (2010).
59. Apps, M. A. J., Lesage, E. & Ramnani, N. Vicarious reinforcement learning signals when instructing others. *J. Neurosci.* **35**, 2904–2913 (2015).
60. Chang, S. W. C., Gariépy, J.-F. & Platt, M. L. Neuronal reference frames for social decisions in primate frontal cortex. *Nat. Neurosci.* **16**, 243–250 (2013).
61. Noritake, A., Ninomiya, T. & Isoda, M. Social reward monitoring and valuation in the

- macaque brain. *Nat. Neurosci.* **21**, 1452–1462 (2018).
62. Falcone, R., Cirillo, R., Ceccarelli, F. & Genovesio, A. Neural representation of others during action observation in posterior medial prefrontal cortex. *Cereb. Cortex* **32**, 4512–4523 (2022).
63. Noritake, A. & Isoda, M. The macaque medial prefrontal cortex simultaneously represents self and others' reward prediction error. *Cell Rep.* **44**, 115368 (2025).
64. Noritake, A., Ninomiya, T. & Isoda, M. Representation of distinct reward variables for self and other in primate lateral hypothalamus. *Proc. Natl. Acad. Sci. U. S. A.* **117**, 5516–5524 (2020).
65. Duerler, P. *et al.* Psilocybin induces aberrant prediction error processing of tactile mismatch responses-A simultaneous EEG-fMRI study. *Cereb. Cortex* **32**, 186–196 (2021).
66. Carhart-Harris, R. L. *et al.* LSD enhances suggestibility in healthy volunteers. *Psychopharmacology (Berl.)* **232**, 785–794 (2015).
67. Duerler, P., Schilbach, L., Stämpfli, P., Vollenweider, F. X. & Preller, K. H. LSD-induced increases in social adaptation to opinions similar to one's own are associated with stimulation of serotonin receptors. *Sci. Rep.* **10**, 12181 (2020).
68. Yao, Y. *et al.* Efficacy and safety of psychedelics for the treatment of mental disorders: A systematic review and meta-analysis. *Psychiatry Res.* **335**, 115886 (2024).
69. Luoma, J. B., Chwyl, C., Bathje, G. J., Davis, A. K. & Lancelotta, R. A meta-analysis of placebo-controlled trials of psychedelic-assisted therapy. *J. Psychoactive Drugs* **52**, 289–299 (2020).

70. Nardou, R. *et al.* Psychedelics reopen the social reward learning critical period. *Nature* **618**, 790–798 (2023).
71. Siegel, J. S. *et al.* Psilocybin desynchronizes the human brain. *Nature* **632**, 131–138 (2024).
72. Barrett, F. S., Doss, M. K., Sepeda, N. D., Pekar, J. J. & Griffiths, R. R. Emotions and brain function are altered up to one month after a single high dose of psilocybin. *Sci. Rep.* **10**, 2214 (2020).
73. Daws, R. E. *et al.* Increased global integration in the brain after psilocybin therapy for depression. *Nat. Med.* **28**, 844–851 (2022).
74. Doss, M. K. *et al.* Psilocybin therapy increases cognitive and neural flexibility in patients with major depressive disorder. *Transl. Psychiatry* **11**, 574 (2021).
75. Bathje, G. J., Majeski, E. & Kudowor, M. Psychedelic integration: An analysis of the concept and its practice. *Front. Psychol.* **13**, 824077 (2022).
76. Carhart-Harris, R. L. *et al.* Canalization and plasticity in psychopathology. *Neuropharmacology* **226**, 109398 (2023).
77. De Gregorio, D. *et al.* Lysergic acid diethylamide (LSD) promotes social behavior through mTORC1 in the excitatory neurotransmission. *Proc. Natl. Acad. Sci. U. S. A.* **118**, e2020705118 (2021).
78. Contreras-Huerta, L. S. *et al.* Neural representations of vicarious rewards are linked to interoception and prosocial behaviour. *Neuroimage* **269**, 119881 (2023).
79. McGovern, H. T. *et al.* An Integrated theory of false insights and beliefs under psychedelics. *Commun Psychol* **2**, 69 (2024).

80. Moujaes, F. *et al.* The emotional architecture of the psychedelic brain. *Trends Cogn. Sci.* (2025) doi:10.1016/j.tics.2025.07.006.
81. Kanen, J. *et al.* Effect of lysergic acid diethylamide (LSD) on reinforcement learning in humans. *Psychol. Med.* **53**, 6434–6445 (2023).
82. Hervig, M. E.-S. *et al.* 5-HT 2A and 5-HT 2C receptor antagonism differentially modulate reinforcement learning and cognitive flexibility: behavioural and computational evidence. *Psychopharmacology (Berl.)* **241**, 1631–1644 (2024).
83. Kanen, J. *et al.* 97. Neural and behavioral evidence from reinforcement learning converge in support of relaxed beliefs under LSD. *Int. J. Neuropsychopharmacol.* **28**, ii82–ii82 (2025).
84. Carhart-Harris, R. L. & Friston, K. J. REBUS and the anarchic brain: Toward a unified model of the brain action of psychedelics. *Pharmacol. Rev.* **71**, 316–344 (2019).
85. Stoliker, D., Egan, G. F. & Razi, A. Reduced precision underwrites ego dissolution and therapeutic outcomes under psychedelics. *Front. Neurosci.* **16**, 827400 (2022).
86. Soares, C., Gonzalo, G., Castelhana, J. & Castelo-Branco, M. The relationship between the default mode network and the theory of mind network as revealed by psychedelics - A meta-analysis. *Neurosci. Biobehav. Rev.* **152**, 105325 (2023).
87. Amodio, D. M. & Frith, C. D. Meeting of minds: the medial frontal cortex and social cognition. *Nat. Rev. Neurosci.* **7**, 268–277 (2006).
88. Saxe, R. & Kanwisher, N. People thinking about thinking people. The role of the temporo-parietal junction in ‘theory of mind’. *Neuroimage* **19**, 1835–1842 (2003).
89. Piva, M. *et al.* The dorsomedial prefrontal cortex computes task-invariant relative

- subjective value for self and other. *Elife* **8**, (2019).
90. Wittmann, M. K. *et al.* Causal manipulation of self-other mergence in the dorsomedial prefrontal cortex. *Neuron* **109**, 2353–2361.e11 (2021).
  91. Báez-Mendoza, R., Mastrobattista, E. P., Wang, A. J. & Williams, Z. M. Social agent identity cells in the prefrontal cortex of interacting groups of primates. *Science* **374**, eabb4149 (2021).
  92. Soutschek, A., Moisa, M., Ruff, C. C. & Tobler, P. N. The right temporoparietal junction enables delay of gratification by allowing decision makers to focus on future events. *PLoS Biol.* **18**, e3000800 (2020).
  93. Soutschek, A., Ruff, C. C., Strombach, T., Kalenscher, T. & Tobler, P. N. Brain stimulation reveals crucial role of overcoming self-centeredness in self-control. *Sci. Adv.* **2**, e1600992 (2016).
  94. Obeso, I., Moisa, M., Ruff, C. C. & Dreher, J.-C. A causal role for right temporo-parietal junction in signaling moral conflict. *Elife* **7**, (2018).
  95. Stoliker, D., Bernasconi, F., Blanke, O. & Razi, A. Psilocybin Modulates TPJ Effective Connectivity during Out-of-Body Experiences. *medRxiv* (2025).
  96. Young, L., Camprodon, J. A., Hauser, M., Pascual-Leone, A. & Saxe, R. Disruption of the right temporoparietal junction with transcranial magnetic stimulation reduces the role of beliefs in moral judgments. *Proc. Natl. Acad. Sci. U. S. A.* **107**, 6753–6758 (2010).
  97. Christian, P. *et al.* Belief updating during social interactions: Neural dynamics and causal role of dorsomedial prefrontal cortex. *J. Neurosci.* **44**, e1669232024 (2024).

98. Luppi, A. I. *et al.* In vivo mapping of pharmacologically induced functional reorganization onto the human brain's neurotransmitter landscape. *Sci. Adv.* **9**, eadf8332 (2023).
99. Mallaroni, P., Singleton, P., Mason, N. L., Satterthwaite, T. D. & Ramaekers, J. G. The forgotten psychedelic: Spatiotemporal mapping of brain organisation following the administration of 2C-B and psilocybin. *bioRxiv* (2024) doi:10.1101/2024.10.22.619393.
100. Zirkel, R. T. *et al.* Psilocybin prolongs the neurovascular coupling response in mouse visual cortex. *bioRxivorg* (2025) doi:10.1101/2025.07.25.666803.
101. Mallaroni, P. *et al.* Assessment of the Acute Effects of 2C-B vs. Psilocybin on Subjective Experience, Mood, and Cognition. *Clin Pharmacol Ther* **114**, 423–433 (2023).
102. Studerus, E., Gamma, A. & Vollenweider, F. X. Psychometric evaluation of the altered states of consciousness rating scale (OAV). *PLoS One* **5**, e12412 (2010).
103. Nour, M. M., Evans, L., Nutt, D. & Carhart-Harris, R. L. Ego-dissolution and psychedelics: Validation of the ego-Dissolution Inventory (EDI). *Front. Hum. Neurosci.* **10**, 269 (2016).
104. Aron, A., Aron, E. N. & Smollan, D. Inclusion of Other in the Self Scale and the structure of interpersonal closeness. *J. Pers. Soc. Psychol.* **63**, 596–612 (1992).
105. Lee, R. M., Draper, M. & Lee, S. Social Connectedness Scale--Revised. *PsycTESTS Dataset* American Psychological Association (APA) <https://doi.org/10.1037/t16389-000> (2013).
106. Watson, D., Clark, L. A. & Tellegen, A. Development and validation of brief measures of positive and negative affect: The PANAS scales. *J. Pers. Soc. Psychol.* **54**, 1063–1070 (1988).
107. Griffiths, R. R. *et al.* Psilocybin-occasioned mystical-type experience in combination with meditation and other spiritual practices produces enduring positive changes in

- psychological functioning and in trait measures of prosocial attitudes and behaviors. *J. Psychopharmacol.* **32**, 49–69 (2018).
108. Bakdash, J. Z. & Marusich, L. R. Repeated measures correlation. *Front. Psychol.* **8**, 456 (2017).
109. Preller, K. H. *et al.* Effective connectivity changes in LSD-induced altered states of consciousness in humans. *Proc. Natl. Acad. Sci. U. S. A.* **116**, 2743–2748 (2019).
110. Fehlbaum, L. V., Borbás, R., Paul, K., Eickhoff, S. B. & Raschle, N. M. Early and late neural correlates of mentalizing: ALE meta-analyses in adults, children and adolescents. *Soc. Cogn. Affect. Neurosci.* **17**, 351–366 (2022).
111. Saxe, R., Moran, J. M., Scholz, J. & Gabrieli, J. Overlapping and non-overlapping brain regions for theory of mind and self reflection in individual subjects. *Soc. Cogn. Affect. Neurosci.* **1**, 229–234 (2006).
112. Arioli, M., Cattaneo, Z., Ricciardi, E. & Canessa, N. Overlapping and specific neural correlates for empathizing, affective mentalizing, and cognitive mentalizing: A coordinate-based meta-analytic study. *Hum. Brain Mapp.* **42**, 4777–4804 (2021).
113. Tor D., W. NeuroSynth: a new platform for large-scale automated synthesis of human functional neuroimaging data. *Front. Neuroinform.* **5**, (2011).
114. Almgren, H. *et al.* Variability and reliability of effective connectivity within the core default mode network: A multi-site longitudinal spectral DCM study. *Neuroimage* **183**, 757–768 (2018).
115. Friston, K. J., Kahan, J., Biswal, B. & Razi, A. A DCM for resting state fMRI. *Neuroimage* **94**,

- 396–407 (2014).
116. Razi, A., Kahan, J., Rees, G. & Friston, K. J. Construct validation of a DCM for resting state fMRI. *Neuroimage* **106**, 1–14 (2015).
117. Friston, K. J. *et al.* Bayesian model reduction and empirical Bayes for group (DCM) studies. *Neuroimage* **128**, 413–431 (2016).
118. Zeidman, P. *et al.* A guide to group effective connectivity analysis, part 2: Second level analysis with PEB. *Neuroimage* **200**, 12–25 (2019).
119. Kass, R. E. & Raftery, A. E. Bayes factors. *J. Am. Stat. Assoc.* **90**, 773–795 (1995).
120. Zeidman, P. *et al.* A guide to group effective connectivity analysis, part 1: First level analysis with DCM for fMRI. *Neuroimage* **200**, 174–190 (2019).
121. Novelli, L., Friston, K. & Razi, A. Spectral dynamic causal modeling: A didactic introduction and its relationship with functional connectivity. *Netw Neurosci* **8**, 178–202 (2024).