

Attractor-like Dynamics in Belief Updating in Schizophrenia

Rick A. Adams,^{1,2*} Gary Napier,^{1*} Jonathan P. Roiser,¹ Christoph Mathys,^{3,4,5†} and James Gilleen^{6,7†}

¹Institute of Cognitive Neuroscience, University College London, London WC1N 3AZ, United Kingdom, ²Division of Psychiatry, University College London, London W1T 7NF, United Kingdom, ³Scuola Internazionale Superiore di Studi Avanzati, 34136 Trieste, Italy, ⁴Translational Neuromodeling Unit, Institute for Biomedical Engineering, University of Zurich and ETH Zurich, 8032 Zurich, Switzerland, ⁵Max Planck Centre for Computational Psychiatry and Ageing Research, University College London, London WC1B 5EH, United Kingdom, ⁶Department of Psychology, University of Roehampton, London SE15 4JD, United Kingdom, and ⁷Department of Psychosis Studies, Institute of Psychiatry, Psychology and Neuroscience, Kings College London, London SE5 8AF, United Kingdom

Subjects with a diagnosis of schizophrenia (Scz) overweight unexpected evidence in probabilistic inference: such evidence becomes “aberrantly salient.” A neurobiological explanation for this effect is that diminished synaptic gain (e.g., hypofunction of cortical NMDARs) in Scz destabilizes quasi-stable neuronal network states (or “attractors”). This attractor instability account predicts that (1) Scz would overweight unexpected evidence but underweight consistent evidence, (2) belief updating would be more vulnerable to stochastic fluctuations in neural activity, and (3) these effects would correlate. Hierarchical Bayesian belief updating models were tested in two independent datasets ($n = 80$ male and $n = 167$ female) comprising human subjects with Scz, and both clinical and nonclinical controls (some tested when unwell and on recovery) performing the “probability estimates” version of the beads task (a probabilistic inference task). Models with a standard learning rate, or including a parameter increasing updating to “disconfirmatory evidence,” or a parameter encoding belief instability were formally compared. The “belief instability” model (based on the principles of attractor dynamics) had most evidence in all groups in both datasets. Two of four parameters differed between Scz and nonclinical controls in each dataset: belief instability and response stochasticity. These parameters correlated in both datasets. Furthermore, the clinical controls showed similar parameter distributions to Scz when unwell, but were no different from controls once recovered. These findings are consistent with the hypothesis that attractor network instability contributes to belief updating abnormalities in Scz, and suggest that similar changes may exist during acute illness in other psychiatric conditions.

Key words: attractor model; Bayesian; beads task; disconfirmatory bias; psychosis; schizophrenia

Significance Statement

Subjects with a diagnosis of schizophrenia (Scz) make large adjustments to their beliefs following unexpected evidence, but also smaller adjustments than controls following consistent evidence. This has previously been construed as a bias toward “disconfirmatory” information, but a more mechanistic explanation may be that in Scz, neural firing patterns (“attractor states”) are less stable and hence easily altered in response to both new evidence and stochastic neural firing. We model belief updating in Scz and controls in two independent datasets using a hierarchical Bayesian model, and show that all subjects are best fit by a model containing a belief instability parameter. Both this and a response stochasticity parameter are consistently altered in Scz, as the unstable attractor hypothesis predicts.

Introduction

Subjects with a diagnosis of schizophrenia (Scz) tend to use less evidence to make decisions in probabilistic tasks than healthy

controls (Garety et al., 1991; Dudley et al., 2016). The paradigm most commonly used to demonstrate this effect is the ‘beads’ or ‘urn’ task, in which subjects are shown two urns, each containing opposite ratios of colored beads (e.g., 85% blue and 15% red and vice versa), which are then hidden. A sequence of beads is then

Received Nov. 2, 2017; revised May 3, 2018; accepted June 27, 2018.

Author contributions: R.A.A. wrote the first draft of the paper; R.A.A., J.P.R., C.M., and J.G. designed research; R.A.A., G.N., and J.G. performed research; C.M. contributed unpublished reagents/analytic tools; R.A.A. and G.N. analyzed data; R.A.A., G.N., J.P.R., C.M., and J.G. wrote the paper.

R.A.A. was supported by Academy of Medical Sciences AMS-SGCL13-Adams and National Institute of Health Research CL-2013-18-003. J.G. was supported by the British Academy. We thank Dr. Emmanuelle Peters for providing Dataset 1.

The authors declare no competing financial interests.

*R.A.A. and G.N. equal contributed equally to this work.

†C.M. and J.G. contributed equally to this work as joint senior authors.

Correspondence should be addressed to Dr. Rick A. Adams, Institute of Cognitive Neuroscience, University College London, 17 Queen Square, London WC1N 3AZ, United Kingdom. E-mail: rick.adams@ucl.ac.uk.

DOI:10.1523/JNEUROSCI.3163-17.2018

Copyright © 2018 the authors 0270-6474/18/389471-15\$15.00/0

drawn (with replacement) from one urn, and the subject either has to stop the sequence when they are sure which urn it is coming from (the ‘draws to decision’ task) or the subject must rate the probability of the sequence coming from either urn after seeing each bead, without having to make any decision (the ‘probability estimates’ task). Bayesian analysis of these tasks has indicated that Scz are more stochastic in their responding (Moutoussis et al., 2011) and that they overweight recent evidence and thus update their beliefs (in the probabilistic sense) more rapidly (Jardri et al., 2017).

Several belief-updating abnormalities have been found in Scz using the ‘probability estimates’ task. The most consistent finding is that Scz (or just Scz with delusions) (Moritz and Woodward, 2005) change their beliefs more than nonpsychiatric controls in response to changes in evidence (Langdon et al., 2010), particularly ‘disconfirmatory’ evidence (i.e., evidence contradicting a current belief) (Garety et al., 1991; Fear and Healy, 1997; Young and Bentall, 1997; Peters and Garety, 2006). Another is that probability ratings at the start of the sequence are higher in currently psychotic (but not in recovered) Scz than in both clinical and healthy controls (Peters and Garety, 2006), similar to the ‘jumping to conclusions’ bias in the ‘draws to decision’ version of the task. Others have also found that Scz update less than controls to more consistent evidence, in this (Baker et al., 2018) and other paradigms (Averbeck et al., 2011).

These findings can potentially be understood in the light of the ‘unstable attractor network’ hypothesis of Scz. An attractor network is a neural network that can occupy numerous stable states that are learned from experience, via adjustments to synaptic weights. It can revisit these states if presented with inputs that resemble previous patterns of synaptic weights, or through spontaneous fluctuations in neural activity: either way, the activity of all nodes is ‘attracted’ to a quasi-stable state because the network energy is lower at these states, and network firing patterns evolve to minimize energy. Attractor networks were originally developed to model the storage and reactivation of memories (Hopfield, 1982), but related network models also offer mechanistic explanations for working memory storage (e.g., Brunel and Wang, 2001), decision-making (Wang, 2013), and interval timing (Standage et al., 2013), as well as Bayesian belief updating (Geperth and Lefort, 2016).

In Scz, attractor states in prefrontal cortex are thought to be less stable, so it is easier for the network to switch between them, but harder to become more confident about (i.e., increase the stability of) any particular one (Rolls et al., 2008). This loss of stable neuronal states, recently demonstrated in two animal models of Scz (Hamm et al., 2017), is thought to be due to hypofunction of NMDARs or cortical dopamine 1 receptors in Scz (Fig. 1). Interestingly, healthy volunteers given ketamine (an NMDAR antagonist) show a decrement in updating to consistent stimulus associations and an increase in decision stochasticity in

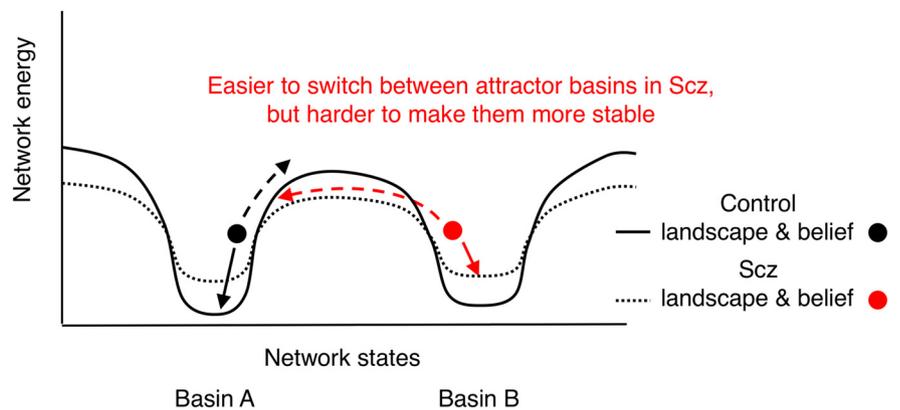


Figure 1. Effects of attractor network dynamics on belief updating. This schematic illustrates the energy landscapes of two Hopfield-type networks each with two basins of attraction. Continuous black line indicates a normal network whose basins of attraction are relatively deep. Dotted black line indicates the effect of NMDAR (or cortical dopamine 1 receptor) (Durstewitz and Seamans, 2008; Redish et al., 2007) hypofunction (Abi-Saab et al., 1998; Javitt et al., 2012) on the energy landscape: the attractor basins become more shallow. We assume that Basins A and B correspond to different inferences about (hidden) states in the world (e.g., one jar or another being the source of beads in the beads task). Dots indicate the networks’ representations of either control or Scz subjects’ beliefs about these hidden states. Such networks are highly reminiscent of Hopfield networks with two stored representations; in this case, the representations correspond to inferences about hidden states, rather than memories. Arrows indicate the changes in network states resulting from sensory evidence for (solid arrows) or against (dashed arrows) the current inference. When the attractor basin is shallower, it is harder for supportive evidence to stabilize the current state much further, but it is easier for contradictory evidence, or just stochastic neuronal firing, to shift the current network state toward an alternative state. These changes in network dynamics may also be reflected in the inferences the network computes (i.e., easier switching between attractor basins may correspond to easier switching between beliefs), although this is yet to be demonstrated experimentally. NMDAR hypofunction could contribute to an increased tendency to switch between beliefs and increased stochasticity in responding in several ways (Rolls et al., 2008): (1) by reducing inhibitory interneuron activity, via weakened NMDAR synapses from pyramidal cells to interneurons, such that other attractor states are less suppressed when one is active (a spiking network model has shown that this leads to more rapid initial belief updating in perceptual tasks) (Lam et al., 2017); (2) by reducing pyramidal cell activity, via weakened recurrent NMDAR synapses on pyramidal cells, such that attractor states are harder to sustain; and (3) by reducing the NMDAR time constant, making states more vulnerable to random fluctuations in neural activity. See also similar schematics elsewhere (Durstewitz and Seamans, 2008; Rolls et al., 2008).

this context (Vinckier et al., 2016). Attractor network perturbations have been linked to working memory problems in Scz using a bistable (i.e., a stable ‘up’ state corresponding to persistent neuronal activity, and a ‘down’ state corresponding to background activity) model (Murray et al., 2014), but not as yet to a computational understanding of belief updating.

We analyzed belief updating in Scz using the Hierarchical Gaussian Filter (HGF) (Mathys et al., 2011), a variational Bayesian model with individual priors, in two independent ‘probability estimates’ beads task datasets. We asked: given the larger belief updates in Scz compared with controls, can these be explained by group differences in (1) general learning rate and/or (2) response stochasticity, or by adding parameters encoding (3) the variance (i.e., uncertainty) of beliefs at the start of the sequence, (4) a propensity to overweight disconfirmatory evidence specifically, or (5) patterns of belief updating typical of unstable attractor states in a Hopfield-type network (i.e., greater instability and stochasticity), which correlate with each other? The HGF does not contain attractor states: the model in (5) is designed to simulate the effects on inference that unstable neuronal attractors may have. Furthermore, are these findings consistent within Scz tested at different illness phases, and are they unique to Scz or also present in other nonpsychotic mood disorders?

Materials and Methods

Subject characteristics. Dataset 1 comprised 23 patients with delusions (18 Scz), 22 patients with nonpsychotic mood disorders, and 35 nonclinical controls (overall, 50 male and 30 female; for details of the groups, see Tables 1, 2); the first two groups were selected from inpatient wards at the

Table 1. Demographic, psychological, and behavioral details of Dataset 1^a

	Nonclinical controls t1	Nonclinical controls t2	Clinical controls t1	Clinical controls t2	Psychotic t1	Psychotic t2
<i>N</i>	35	20	22	18	23	17
Age ^b (yr)	27.77 (6.74)	27.9 (6.37)	40.91 (13.57)	40.1 (13)	31.22 (7.28)	29.9 (7.83)
Gender	18 M, 17 F	12 M, 8 F	11 M, 11 F	8 M, 10 F	21 M, 2 F	17 M, 0 F
Cognitive measures						
IQ ^c	107.5 (11.6)	108.6 (10.3)	97.4 (13.8)	99.8 (10.2)	88.1 (12.7)	87.8 (14.2)
Delusion proneness						
PDI (total) ^d	54.6 (43.1)	43.6 (42.5)	87.1 (55.2)	64.3 (57.3)	138.1 (74.2)	96.7 (42.6)
DSSI ^e	2.3 (4.9)	2.9 (5.3)	4.8 (4.5)	4.5 (5.6)	15.2 (6.3)	8.1 (6.6)
Diagnosis/symptoms						
Diagnoses			16 Depression, 3 anxiety and depression, 3 SAD	12 Depression, 3 anxiety and depression, 3 SAD	18 Scz, 5 bipolar/schizo-affective	13 Scz, 4 bipolar/schizo-affective
MS affective	—	—	4.6 (1.7)	1.0 (1.2)	1.8 (1.5)	1.5 (1.3)
MS positive	—	—	0.3 (0.8)	0 (0)	6.0 (2.4)	1.4 (1.7)
MS negative	—	—	0.7 (1.6)	1.8 (3.19)	1.3 (2.0)	0.9 (1.6)
MS total ^f	—	—	5.5 (2.6)	2.8 (3.39)	9.1 (3.76)	3.7 (3.9)
Beads task						
Initial certainty (1 bead) ^g	0.58 (0.15)	0.59 (0.12)	0.68 (0.19)	0.63 (0.16)	0.76 (0.17)	0.68 (0.29)
Initial certainty (3 beads) ^h	0.65 (0.14)	0.67 (0.1)	0.69 (0.15)	0.64 (0.16)	0.78 (0.15)	0.74 (0.15)
Disconfirmatory updating ⁱ	−0.06 (0.14)	−0.03 (0.13)	−0.19 (0.3)	−0.11 (0.22)	−0.29 (0.33)	−0.2 (0.3)
Final certainty ^j	0.85 (0.2)	0.94 (0.11)	0.82 (0.16)	0.79 (0.23)	0.88 (0.11)	0.85 (0.23)

^aDataset 1 includes measures at both baseline (t1) and follow-up (t2). In Dataset 1, verbal IQ was estimated using the Quick Test (Ammons and Ammons, 1962) and delusion proneness using the Peters Delusion Inventory (PDI) (Peters et al., 1999) and Delusions-Symptoms-States Inventory (DSSI) (Foulds and Bedford, 1975). Symptoms were assessed using the Manchester Scale (MS) (Krawiecka et al., 1977). In the tests below, ‘Scz’ refers to the whole Psychotic group. Results are given for ‘Initial certainty’ using both the measure in the original analysis of Dataset 1 (Peters and Garety, 2006), the mean response to the first three beads (3 beads); in Dataset 2, this had to be the mean response to the first three beads in Sequences B and C and two beads in Sequences A and D (2–3 beads), and using the response to the first bead (1 bead).

^bAt t1: one-way ANOVA $F_{(2,77)} = 13.9, p = 10^{-5}$; Tukey’s HSD: Scz versus Nonclinical controls diff = 3.45, $p(adj) = 0.35$; Clinical versus Nonclinical controls diff = 13.1, $p(adj) = 10^{-5}$; Clinical controls versus Scz diff = 9.69, $p(adj) = 0.002$. At t2: one-way ANOVA $F_{(2,52)} = 8.85, p = 0.0005$. Tukey’s HSD: Scz versus Nonclinical controls diff = 1.98, $p(adj) = 0.8$; Clinical versus Nonclinical controls diff = 12.2, $p(adj) = 0.0006$; Clinical controls versus Scz diff = 10.2, $p(adj) = 0.007$.

^cAt t1: one-way ANOVA $F_{(2,75)} = 16.2, p = 10^{-6}$; Tukey’s HSD: Scz versus Nonclinical controls diff = −19.5, $p(adj) = 10^{-6}$; Clinical versus Nonclinical controls diff = −10.1, $p(adj) = 0.011$; Clinical controls versus Scz diff = 9.36, $p(adj) = 0.043$. At t2: one-way ANOVA $F_{(2,51)} = 14.5, p = 10^{-5}$; Tukey’s HSD: Scz versus Nonclinical controls diff = −20.8, $p(adj) = 10^{-5}$; Clinical versus Nonclinical controls diff = −8.8, $p(adj) = 0.057$; Clinical controls versus Scz diff = 12, $p(adj) = 0.01$.

^dAt t1: one-way ANOVA $F_{(2,68)} = 12.6, p = 0.00002$; Tukey’s HSD: Scz versus Nonclinical controls diff = 83.5, $p(adj) = 10^{-5}$; Clinical versus Nonclinical controls diff = −32.5, $p(adj) = 0.094$; Clinical controls versus Scz diff = −51, $p(adj) = 0.016$. At t2: one-way ANOVA $F_{(2,52)} = 4, p = 0.024$; Tukey’s HSD: Scz versus Nonclinical controls diff = 53.1, $p(adj) = 0.018$; Clinical versus Nonclinical controls diff = −20.7, $p(adj) = 0.5$; Clinical controls versus Scz diff = −32.4, $p(adj) = 0.22$.

^eAt t1: one-way ANOVA $F_{(2,76)} = 43, p = 10^{-13}$; Tukey’s HSD: Scz versus Nonclinical controls diff = 12.9, $p(adj) = 10^{-10}$; Clinical versus Nonclinical controls diff = 2.52, $p(adj) = 0.19$; Clinical controls versus Scz diff = −10.4, $p(adj) = 10^{-8}$. At t2: one-way ANOVA $F_{(2,51)} = 3.7, p = 0.032$; Tukey’s HSD: Scz versus Nonclinical controls diff = 5.2, $p(adj) = 0.026$; Clinical versus Nonclinical controls diff = 1.65, $p(adj) = 0.66$; Clinical controls versus Scz diff = −3.56, $p(adj) = 0.18$.

^fAt t1: Welch’s $t_{(38.4)} = -3.62, p = 0.00086$, Cohen’s $d = 1.1$. At t2: Welch’s $t_{(17.8)} = -2.55, p = 0.02$, Cohen’s $d = 1.0$.

^gAt t1: one-way ANOVA $F_{(2,77)} = 8.7, p = 0.0004$; Tukey’s HSD: Scz versus Nonclinical controls diff = 0.18, $p(adj) = 0.0003$; Clinical versus Nonclinical controls diff = 0.11, $p = 0.06$; Clinical controls versus Scz diff = −0.08, $p(adj) = 0.25$. At t2: one-way ANOVA $F_{(2,52)} = 0.9, p = 0.4$.

^hAt t1: one-way ANOVA $F_{(2,77)} = 6.2, p = 0.003$; Tukey’s HSD: Scz versus Nonclinical controls diff = −0.14, $p(adj) = 0.002$; Clinical versus Nonclinical controls diff = 0.04, $p = 0.57$; Clinical controls versus Scz diff = −0.096, $p(adj) = 0.074$. At t2: one-way ANOVA $F_{(2,52)} = 2.35, p = 0.11$; Tukey’s HSD: Scz versus Nonclinical controls diff = 0.07, $p(adj) = 0.28$; Clinical versus Nonclinical controls diff = −0.03, $p = 0.8$; Clinical controls versus Scz diff = −0.1, $p(adj) = 0.1$.

ⁱAt t1: one-way ANOVA $F_{(2,77)} = 6, p = 0.004$; Tukey’s HSD: Scz versus Nonclinical controls diff = −0.23, $p(adj) = 0.003$; Clinical versus Nonclinical controls diff = −0.14, $p = 0.13$; Clinical controls versus Scz diff = 0.097, $p(adj) = 0.41$. At t2: one-way ANOVA $F_{(2,52)} = 2.9, p = 0.062$; Tukey’s HSD: Scz versus Nonclinical controls diff = −0.18, $p(adj) = 0.049$; Clinical versus Nonclinical controls diff = −0.08, $p = 0.51$; Clinical controls versus Scz diff = 0.098, $p(adj) = 0.4$.

^jAt t1: one-way ANOVA $F_{(2,77)} = 0.71, p = 0.5$. At t2: one-way ANOVA $F_{(2,52)} = 2.79, p = 0.07$; Tukey’s HSD: Scz versus Nonclinical controls diff = −0.082, $p(adj) = 0.41$; Clinical versus Nonclinical controls diff = −0.15, $p = 0.057$; Clinical controls versus Scz diff = −0.066, $p(adj) = 0.57$. As reported previously, there were consistent negative correlations between initial certainty (2–3 beads) and disconfirmatory updating in the clinical controls (baseline: $\rho = -0.68, p = 0.0005$; follow-up: $\rho = -0.75, p = 0.0003$) and the nonclinical controls (baseline: $\rho = -0.52, p = 0.001$; follow-up: $\rho = -0.43, p = 0.06$), but not in the psychotic group (baseline: $\rho = -0.30, p = 0.17$; follow-up: $\rho = 0.17, p = 0.5$). There was no consistent correlation between final certainty and either of the other two measures at either time point ($p \geq 0.1$ in 11 of 12 comparisons).

Maudsley and the Bethlem Royal Hospitals. All groups were tested twice (with loss of $n = 25$ from the groups; Tables 1, 2); the clinical groups were tested once when they were unwell (‘baseline’), and again once they had recovered (‘follow-up’). The mean time between testing sessions was 17.4 (range 6–41) weeks in the deluded group, 33.4 (range 4–68) weeks in the clinical control group, and 35.6 (range 27–46) weeks in the nonclinical control group. The deluded group’s shorter intertest interval was due to their shorter admission period and to the prioritization of their follow-up over the nonclinical control group. Dataset 1 was described in detail previously (Peters and Garety, 2006).

Dataset 2 comprised 56 subjects with a diagnosis of Scz and 111 controls (overall, 83 male and 84 female; Tables 1, 2). All subjects provided informed, written consent, and ethical permission for the study was obtained from the local NHS Research Ethics Committee (Reference 14/LO/0532). Given the National Adult Reading Test (Nelson, 1982) was used to estimate IQ in these participants, a recruitment condition was that English was their first language.

Measures of cognitive function and delusion-proneness (or schizotypy) were collected in all subjects; clinical symptom ratings were collected in clinical subjects only (for details, see Tables 1, 2).

Experimental design. Subjects in Dataset 1 performed the ‘probability estimates’ beads task as used previously (Garety et al., 1991), with two urns with ratios of 85:15 and 15:85 blue and red beads, respectively, and viewing a single sequence of 10-beads (Fig. 2); after each bead, they had to mark an analog scale (from 1 to 100) denoting the probability the urn was 85% red.

Subjects in Dataset 2 performed the ‘probability estimates’ beads task, with two urns with ratios of 80:20 and 20:80 red and blue beads, respectively. They each viewed four separate sequences (two identical pairs of sequences with the colors swapped within each pair) of 10-beads (Fig. 2); after each bead, they had to mark a Likert scale (from 1 to 7) denoting the probability the urn was the 80% blue one. Two sequences contained an apparent change of jar. The order of the four sequences was randomized.

We used some of the behavioral measures used in the original analysis of Dataset 1 (Peters and Garety, 2006) to analyze Dataset 2. These were ‘disconfirmatory updating,’ the mean change in belief on seeing a bead of a different color to the ≥ 2 beads preceding it and ‘final certainty’ (the response to the last bead). We altered their ‘initial certainty’ measure from the mean response to the first three beads to the response to the first bead, which comes closer to capturing the classic ‘jumping to conclu-

Table 2. Demographic, psychological, and behavioral details of Dataset 2^a

	Controls (all)	Scz	Controls (subset)
<i>N</i>	111	56	60
Age (yr)	32.8 (11.5)	45.3 (8.8)	39.5 (11.4)
Gender	45 M, 66 F	38 M, 18 F	40 M, 20 F
NART ^b	112 (6.9)	109 (8.2)	112 (7.5)
Working memory (LNS) ^c	16.2 (2.8)	10.3 (4.2)	16.4 (2.7)
Schizotypy			
SPQ, cognitive	2.8 (1.9)	4.0 (2.6)	3.1 (2)
SPQ, interpers	3.2 (2.2)	5.3 (2.6)	3.2 (2.2)
SPQ, disorg	2.1 (1.7)	2.7 (1.9)	1.9 (1.8)
SPQ, total ^d	8.2 (1.3)	12 (5.3)	8.2 (4.4)
Diagnoses	—	56 Scz	—
PANSS, gen	—	32.6 (9.2)	—
PANSS, pos	—	15.9 (5.8)	—
PANSS, neg	—	15.9 (6.2)	—
PANSS, total	—	64.4 (17.3)	—
Initial certainty (all, 1 bead) ^e	0.67 (0.13)	0.71 (0.14)	0.68 (0.14)
Initial certainty (all, 2–3 beads) ^f	0.7 (0.12)	0.71 (0.12)	0.71 (0.13)
Disconfirmatory updating (all sequences) ^g	−0.16 (0.17)	−0.23 (0.22)	−0.19 (0.2)
Final certainty Sequence A ^h	0.88 (0.16)	0.77 (0.25)	0.86 (0.18)
Final certainty Sequence D ⁱ	0.12 (0.18)	0.25 (0.24)	0.16 (0.2)

^aIn dataset 2, IQ was estimated using the National Adult Reading Test (NART) (Nelson, 1982) and working memory using the Letter Number Sequencing task (LNS) from the Wechsler Adult Intelligence Scale-III (Wechsler, 1997). Schizotypy was assessed using the Schizotypal Personality Questionnaire (SPQ) (Raine, 1991), and symptoms using the Positive and Negative Syndrome Scale (PANSS) (Kay et al., 1987). As can be seen in Figure 2, the Scz group showed greater initial certainty (1 bead) in Sequences A and B (Welch's $t_{(94)} = 2.8, p = 0.007$, Cohen's $d = 0.47$; Welch's $t_{(97)} = 3, p = 0.004$, Cohen's $d = 0.5$, respectively) but not Sequences C and D (Welch's $t_{(87)} = 0.5, p = 0.6$, Cohen's $d = 0.09$; Welch's $t_{(90)} = -0.34, p = 0.73$, Cohen's $d = 0.06$, respectively).

^bControls (all): Welch's $t_{(95,1)} = 2.27, p = 0.026$, Cohen's $d = 0.38$; Controls (subset): Welch's $t_{(111)} = 1.95, p = 0.053$, Cohen's $d = 0.36$.

^cControls (all): Welch's $t_{(81)} = 9.57, p = 10^{-14}$, Cohen's $d = 1.66$; Controls (subset): Welch's $t_{(93,6)} = 9.25, p = 10^{-15}$, Cohen's $d = 1.73$.

^dControls (all): Welch's $t_{(92,4)} = -4.64, p = 10^{-5}$, Cohen's $d = 0.78$; Controls (subset): Welch's $t_{(107)} = -4.19, p = 10^{-5}$, Cohen's $d = 0.78$.

^eControls (all): Welch's $t_{(110)} = -1.9, p = 0.059$, Cohen's $d = 0.32$; Controls (subset): Welch's $t_{(110)} = -1.1, p = 0.28$, Cohen's $d = 0.2$.

^fControls (all): Welch's $t_{(109,1)} = -0.76, p = 0.45$, Cohen's $d = 0.12$; Controls (subset): Welch's $t_{(113,9)} = -0.19, p = 0.85$, Cohen's $d = 0.03$.

^gControls (all): Welch's $t_{(88,2)} = 2.09, p = 0.04$, Cohen's $d = 0.36$; Controls (subset): Welch's $t_{(110,4)} = -0.94, p = 0.35$, Cohen's $d = 0.18$.

^hControls (all): Welch's $t_{(80,1)} = 2.99, p = 0.0038$, Cohen's $d = 0.56$; Controls (subset): Welch's $t_{(98,7)} = 2.18, p = 0.032$, Cohen's $d = 0.41$.

ⁱControls (all): Welch's $t_{(85,5)} = -3.41, p = 0.001$, Cohen's $d = 0.62$; Controls (subset): Welch's $t_{(106)} = -2.21, p = 0.029$, Cohen's $d = 0.42$.

sions' bias (in which ~50% of Scz decide on the jar color after seeing only one bead) (Garety et al., 1991), although the results of both measures are presented below.

Computational modeling. The optimal way to use sensory information to update one's beliefs under conditions of uncertainty is to use Bayesian inference. Neural systems are likely to approximate Bayesian inference using schemes of simple update equations (Rao and Ballard, 1999; Friston, 2005); one such model is the HGF. The HGF is a hierarchical Bayesian inference scheme that gives a principled account of how beliefs are updated on acquiring new data, using variational Bayes and individual priors. Variational Bayesian schemes (e.g., Beal, 2003) use analytic equations to derive an exact solution to an approximation of the posterior distribution over the latent variables and parameters (as opposed to sampling methods, which approximate a solution to the exact posterior). The HGF has been used as a generic state model for learning under uncertainty and has repeatedly been shown to outperform similar approaches, such as reinforcement learning models with fixed (e.g., Rescorla-Wagner) or dynamic (e.g., Sutton, 1992) learning rates (Iglesias et al., 2013; Diaconescu et al., 2014; Hauser et al., 2014; Vossel et al., 2014). One advantage of the HGF is that it contains subject-specific parameters (and prior beliefs) that can account for between-subject differences in learning while preserving the (Bayes) optimality of any individual's learning (relative to his/her model parameters and prior beliefs). These parameters

may be encoded by tonic levels of neuromodulators, such as dopamine (Marshall et al., 2016), or by the intrinsic properties of neuronal networks (e.g., the ratio of excitatory to inhibitory neural activity can affect the speed of evidence accumulation) (Lam et al., 2017), analogous to the evolution rate in the HGF and also response stochasticity (Murray et al., 2014). Differences in model parameters between Scz and controls may therefore explain, in computational terms, how pathophysiology leads to abnormal inference (Adams et al., 2016).

In general, when modeling behavior under Bayesian assumptions, it is necessary to distinguish between the model of the world used by the subject (the perceptual model) and a model of how a subject's beliefs translated into observed behavior (the observation or response model). Most of the parameters pertain to the perceptual model (here, all parameters except response stochasticity ν ; Table 3) and reflect (inferred) neuronal processing. In contrast, the parameters of the response model link subjective states to behavioral outcomes, and thus may reflect stochasticity in neuronal processing, measurement noise (in some paradigms), or nonrandom effects that have not been captured by the perceptual model. This and related learning models are freely available from <http://www.translationalneuromodeling.org/tapas/> (version 5.1.0): this analysis used the perceptual models 'hgf_binary' or 'hgf_ar1_binary' and the response model 'beta_obs.'

At the bottom of the model (Fig. 3 shows some simulated responses) is the bead drawn $u^{(k)}$ on trial k and the probability $x_1^{(k)}$ that draws are coming from the blue jar. At the level above this is x_2 , the tendency toward the blue jar (a transform of the probability, bounded by $\pm\infty$); by definition, $x_1 = s(x_2)$, where $s(\bullet)$ is the logistic sigmoid function. As x_2 approaches infinity, the probability of the blue jar approaches 1; as it approaches minus infinity, the probability of the blue jar approaches 0. For $x_2 = 0$, both jars are equally probable. This quantity is hidden from the subject and must be inferred: the subject's posterior estimate of x_2 is μ_2 , and the subject's posterior estimate of the probability of the jar being blue on trial k is $s(\mu_2^{(k)})$, equivalent to the prediction (denoted by \wedge) on the next trial $\hat{\mu}_1^{(k+1)}$.

Before seeing any new input on trial k , the model's expected jar probability $\hat{\mu}_1^{(k)}$ and precisions (inverse variances) $\hat{\pi}_1^{(k)}, \hat{\pi}_2^{(k)}$ of the expectations at each level are given by the following:

$$\hat{\mu}_1^{(k)} \equiv s(\kappa_1 \mu_2^{(k-1)})$$

$$\hat{\pi}_1^{(k)} \equiv \frac{1}{\hat{\mu}_1^{(k)}(1 - \hat{\mu}_1^{(k)})}$$

$$\hat{\pi}_2^{(k)} \equiv \frac{1}{\sigma_2^{(k-1)} + \exp(\omega)}$$

In Models 1–4, κ_1 is fixed to 1. A new input $u^{(k)} \equiv \mu_1^{(k)}$ generates a prediction error $\delta_1^{(k)}$, and the model updates and generates a new prediction as follows:

$$\delta_1^{(k)} \equiv \mu_1^{(k)} - \hat{\mu}_1^{(k)}$$

$$\pi_2^{(k)} = \hat{\pi}_2^{(k)} + \frac{\kappa_1^2}{\hat{\pi}_1^{(k)}}$$

$$\mu_2^{(k)} = \mu_2^{(k-1)} + \frac{\kappa_1}{\pi_2^{(k)}} \delta_1^{(k)}$$

$$\hat{\mu}_1^{(k+1)} \equiv s(\kappa_1 \mu_2^{(k)})$$

The subject's response $y^{(k)}$ (i.e., where on the continuous or Likert scale they responded) is determined by $\hat{\mu}_1^{(k+1)}$ and the precision of the response model's β distribution ν .

We parameterize the β distribution in terms of its mean μ and precision ν . These sufficient statistics relate to the conventional parameterization in terms of the sufficient statistics α and β by the following bijection:

$$\mu: = \frac{\alpha}{\alpha + \beta}$$

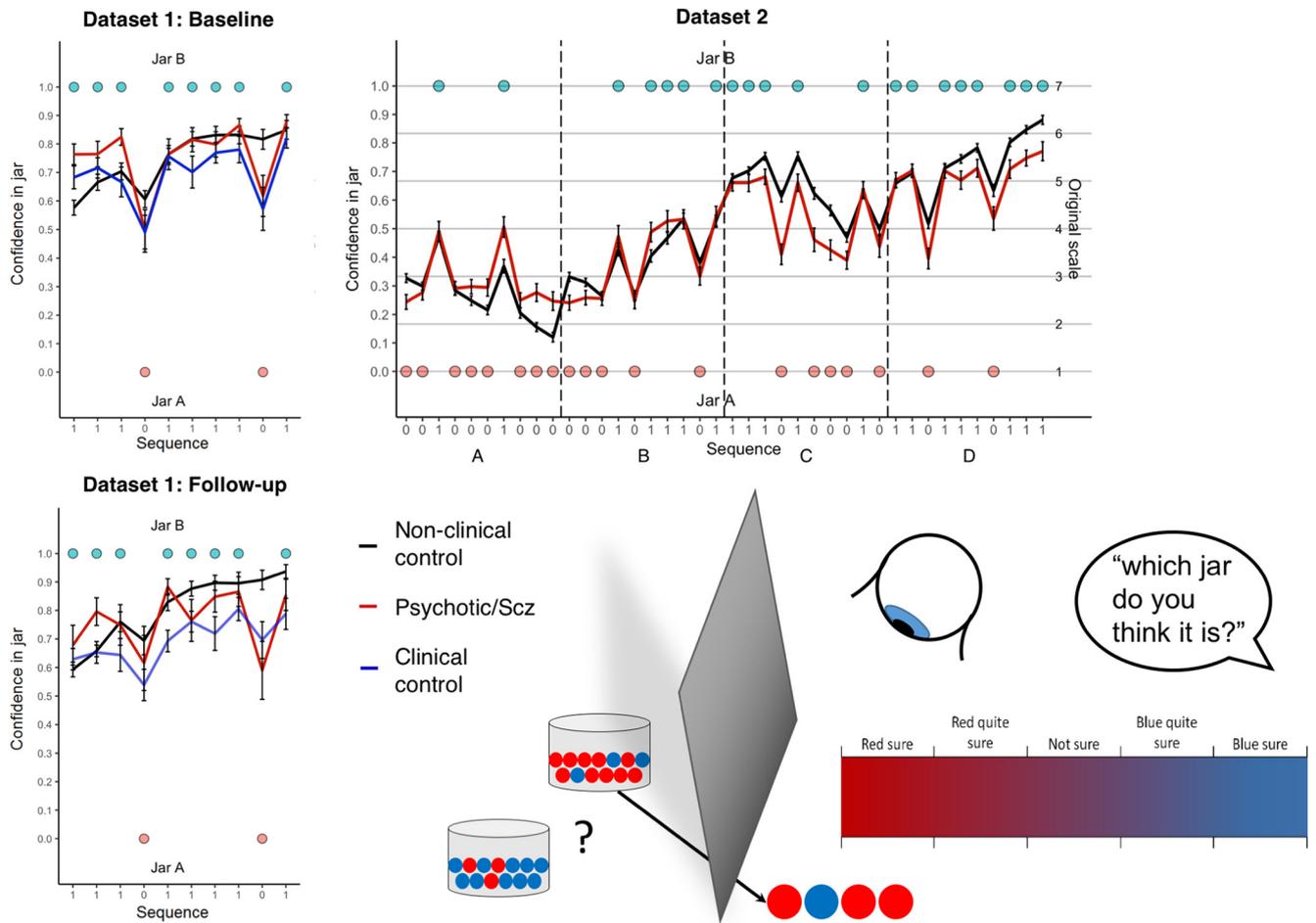


Figure 2. Beads task schematic and group average confidence ratings in Datasets 1 and 2. Bottom right, Schematic of the beads task: two jars containing opposite proportions of beads are concealed from view, and a subject is asked to rate the probability of either jar being the source of a sequence of beads he/she is viewing (after each bead in turn). Top left, Mean (\pm SE) confidence ratings in the blue jar over the 10-bead sequence averaged across each group at baseline in Dataset 1. Bottom left, Same quantities at follow-up in Dataset 1. Top right, Quantities in four 10-bead sequences concatenated together (they were presented to the subjects separately during testing) in Dataset 2.

Table 3. Models, parameters, and their prior distributions

Model	Perceptual model parameters (prior mean in native space, prior variance in estimation space)				Response model parameter: response stochasticity
	Evolution rate	Initial variance of belief regarding jars	Disconfirmatory bias	Belief instability	
1	$\omega (-2, 16)$				$\nu (\exp(4.85), 1)$
2	$\omega (-2, 16)$	$\sigma_2^{(0)} (0.8, 0.5)$			$\nu (\exp(4.85), 1)$
3	$\omega (-2, 16)$		$\varphi (0.1, 2)$		$\nu (\exp(4.85), 1)$
4	$\omega (-2, 16)$	$\sigma_2^{(0)} (0.8, 0.5)$	$\varphi (0.1, 2)$		$\nu (\exp(4.85), 1)$
5	$\omega (-2, 16)$			$\kappa_1 (1, 1)$	$\nu (\exp(4.85), 1)$
6	$\omega (-2, 16)$	$\sigma_2^{(0)} (0.8, 0.5)$		$\kappa_1 (1, 1)$	$\nu (\exp(4.85), 1)$

$$\nu = \alpha + \beta$$

Updates to μ_2 are driven by the product of the prediction error from Bayesian updating explained above and a learning rate which, crucially, can change over time: this is an important aspect of the HGF in contrast to learning models, such as Rescorla-Wagner, which have a fixed learning rate. Parameters that affect the degree to which μ_2 can change during the experiment include ω , φ , κ_1 , and $\sigma_2^{(0)}$. The contributions of φ and κ_1 are illustrated in Figure 4 (left panels).

The model usually has a third level, at which x_3 encodes the phasic volatility of x_2 (this determines the probability of the jar changing at any point): given the very short sequences used in our datasets, from which volatility cannot be reliably estimated, we omitted this level. In any case, volatility could not account for the rapid changes in learning rate (from trial to trial, following confirmatory vs disconfirmatory evidence) present in the Scz group in these datasets.

In Models 1 and 2, changes in x_2 from trial to trial occur only according to the evolution rate ω , the variance of the random process at the second level. These models were equivalent to the subsequent models with either φ (Models 3 and 4) fixed to 0 or κ_1 (Models 5 and 6) fixed to 1.

In Models 3 and 4, changes in x_2 from trial to trial occur according to an autoregressive (AR(1)) process that is controlled by three parameters: m , the level to which x_2 is attracted; φ , the rate of change of x_2 toward m ; and ω , the variance of the random process as follows:

$$p(x_2^{(k+1)}) \sim N(x_2^{(k)} + \varphi(m - x_2^{(k)}), \exp(\omega))$$

After inversion, the evolution of x_2 according to this equation is reflected in the prediction of μ_2 as follows:

$$\hat{\mu}_2^{(k+1)} = \mu_2^{(k)} + \varphi(m - \mu_2^{(k)})$$

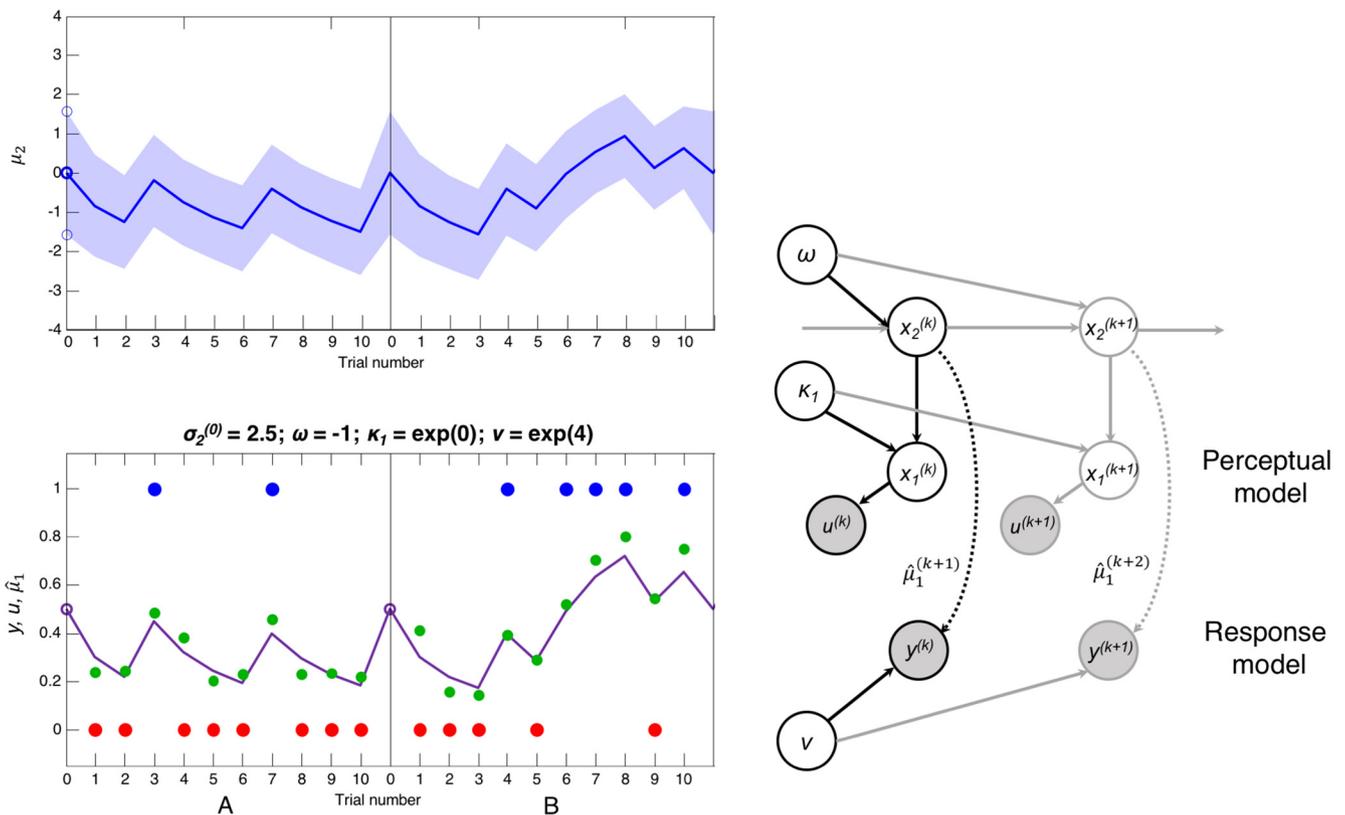


Figure 3. The structure of the HGF (Model 6) and some simulated data. Top left, The evolution of μ_2 , the posterior estimate of tendency x_2 toward the blue (positive) or red (negative) jar, is plotted over two concatenated series of 10 trials (the first two in Dataset 2). The estimate of the tendency on trial $k + 1$, $\mu_2^{(k+1)}$, is selected from a Gaussian distribution with mean $\mu_2^{(k)}$ (blue line) and variance $\sigma_2^{(k)} + \exp(\omega)$ (blue shading). ω is a static source of variance at this level. The initial variance $\sigma_2^{(0)}$ (along with ω) affects the size of initial updates, so we estimated this parameter (which is often fixed). Bottom left, The beads seen by the subjects, $u^{(k)}$ (blue and red dots) and the response model. The response model maps from $\mu_1^{(k+1)}$ (purple line), the prediction of x_1 on the next trial, which is a sigmoid function s of $\mu_2^{(k)}$ (or of $\kappa_1 \mu_2^{(k)}$ in Models 5 and 6), to $y^{(k)}$, the subject’s indicated estimate of the probability the jar is blue (green dots). Variation in this mapping is modeled as the precision ν of a β distribution. Right, Schematic representation of the generative model in Models 5 and 6 (i.e., including κ_1). Black arrows indicate the probabilistic network on trial k . Gray arrows indicate the network at other points in time. The perceptual model lies above the dotted arrows, and the response model below them. Shaded circles represent known quantities. Unshaded circles represent estimated parameters and states. Dotted line indicates the result of an inferential process (the response model builds on a perceptual model inference). Solid lines indicate generative processes.

In this study, given there was no bias toward one jar or the other, m was fixed to 0, so φ always acted to shift the model’s beliefs back toward maximum uncertainty (i.e., disconfirm the current belief) about the jars. Figure 4 (top left) illustrates the effect of φ on $s(\mu_2^{(k)})$ over time.

In Models 5 and 6, changes in μ_2 from trial to trial occur according to two parameters: ω , the variance of the random process; and κ_1 , a scaling factor that changes the size of updates when $\hat{\mu}_1 = 0.5$, or maximum uncertainty, relative to when $\hat{\mu}_1$ is closer to 0 or 1 (i.e., when the subject is more confident about either jar). Figure 4 (bottom left) illustrates the effect of κ_1 on $\hat{\mu}_1$ over time. Formally, the scaling occurs as follows:

$$\hat{\mu}_1^{(k+1)} \equiv s(\mu_2^{(k)} \kappa_1)$$

When $\kappa_1 > 1$, updating toward 1 on observing a blue bead ($\mu = 1$) is greatest (i.e., switching between jars becomes more likely) when $\hat{\mu}_1 < 0.3$; when $\kappa_1 < 1$, updating is comparatively far lower when $\hat{\mu}_1 < 0.3$. This is illustrated in Figure 4 (middle): for high values of κ_1 (brown line), belief updates that cross the $\hat{\mu}_1 = 0.5$ line encounter little resistance (i.e., little evidence is required to cause a large shift), whereas approaching the extremes of $\hat{\mu}_1 = 0$ and $\hat{\mu}_1 = 1$ in response to confirmatory evidence is resisted (belief shifts are very small for $\hat{\mu}_1$ near 1). By contrast, for low values of κ_1 (Fig. 4 middle, black line), there is relatively less resistance against approaching the extremes while it takes more evidence for beliefs to cross the $\hat{\mu}_1 = 0.5$ line.

Figure 4 (right) illustrates the average absolute shifts in beliefs on observing beads of either color. This ‘vulnerability to updating’ is highly reminiscent of the ‘energy state’ of a neural network model (i.e., in low-

energy states) less updating occurs. The effect of increasing κ_1 is to convert confident beliefs about the jar (near 0 and 1) from low to high ‘energy states’ (i.e., to make them much more unstable). This recapitulates the attractor network properties illustrated in Figure 1: an unstable network easily switches from one state to another but has difficulty stabilizing any one state, whereas a stable network requires more energy (here, information) to overcome the boundary between two states (here, beliefs). Models 5 and 6 therefore capture the effects of attractor (in)stability on belief updating, or at least the kind of updating for which (un)stable attractor states are a good analogy.

As group differences in initial updating had been observed in Dataset 1, we also estimated the SD of μ_2 before the sequence begins, $\sigma_2^{(0)}$, in Models 2, 4, and 6.

NB for intermediate values of κ_1 , Models 5 and 6 produce similar belief updating trajectories to Models 3 and 4 (containing the disconfirmatory updating parameter φ): both make greater updates following disconfirmatory evidence. For more extreme values of κ_1 , however, Models 5 and 6 produce trajectories that Models 3 and 4 cannot: φ cannot pull beliefs far toward certainty in the opposite jar (compare Fig. 4, bottom left, brown line), and neither can it make it more difficult to update to disconfirmatory evidence (compare Fig. 4, bottom left, black line).

The parameters ω and $\nu \pm \sigma_2^{(0)} \pm \varphi$ or κ_1 were estimated individually for each subject. If estimated, the prior probability distributions for their values are given in Table 3. The means given here refer to the parameters’ native space, but the variances refer not to the parameters’ native space, which in many cases is bounded, but to the unbounded space they were

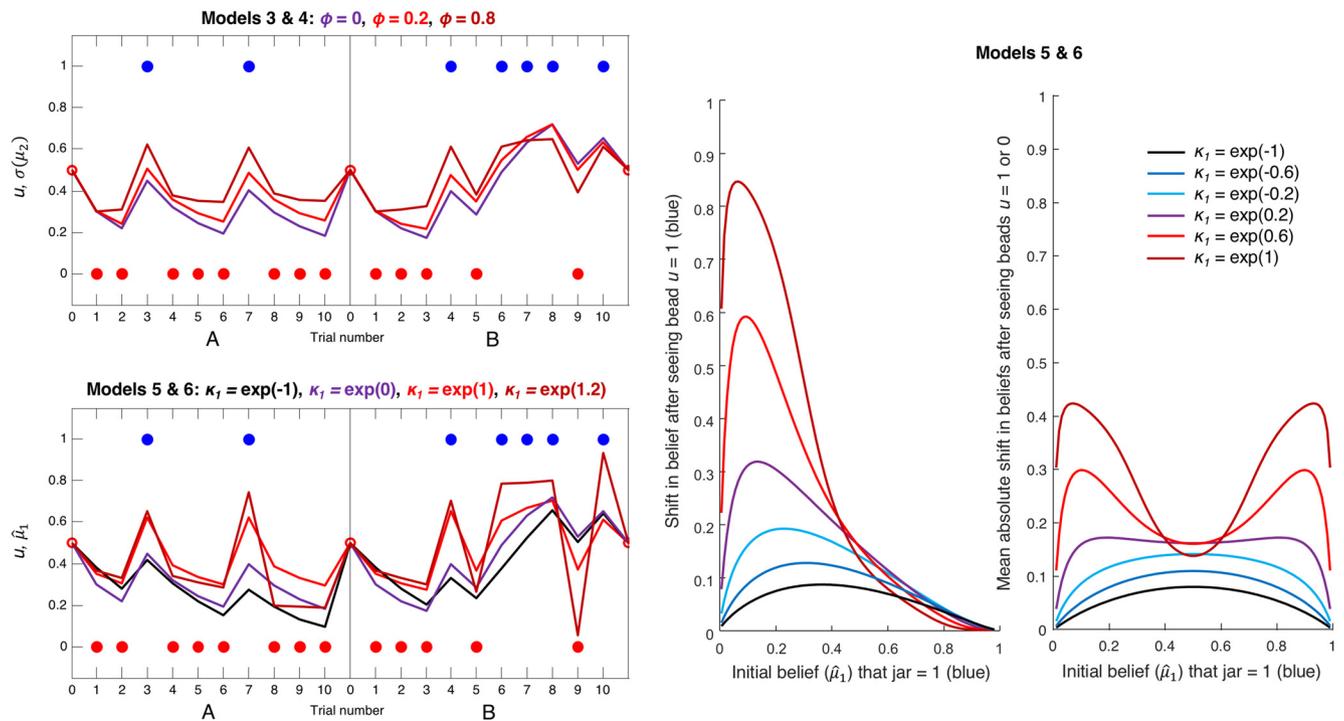


Figure 4. Simulated data illustrating the effects of φ (Models 3 and 4) and κ_1 (Models 5 and 6) on inference. Both panels represent simulated perceptual model predictions in the same format as before, with $\sigma_2^{(0)}$ and ω set to their previous values; hence, the purple line in these plots is identical to that in Figure 3. The second level and simulated responses y have been omitted for clarity. Top left, Simulations of a perceptual model incorporating an autoregressive order (1) process at the second level, using three different values of AR(1) parameter φ : 0, 0.2, and 0.8. The estimate of the tendency on trial $k + 1$, $\mu_2^{(k+1)}$, is selected from a Gaussian distribution with mean $\mu_2^{(k)} + \varphi(m - \mu_2^{(k)})$ and variance $\sigma_2^{(k)} + \exp(\omega)$. Over time, μ_2 is therefore attracted toward level m (fixed to 0, i.e., at $\sigma(\mu_2) = 0.5$) at a rate determined by φ . In effect, this gives the model a ‘disconfirmatory bias,’ such that as φ increases, $\sigma(\mu_2)$ is pulled further away from a belief in either jar, and toward 0.5 (maximum uncertainty about the jars). Bottom left, Simulations of a perceptual model using four different values of scaling factor κ_1 , which alters the sigmoid transformation: $\hat{\mu}_1^{(k+1)} = s(\kappa_1 \cdot \mu_2^{(k)})$. When $\kappa_1 > \exp(0)$, updating is greater to unexpected evidence and lower to consistent evidence; when $\kappa_1 < \exp(0)$, the reverse is true. Red and brown lines ($\kappa_1 > \exp(0)$) indicate the effects of increasingly unstable attractor networks; that is, switching between states (jars) becomes more likely (a concomitant increase in vulnerability to noise, i.e., response stochasticity, is not shown). Green line ($\kappa_1 = \exp(-1)$) indicates slower updating around $\hat{\mu}_1 = 0.5$, as was found in controls. κ_1 permits a greater range of updating patterns than φ (the green and brown trajectories in the bottom cannot be produced by Model 4), which may be why Model 6 can fit both controls and Scz groups well. Middle, Plot represents the effects of κ_1 on belief updating, as a function of the initial belief $\hat{\mu}_1$ ($\sigma_2^{(0)}$ and ω were set to 1.5 and -1 , respectively, as in Fig. 5; changing these parameters does not qualitatively alter the effects of κ_1 shown here). For values of $\kappa_1 < \exp(0) = 1$ (bottom three curves) and initial beliefs to the left of these curves’ maxima (i.e., that the jar is probably red), relatively small increases in $\hat{\mu}_1$ are made if one blue bead ($u = 1$) is observed, such that the subject still believes the jar is most likely red. For values of $\kappa_1 > \exp(0.5)$ (top two curves), observing one blue bead causes such a large update for all but the most certain initial beliefs in a red jar that the subject’s posterior belief is that the jar is probably blue. These subjects’ beliefs are no longer stable, but neither can they reach certainty: only tiny updates toward 1 are possible for $\hat{\mu}_1 > 0.8$. Right, Plot represents the average absolute shifts in beliefs on observing beads of either color. This ‘vulnerability to updating’ is highly reminiscent of the ‘energy state’ of a neural network model (schematically illustrated in Fig. 1) (i.e., in low energy states); less updating is expected. The effect of increasing κ_1 is to convert confident beliefs about the jar (near 0 and 1) from low to high ‘energy states’ (i.e., to make them much more unstable).

transformed to for estimation purposes. Otherwise, they were fixed as $\varphi = 0$ (Models 1 and 2) and $\sigma_2^{(0)} = 0.006$ (Models 1, 3, and 5). The model’s prior beliefs about the jars at the start of the sequence were fixed at $\mu_2^{(0)} = 0$ (i.e., believing each to be equally likely). The priors were sufficiently uninformative to be easily updated by the data: all prior means are standard for the HGF, except $\sigma_2^{(0)}$, which had to be increased from 0.006 to 0.8 to allow the data to change it. The latter change ensured that group differences in initial belief updating alone would cause group differences in $\sigma_2^{(0)}$ rather than κ_1 .

Model fitting and statistical analysis. We tested models with different combinations of parameters ω , ν , φ , or κ_1 and $\sigma_2^{(0)}$ (Table 3). In analyzing Dataset 2, we concatenated all four sequences for each subject to estimate the model parameters as accurately as possible (resetting the beliefs about the jars at the start of each sequence).

After fitting the six models to each subject’s data, we performed Bayesian model selection on all groups separately in both Dataset 1 (at baseline and follow-up) and Dataset 2. This procedure weights models according to their accuracy but penalizes them for complexity (i.e., unnecessary extra parameters) to prevent overfitting (Stephan et al., 2009; Rigoux et al., 2014). The winning model in all eight groups was Model 6 (see Fig. 6), although approximately one-third of psychotic subjects and nonclinical controls in Dataset 1 (at baseline) and in Dataset 2 were better fit by Model 4. It is unclear why this change occurs; but given that Model 6 can

produce very similar trajectories to Model 4 for intermediate values of κ_1 (Fig. 4), any increase in response stochasticity is likely to diminish the strength of evidence for one model over a similar one.

To confirm we could reliably estimate the parameters of the winning model, Model 6, we simulated 100 datasets using the modal values of the parameters for both control and Scz groups (Fig. 5, top and bottom rows, respectively; an example simulated dataset is shown in Fig. 3). We then estimated the parameters for the simulated data and showed that, in most cases, the parameters are recovered reasonably accurately. The exception was $\sigma_2^{(0)}$ in the Scz group simulation, which was distributed around the prior mean of 0.8 rather than the true value of 1.5. We retained a prior mean of 0.8 for $\sigma_2^{(0)}$ because using a higher prior mean led to overestimation of $\sigma_2^{(0)}$ in other simulations (data not shown).

Results

Behavioral results: Dataset 1

Each group’s mean responses are plotted in Figure 2A, and statistical tests detailed in Tables 1 and 2 (p (adj) refer to the adjusted p value of Tukey’s HSD *post hoc* test). As described previously (Peters and Garety, 2006), at baseline there was a significant difference in disconfirmatory updating between the groups

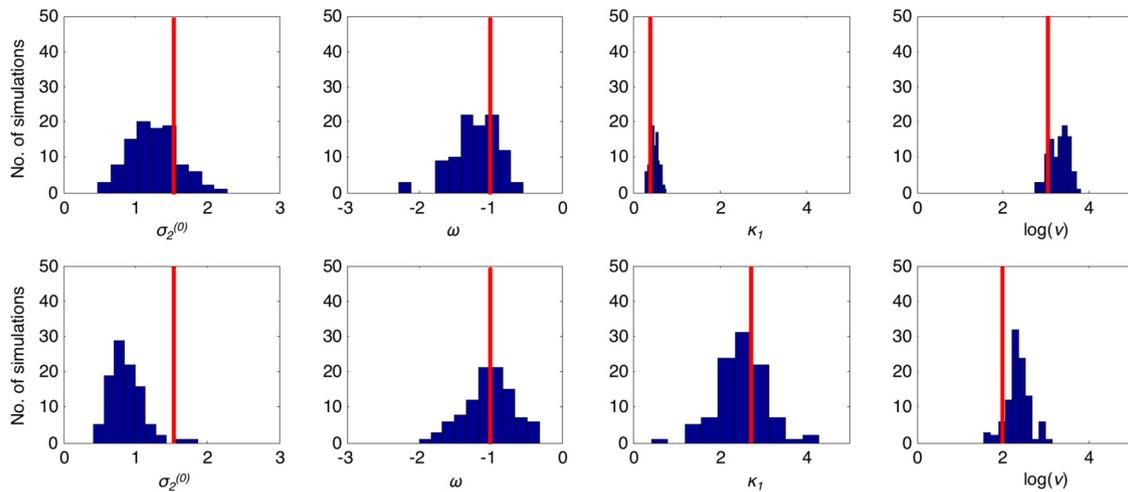


Figure 5. Recovery of model parameters from simulated data. The 200 datasets were simulated using Model 6: 100 using modal parameter values for the control group (Dataset 2) and 100 using modal values for the Scz group (also Dataset 2). Red lines indicate the values. Both used settings of $\sigma_2^{(0)} = 1.5$, $\omega = -1$. The control group used $\kappa_1 = 0.37$ (i.e., $\exp(-1)$) and $\nu = \exp(3)$. The Scz group used $\kappa_1 = 2.7$ (i.e., $\exp(1)$) and $\nu = \exp(2)$. Histograms represent the parameter estimates from model inversion using the same priors as were used in the main analysis shown above: the modal control and Scz simulation results are in the top and bottom rows, respectively.

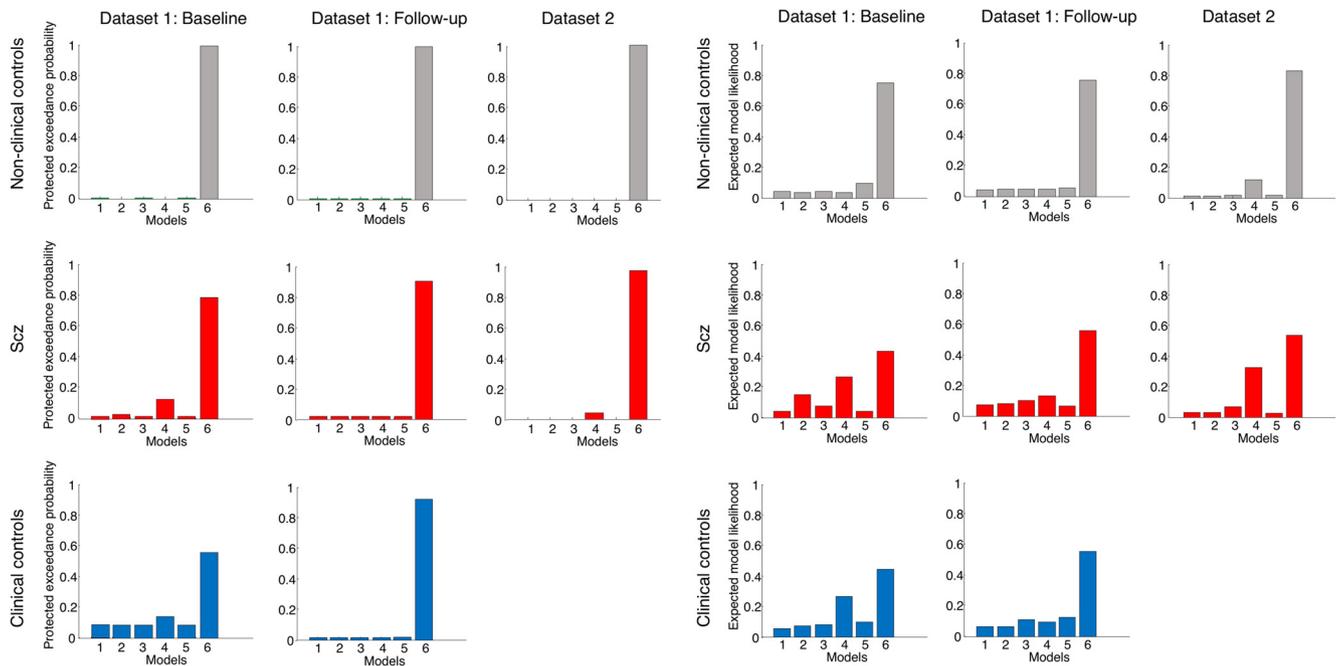


Figure 6. Bayesian model selection results for both datasets. Left, Protected exceedance probabilities for the six models in each group in each dataset. The protected exceedance probability is the probability a particular model is more likely than any other tested model, above and beyond chance, given the group data (Rigoux et al., 2014). Model 6 wins in all groups in both datasets (top row, controls; middle row, Scz; bottom row, clinical controls). Right, Model likelihoods for the six models in each group in each dataset. The model likelihood is the probability of that model being the best for any randomly selected subject (Stephan et al., 2009). Model 4 is a clear runner-up in the psychotic (Scz) and clinical control groups at baseline in Dataset 1, and in the Scz group in Dataset 2.

($F_{(2,77)} = 6$, $p = 0.004$, ANOVA), and the psychotic group had greater disconfirmatory updating than the nonclinical controls ($p(\text{adj}) = 0.003$) but not the clinical controls ($p(\text{adj}) = 0.4$). There was no difference between the clinical and nonclinical controls ($p(\text{adj}) = 0.13$). There were also significant differences in initial certainty across the three groups ($F_{(2,77)} = 8.7$, $p = 0.0004$, ANOVA); the psychotic group’s initial certainty was higher than the nonclinical controls’ ($p(\text{adj}) = 0.0003$) but not the clinical controls’ ($p(\text{adj}) = 0.25$). There was not a significant difference between the clinical and nonclinical control groups ($p(\text{adj}) = 0.06$). There were no group differences in final certainty ($F_{(2,77)} = 0.7$, $p = 0.5$, ANOVA).

At follow-up, the difference in disconfirmatory updating between the groups was no longer significant ($F_{(2,52)} = 2.9$, $p = 0.06$, ANOVA); the psychotic group had greater disconfirmatory updating than the nonclinical controls ($p(\text{adj}) = 0.049$) but not the clinical controls ($p(\text{adj}) = 0.4$). There was no significant difference in initial certainty across the groups ($F_{(2,52)} = 0.9$, $p = 0.4$, ANOVA). Differences in final certainty were no longer significant ($F_{(2,52)} = 2.8$, $p = 0.07$, ANOVA); the biggest difference was the nonclinical controls’ final certainty, which was numerically higher than the clinical controls’ ($p(\text{adj}) = 0.057$).

There were negative correlations between initial certainty and disconfirmatory updating at both baseline ($\rho = -0.41$, $p =$

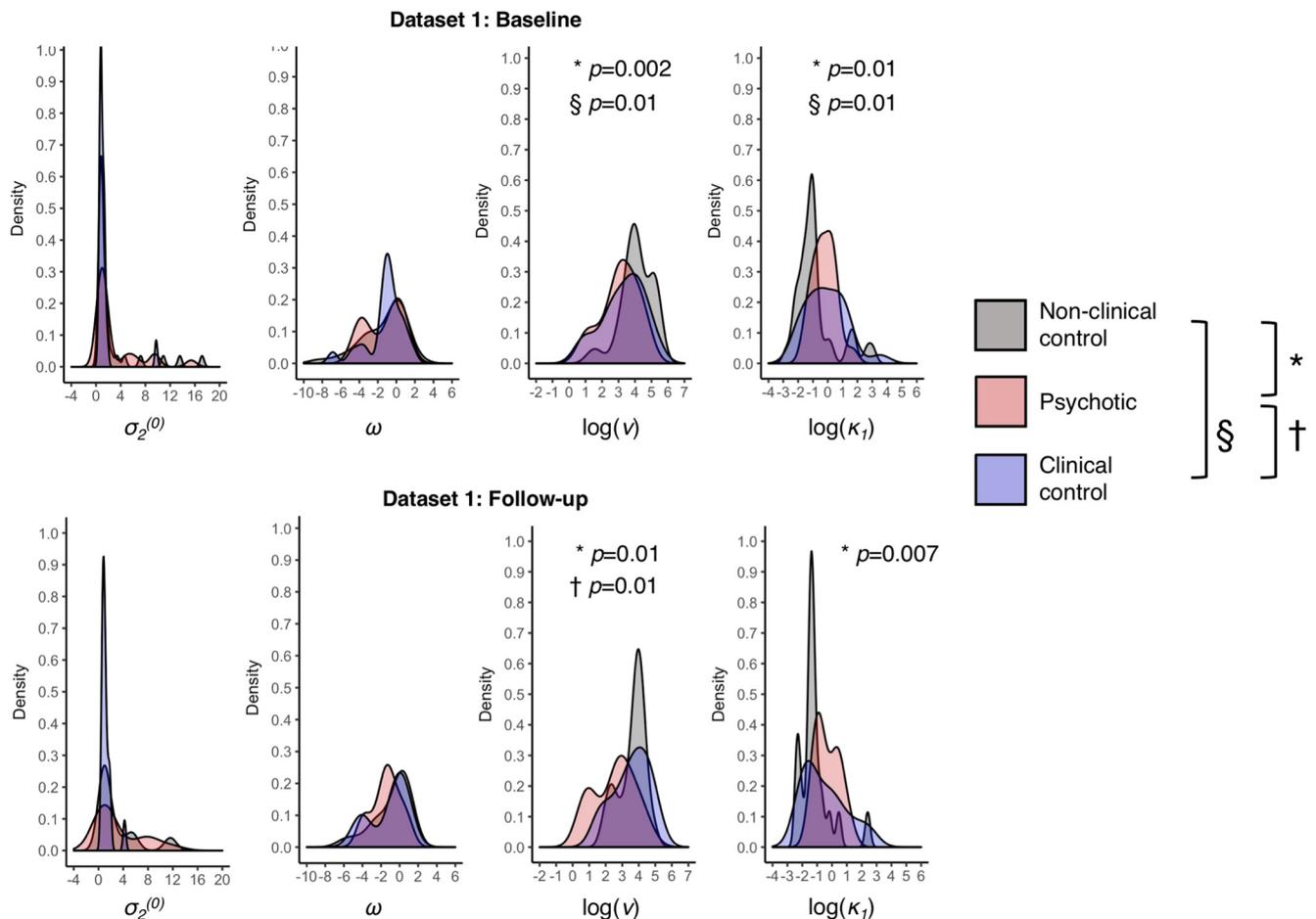


Figure 7. Probability density plots for Model 6 parameters in Dataset 1. The distributions of parameter values for $\sigma_2^{(0)}$, ω , $\log(\nu)$, and $\log(\kappa_1)$ are plotted for Dataset 1 at baseline (top row) and Dataset 1 at follow-up (bottom row). Symbols represent significant group differences: §between nonclinical controls and clinical controls; *between nonclinical controls and Scz; †between Scz and clinical controls.

0.00015) and follow-up ($p = -0.41, p = 0.002$), but not between final certainty and the other two measures ($p > 0.1$ in all four comparisons).

Behavioral results: Dataset 2

The mean responses of subjects in each group are plotted in Figure 2B. There was a significant increase in disconfirmatory updating in Scz compared with controls ($t_{(88.6)} = 2.1, p = 0.04$, Welch’s t test). There was mixed evidence for a difference in initial certainty between Scz and controls: Scz were more certain after the first bead in Sequences A and B but not Sequences C or D (Fig. 2; Table 3), but the difference in mean initial certainty fell short of statistical significance ($t_{(110)} = -1.9, p = 0.059$, Cohen’s $d = 0.32$, Welch’s t test). Final certainty was only assessed in Sequences A and D (B and C contained two changes of color in the last three beads): in both sequences, Scz were less certain than controls (Sequence A: $t_{(80.1)} = 3.0, p = 0.004$; Sequence D: $t_{(85.5)} = 3.4, p = 0.001$, Welch’s t tests).

Initial certainty and disconfirmatory updating negatively correlated within both Scz ($\rho = -0.46, p = 0.0003$) and control ($\rho = -0.57, p = 10^{-11}$) groups. Final certainty did not correlate with either measure in either group ($p > 0.4$ in four comparisons).

Modeling results: Dataset 1

Model selection results for the three groups analyzed separately at both baseline and follow-up are plotted in Figure 6 (columns 1, 2,

4, and 5); the probability of each model being best for any given subject is shown in the left panel, and the probability of each model being the best overall is shown in the right panel. Model 6 is the clear winner at each time point, although a minority of psychotic and clinical controls are best fit by Model 4.

Model 6’s parameter distributions are shown in Figure 7; they are skewed; hence, nonparametric tests were used to determine group differences (full details in Table 4; $p(adj)$ refers to the adjusted p value of Dunn’s *post hoc* test). At baseline there were large group differences in belief instability κ_1 ($\chi^2(2, n = 80) = 9.64, p = 0.008, \eta^2 = 0.12$, Kruskal–Wallis’ one-way ANOVA on ranks) and response stochasticity ν ($\chi^2(2, n = 80) = 11.9, p = 0.003, \eta^2 = 0.15$) but not in $\sigma_2^{(0)}$ or ω . There were statistically significant differences in κ_1 between the nonclinical controls and both the psychotic group ($p(adj) = 0.01$, Dunn’s test) and the clinical control group ($p(adj) = 0.01$), but not between the latter two groups ($p(adj) = 0.4$). Similarly, there were statistically significant differences in ν between the nonclinical controls and both the psychotic group ($p(adj) = 0.002$, Dunn’s test) and the clinical control group ($p(adj) = 0.01$), but not between the latter two groups ($p(adj) = 0.3$).

At follow-up, there were still large group differences in κ_1 ($\chi^2(2, n = 55) = 8.0, p = 0.02, \eta^2 = 0.15$, Kruskal–Wallis’ one-way ANOVA on ranks) and ν ($\chi^2(2, n = 55) = 8.5, p = 0.01, \eta^2 = 0.16$), but not in $\sigma_2^{(0)}$ or ω . There was a significant difference in κ_1 between the psychotic and nonclinical control groups

Table 4. Parameter distributions and statistical tests in Datasets 1 and 2

	$\sigma_2^{(0)}$	ω	$\log(\nu)$	$\log(\kappa_1)$
Dataset 1 (baseline, $n = 80$)				
Nonclinical controls: mean (SD)	2.5 (3.9)	-1.3 (2.4)	4.1 (1.0)	-0.8 (1.4)
Psychotic: mean (SD)	3.0 (3.9)	-1.4 (2.0)	3.1 (1.1)	-0.2 (0.8)
Clinical controls: mean (SD)	1.4 (1.9)	-1.2 (2.0)	3.3 (1.3)	-0.1 (1.4)
Kruskal–Wallis $\chi^2_{(2,80)}$	2.33, $p = 0.31$ $\eta^2 = 0.02$	0.22, $p = 0.9$ $\eta^2 = 0.0$	11.9, $p = 0.003$ $\eta^2 = 0.15$	9.6, $p = 0.008$ $\eta^2 = 0.12$
Post hoc Dunn tests				
Psychotic versus nonclinical controls	$p(\text{adj}) = 0.3$	$p(\text{adj}) = 1$	$p(\text{adj}) = 0.002$	$p(\text{adj}) = 0.01$
Clinical versus nonclinical controls	$p(\text{adj}) = 0.2$	$p(\text{adj}) = 0.7$	$p(\text{adj}) = 0.01$	$p(\text{adj}) = 0.01$
Psychotic versus clinical controls	$p(\text{adj}) = 0.2$	$p(\text{adj}) = 0.5$	$p(\text{adj}) = 0.3$	$p(\text{adj}) = 0.4$
Dataset 1 (follow-up, $n = 55$)				
Nonclinical controls: mean (SD)	2.8 (3.4)	-0.9 (2.0)	3.6 (0.8)	-1.2 (1.1)
Psychotic: mean (SD)	3.2 (3.7)	-1.4 (1.5)	2.5 (1.2)	-0.3 (0.8)
Clinical controls: mean (SD)	1.2 (0.9)	-1.1 (2.0)	3.5 (1.1)	-0.5 (1.4)
Kruskal–Wallis $\chi^2_{(2,80)}$	2.35, $p = 0.3$ $\eta^2 = 0.04$	2.32, $p = 0.3$ $\eta^2 = 0.04$	8.5, $p = 0.01$ $\eta^2 = 0.16$	8.0, $p = 0.02$ $\eta^2 = 0.15$
Post hoc Dunn tests				
Psychotic versus nonclinical controls	$p(\text{adj}) = 0.4$	$p(\text{adj}) = 0.2$	$p(\text{adj}) = 0.01$	$p(\text{adj}) = 0.007$
Clinical versus nonclinical controls	$p(\text{adj}) = 0.2$	$p(\text{adj}) = 0.3$	$p(\text{adj}) = 0.5$	$p(\text{adj}) = 0.1$
Psychotic versus clinical controls	$p(\text{adj}) = 0.3$	$p(\text{adj}) = 0.3$	$p(\text{adj}) = 0.01$	$p(\text{adj}) = 0.1$
Dataset 2 ($n = 167$)				
Nonclinical controls: mean (SD)	3.1 (2.6)	-2.3 (2.0)	2.8 (1.0)	-0.8 (0.9)
Scz: mean (SD)	1.9(1.5)	-2.1 (1.8)	2.1 (1.2)	0.2 (1.0)
Mann–Whitney U test	$Z = 3.1, p = 0.002, r = 0.24$	$Z = -0.6, p = 0.6, r = 0.04$	$Z = 3.9, p = 0.0001, r = 0.3$	$Z = -5.6, p = 3 \times 10^{-8}, r = 0.43$
Dataset 2 (better-matched controls, $n = 116$)				
Nonclinical controls: mean (SD)	2.8 (2.7)	-2.2 (2.1)	2.9 (1.1)	-0.6 (1.0)
Scz: mean (SD)	1.9 (1.5)	-2.1 (1.8)	2.1 (1.2)	0.2 (1.0)
Mann–Whitney U test	$Z = 1.9, p = 0.056, r = 0.18$	$Z = 0.12, p = 0.9, r = 0.01$	$Z = 3.4, p = 0.0007, r = 0.31$	$Z = -4.1, p = 0.00004, r = 0.38$

($p(\text{adj}) = 0.007$, Dunn’s test) but not the clinical and nonclinical control groups ($p(\text{adj}) = 0.1$); ν remained significantly different between the nonclinical controls and both the psychotic group ($p(\text{adj}) = 0.01$, Dunn’s test) and now also between the psychotic and clinical control groups ($p(\text{adj}) = 0.01$), but not between the clinical and nonclinical controls ($p(\text{adj}) = 0.5$).

We explored whether group differences in κ_1 or ν at baseline and follow-up might be ascribable to IQ (Quick Test score) (Ammons and Ammons, 1962), as the groups’ IQ scores were not equivalent (Tables 1, 2). Including both IQ and group status within one regression model is an unsound method of testing for confounding by IQ because group and IQ are clearly not independent here (Miller and Chapman, 2001), so we tested for relationships between the parameters and IQ separately within each group at each time point. No relationships reached statistical significance (all $p > 0.1$), the closest being a trend between κ_1 and IQ in nonclinical controls only ($r = -0.30, p = 0.08$); nevertheless, given the smaller group sizes and larger between- versus within-group variances, it remains plausible that IQ differences contribute to group parameter differences.

We tested whether κ_1 or ν at baseline related to delusion-proneness (Peters Delusion Inventory score [PDI]; Peters et al., 1999) across all groups, after first excluding any interaction between PDI and group; PDI significantly correlated with ν ($F_{(1,67)} = 7.1, p = 0.01$, ANCOVA) but not κ_1 ($F_{(1,67)} = 3.2, p = 0.079$, ANCOVA). We did not analyse the Delusions-Symptoms-States Inventory (Foulds and Bedford, 1975) as it is a less specific measure of delusions. We tested whether κ_1 or ν at baseline was correlated with any particular subgroup of symptoms (measured using the Manchester Scale) (Krawiecka et al., 1977) in both clinical groups only, using the regression models $\kappa_1[\text{ or } \nu] \sim \text{const} + \nu_1 * \text{MSaffective} + \nu_2 * \text{MSpositive} + \nu_3 * \text{MSnegative}$: none of the models was significant, however (all $p > 0.1$).

At baseline, there was no evidence of a correlation between κ_1 and antipsychotic medication dose ($p = 0.3$), but the correlation between ν and medication dose approached significance ($\rho = -0.4, p = 0.067$).

We tested for correlations between the Model 6 parameters (Spearman’s ρ was used where distributions were not parametric): κ_1 and ν were negatively correlated both at baseline ($\rho = -0.38, p = 0.0004$) and at follow-up ($\rho = -0.52, p = 0.0001$), as were κ_1 and ω at baseline ($\rho = -0.47, p = 10^{-5}$) and follow-up ($\rho = -0.53, p = 10^{-5}$). In estimating the parameters from simulated data, the only correlation present in both simulations (indicating some consistent trading-off between these parameters during estimation) was between κ_1 and ω , with $r = -0.5$ in each case. This is not surprising, as both κ_1 and ω affect updating to new information throughout the sequence (unlike $\sigma_2^{(0)}$) in a deterministic way (unlike ν). Nevertheless, κ_1 was estimated very reliably in the first simulation (Fig. 5, top row) and with reasonable accuracy in the second (Fig. 5, bottom row), so we are confident that the group differences in κ_1 are genuine. The correlations of $\rho \approx -0.5$ between ω and κ_1 in Dataset 1 are unlikely to be reliable, however.

Modeling results: Dataset 2

We tested the same six models and performed Bayesian model selection as before. As in Dataset 1, the winning model was Model 6 overall and in each group separately (Fig. 6), although in the Scz group a minority were best captured by Model 4. Model 6’s parameter distributions are shown in Figure 8; they are skewed, so nonparametric tests were used (for full details, see Table 4).

As in Dataset 1, belief instability κ_1 was significantly higher in Scz than in controls ($Z = -5.6, p = 10^{-8}$, Mann–Whitney U test) with a medium-to-large effect size ($r = 0.43$); also response sto-

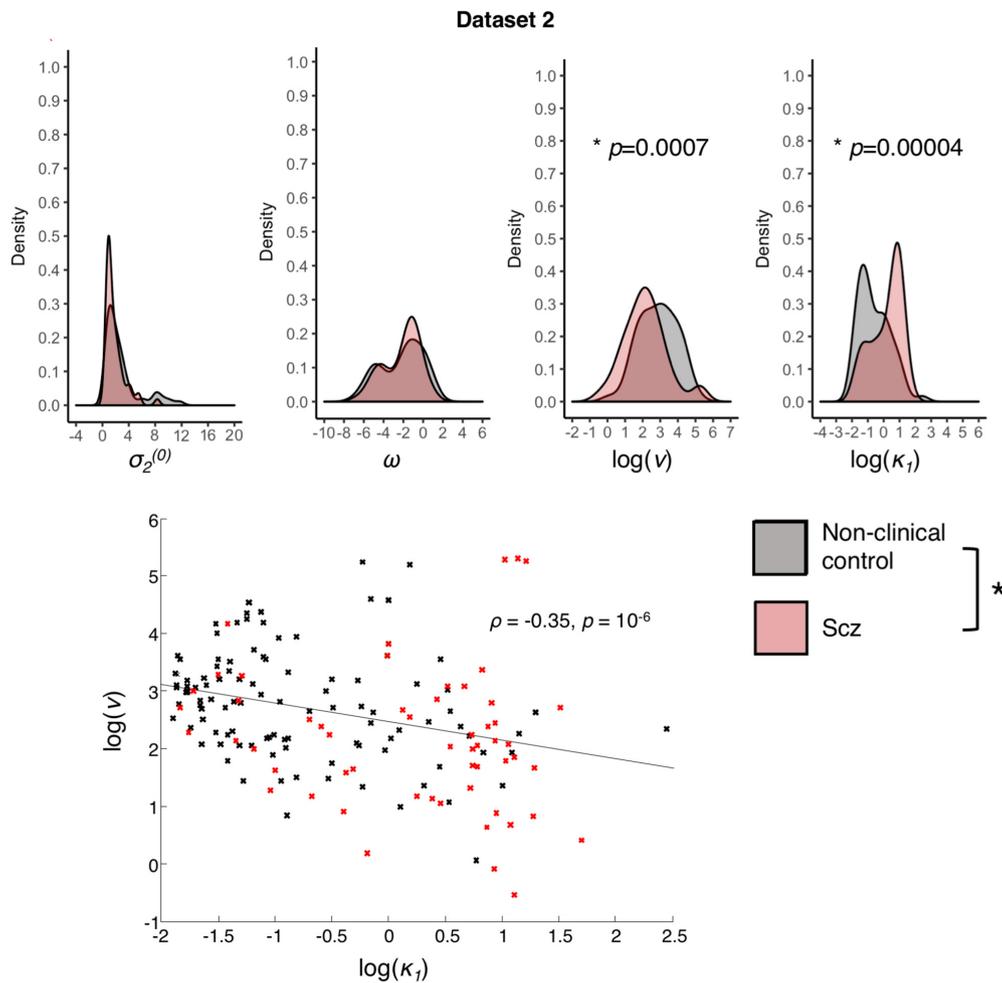


Figure 8. Model 6 parameters in Dataset 2: distributions and correlation. Top, The distributions of parameter values for $\sigma_2^{(0)}$, ω , $\log(\nu)$, and $\log(\kappa_1)$ are plotted for Dataset 2. *Significant group differences between the Scz group and nonclinical control subgroup (well matched in age and sex); the group difference in $\sigma_2^{(0)}$ is not indicated because it was nonsignificant ($p = 0.056$) in the well-matched comparison. Bottom, The significant correlation between $\log(\nu)$ and $\log(\kappa_1)$ in Dataset 2 is plotted, with controls’ parameters in black and Scz in red. Similar correlations were also found in Dataset 1 at both time points.

chasticity ν was lower in Scz than in controls ($Z = 3.9, p = 0.0001, r = 0.3$, Mann–Whitney U test), as was initial belief variance $\sigma_2^{(0)}$ ($Z = 3.1, p = 0.002, r = 0.24$, Mann–Whitney U test). There were no statistically significant group differences in evolution rate ω . See Figures 6 and 7 for examples of model fits in subjects with lower κ_1 values (two controls in Fig. 9) and higher κ_1 values (2 Scz subjects in Fig. 10); each figure also illustrates the effects of lower and higher ω values (in the top and bottom rows, respectively). We repeated the analysis using a subset of the controls ($n = 60$) that were better matched in age and sex, as the original control group was younger and more female than the patient group (Tables 1, 2). The group differences in κ_1 and ν were unchanged in this analysis ($Z = -4.1, p = 0.00004; Z = 3.4, p = 0.0007$, respectively, Mann–Whitney U tests), but that in $\sigma_2^{(0)}$ was no longer significant ($Z = 1.9, p = 0.056$, Mann–Whitney U test).

Although IQ (National Adult Reading Test score) (Nelson, 1982) was evenly matched in these groups, working memory (Letter Number Sequencing score) (Wechsler, 1997) was lower in Scz than in controls (Tables 1, 2). We explored whether the group parameter differences might be related to working memory, by testing for correlations between κ_1 or ν and working memory in each group separately (Miller and Chapman, 2001): none was statistically significant (all $p > 0.1$). We also tested for relationships between κ_1 or ν and IQ (National Adult Reading Test) in

each group: ν and IQ (National Adult Reading Test) were correlated in Scz ($r = 0.33, p = 0.014$), but no other relationships were significant (all $p > 0.1$).

We tested whether κ_1 or ν related to schizotypy (Schizotypal Personality Questionnaire score; Raine, 1991) across all groups, but neither did so (both $p = 0.4$, ANCOVA). We tested whether κ_1 or ν was predicted by any particular subgroup of symptoms (measured using the Positive and Negative Symptom Scale) (Kay et al., 1987) in the Scz group only, using the regression model κ_1 [or ν] \sim const + $\nu_1 * \text{PANSS}_{\text{general}} + \nu_2 * \text{PANSS}_{\text{positive}} + \nu_3 * \text{PANSS}_{\text{negative}}$: the κ_1 model was not significant ($F = 0.9, p = 0.4$), but ν was weakly predicted by negative symptoms (overall $F = 2.76, p = 0.051$; for $\nu_3, t = -2.1, p = 0.04$). We had no record of medication dose in Dataset 2.

We tested for correlations between the Model 6 parameters: as in Dataset 1, κ_1 and ν were negatively correlated (Fig. 8; $\rho = -0.35, p = 10^{-6}$), but unlike Dataset 1, the only other statistically significant correlation was between κ_1 and $\sigma_2^{(0)}$ ($\rho = -0.54, p = 10^{-13}$). There was a correlation of $r = -0.2$ between κ_1 and ν in the data simulated from modal Scz parameter values (Fig. 5, bottom row), but no correlation in the first. This implies that the consistent correlations between these parameters of $\rho = -0.38, \rho = -0.52$ (Dataset 1 baseline and follow-up) and $\rho = -0.35$ (Dataset 2) are unlikely to be just estimation artifacts. The only

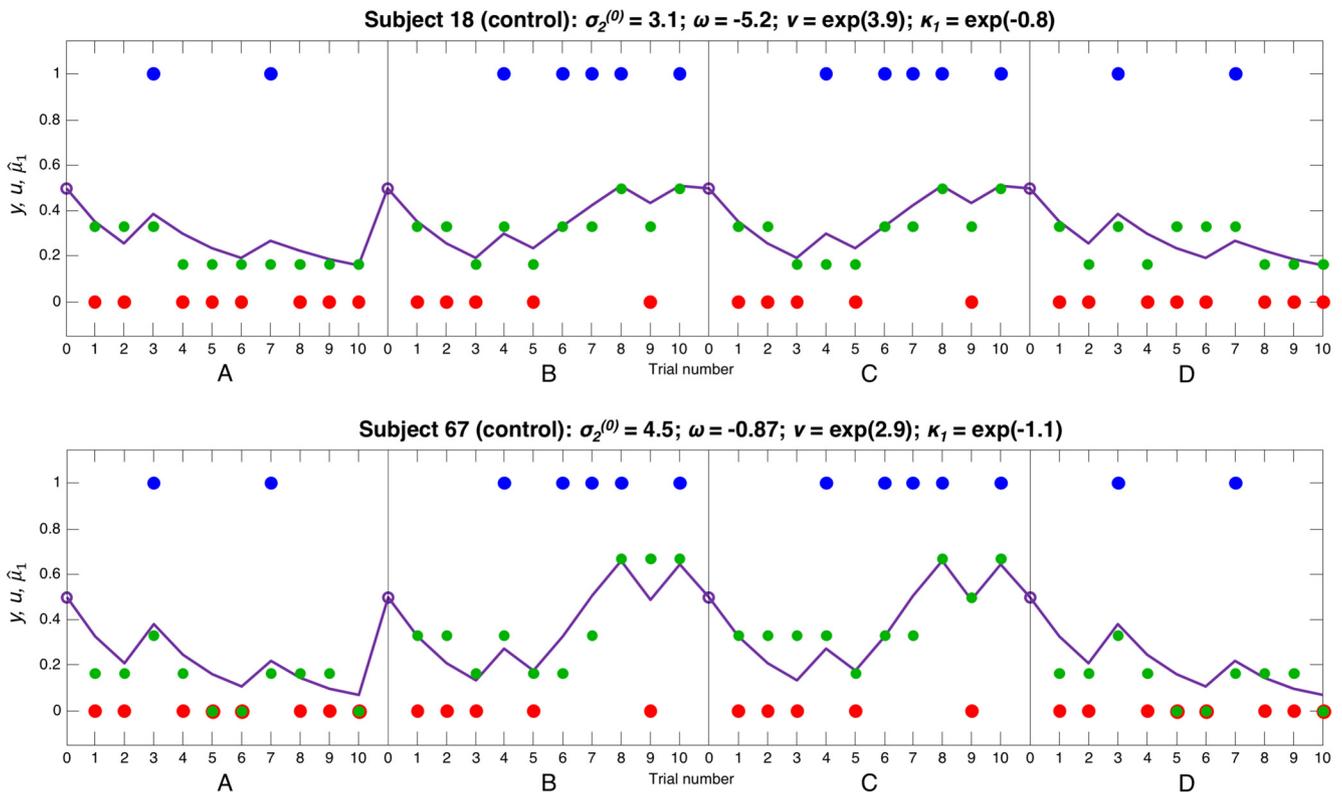


Figure 9. Responses and model fits for 2 control subjects. These plots show 2 control subjects’ responses to four 10-bead sequences concatenated together, in the same format as Figure 3 (but without the second level, due to space constraints); in the latter two sequences, blue and red were swapped around for model-fitting purposes. Each plot shows $u^{(k)}$, the beads seen by the subjects on trials $k = 1, \dots, 10$ (blue and red dots), y , the subject’s (Likert scale) response about the probability the jar is blue (green dots), and $\hat{\mu}_1^{(k+1)}$, the model’s estimate of the subject’s prediction the jar is blue (purple line). The parameter estimates for each subject are shown above their graphs. These subjects have fairly similar initial variance $\sigma_2^{(0)}$, (inverse) response stochasticity ν , and instability factor κ_1 . Subject 18 in the top has a much lower overall evolution rate ω than Subject 67 in the bottom; therefore, Subject 18 never reaches certainty about either jar, and makes relatively small changes to her beliefs in response to beads of varying colors. Both subjects have a low κ_1 , and so they make relatively small adjustments to their beliefs following unexpected evidence (this behavior can best be captured by the models containing κ_1 ; see Fig. 4). Subject 18’s responses are very close to those predicted by the model, and this is reflected in her relatively high value of ν .

other correlation between parameters in the simulated data was between $\sigma_2^{(0)}$ and κ_1 , of $r = -0.25$, in the first simulation only. These parameters were correlated in Dataset 2 but not Dataset 1.

Discussion

Scz tend to update their beliefs more to unexpected information and less to consistent information, compared with controls. We have replicated these behavioral effects, and demonstrated a computational basis for them that is informed by the unstable attractor hypothesis of Scz. In computational models of two ‘beads task’ datasets, Scz had consistently greater belief instability (κ_1) and response stochasticity (ν) than controls, as the unstable attractor hypothesis predicts. Furthermore, ν correlated with κ_1 in all three experiments, supporting the idea that ν is measuring a stochasticity that is related to κ_1 by an underlying neurobiological process, rather than simply an unmodeled effect.

These findings are important because they connect numerous reasoning biases previously found in Scz (e.g., a disconfirmatory bias) (Garety et al., 1991; Fear and Healy, 1997; Young and Bental, 1997; Peters and Garety, 2006), increased initial certainty (Peters and Garety, 2006), and decreased final certainty (Baker et al., 2018), and its associated stochasticity in responding (Moutoussis et al., 2011; Schlagenhaut et al., 2014) to model parameters that describe how belief updating in cortex could be perturbed by unstable attractor states due to NMDA (or dopamine 1) receptor hypofunction (Fig. 1).

The unique features of Model 6 that make attractor dynamics a compelling neurobiological explanation for its dominance are both Scz and controls’ nonlinearities in belief updating to confirmatory versus disconfirmatory evidence. The Scz group updated its beliefs (sometimes much) more to disconfirmatory than confirmatory evidence, particularly at points of relative certainty about the jar, and the controls were the opposite. Models with uniformly high or low learning rates cannot reproduce these effects; and adding high- or low-level (sensory) uncertainty to a hierarchical model would lead to uniformly high or low learning rates, respectively. Although Models 3 and 4 do show differential updating to confirmatory versus disconfirmatory evidence, this results in beliefs in either jar hovering at ~ 0.5 (as in Fig. 4, top left) rather than making large updates from belief in one jar to the other (as when $\kappa_1 = \exp(1.2)$: Fig. 4, bottom left). Furthermore, degraded neuronal ensemble firing (consistent with unstable attractor states) has recently been shown to be common to two different mouse models of Scz (Hamm et al., 2017).

In Dataset 1, belief instability κ_1 and response stochasticity ν were also significantly different between the clinical (mood disorder) and nonclinical control groups when the former were unwell, but not at follow-up, whereas the differences between the psychotic group and nonclinical controls persisted. This indicates that the same computational parameters can be perturbed in either a trait- or state-like manner, perhaps by different mech-

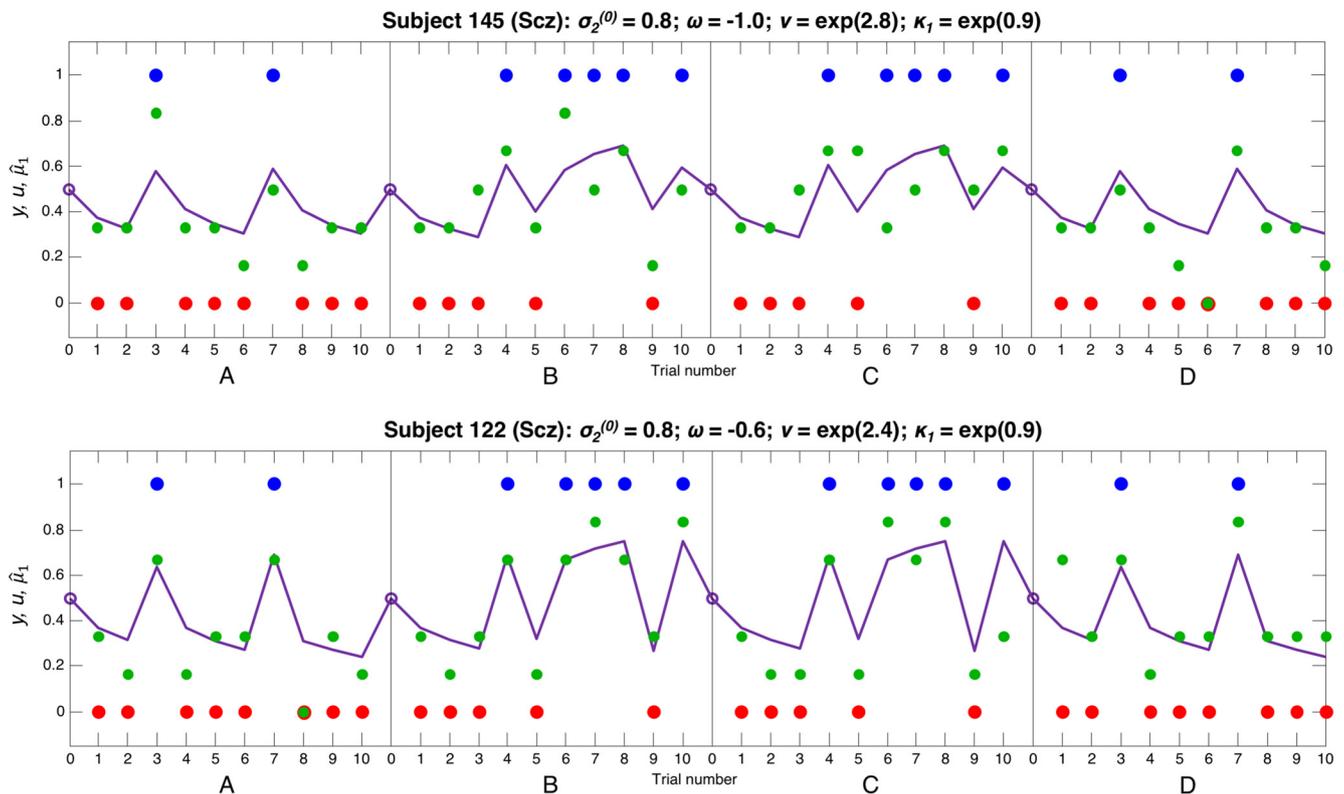


Figure 10. Responses and model fits for 2 Scz subjects. These plots show 2 Scz subjects’ responses to four 10-bead sequences in the same format as Figure 9. These subjects have similar evolution rate ω to the control subjects in Figure 9, but they both have a much higher κ_1 , meaning that they make much greater changes to their beliefs when presented with unexpected evidence, but do not reach certainty when faced with consistent evidence. Subject 122 (bottom) has a slightly higher evolution rate ω than Subject 145 (top), and so his switching between jars is even more pronounced. These subjects also have slightly lower (inverse) response stochasticity ν than the control subjects in Figure 9, and so their responses tend to be further from the model predictions.

anisms. It seems unlikely that these parameter changes simply reflect a lack of engagement with the task in clinical groups (especially when unwell) because the consistent changes in κ_1 , with which the changes in ν consistently correlate, reflect specific patterns of belief updating.

Parameter relationships with cognition and symptoms

Neither κ_1 nor ν showed significant relationships with IQ (in Dataset 1) or working memory (in Dataset 2) within the groups, giving some indication that the group differences in these cognitive measures were unlikely to be the main drivers of group differences in the parameters. Nevertheless, aside from the correlation between response stochasticity ν and IQ in Dataset 2, it is perhaps surprising that there were not more relationships between κ_1 or ν and cognitive measures in Scz, given it is likely that abnormal prefrontal dynamics have profound effects on all these variables. We may have lacked power to detect them, although Dataset 2 had 80% power to detect a correlation of 0.33, or perhaps different prefrontal regions contribute to working memory, IQ, and belief updating.

One might also question why there were no strong relationships between κ_1 or ν and positive or negative symptom domains (negative symptoms were weakly associated with ν in Dataset 2 only). Again, power may have been an issue, although across all subjects in Dataset 1, response stochasticity ν was associated with PDI score, even after including group in the model, indicating a potential relationship with delusions, but not with the broader concept of schizotypy (assessed in Dataset 2). It is also likely that other pathological factors contribute to symptoms, beyond those measured here (e.g., striatal dopamine availability and positive

symptoms). Of note, two other computational studies demonstrating clear working memory parameter differences between Scz and controls also failed to detect any relationship between those parameters and symptom domains (Collins et al., 2014, 2017). Both their and our Scz groups were taking antipsychotic medication, which is also likely to weaken correlations of parameters to positive symptoms.

Although replicated numerous times in the beads task, a ‘disconfirmatory bias’ is perhaps surprising in Scz, given one might expect delusional subjects to show a bias against disconfirmatory evidence (as indeed they do in tasks involving scenario interpretation) (Woodward et al., 2006). Indeed, the disconfirmatory bias is misleadingly named, as Scz make large shifts in beliefs both away from and back toward the current hypothesis (there are numerous examples in both datasets in Fig. 2). This pronounced switching behavior in the beads task is likely to illustrate a more fundamental instability of cognition and prefrontal dynamics in Scz, rather than being related to delusions specifically; indeed, the latter may be an attempt to remedy the former.

It is interesting that nonclinical controls’ data were also best fit by Model 6 in both datasets, implying that even healthy subjects show some asymmetry in their belief updating to expected versus unexpected evidence. Most nonclinical control subjects had $\kappa_1 < 1$ (i.e., reduced updating to changing evidence).

Related modeling studies

How do these findings relate to other computational modeling work in Scz? A study of unmedicated, mainly first-episode Scz performing a reversal learning task (Schlagenhauf et al., 2014) also demonstrated an increased tendency to switch that was not accounted for by re-

ward sensitivity (which would be affected by more stochastic behavior), and increased switching also occurs in chronic Scz (Waltz et al., 2013), although not always (Pantelis et al., 1999).

Two recent studies of similar tasks in Scz populations have also demonstrated evidence of nonlinear belief updating. Jardri et al. (2017) showed that the Scz group on average “overcount” the likelihood in a single belief update, an effect they attribute to reverberating cortical message-passing but could also be due to the belief instability shown by Model 6. Stuke et al. (2017) showed, in a very similar task, that all subjects showed evidence of nonlinear updating, but the Scz group updated more than controls to “irrelevant information” (i.e., disconfirmatory evidence). Some differences between their model and ours are that they did not estimate response stochasticity in their subjects (neither did Jardri et al., 2017), and their ‘nonlinearity’ parameter was bounded by linear updating on one side, approximately equivalent to belief instability κ_1 being constrained to being <1 in our model, whereas we have shown (as in Jardri et al., 2017) that Scz belief updating is often beyond this bound (Fig. 7) and more stochastic. Conversely, Moutoussis et al. (2011) demonstrated increased response stochasticity in acutely psychotic subjects but did not test for differences in belief updating.

The extent to which a loss of belief stability in Scz is apparent depends critically on the strength (precision) of incoming sensory evidence relative to the current belief (prior): if the former is less precise, no belief switching may occur, and instead the percept may be weighted toward the prior. In the beads task, sensory evidence (i.e., the color of the bead drawn) is unambiguous, but a task using very imprecise auditory sensory evidence (Powers et al., 2017) demonstrated some interesting heterogeneity in Scz: nonhallucinating Scz showed greater belief updating relative to controls, whereas in hallucinating Scz, percepts were driven by prior expectations, leading to a reduction in the updating of their beliefs (relative to controls).

Further evidence for heterogeneity in Scz is that those with delusions have greater certainty about the hypothesis that matches the evidence at every stage (Speechley et al., 2010), unlike the reduced final certainty we observed in Scz in Dataset 2. On the other hand, Scz with high negative symptoms have difficulty choosing the most rewarding option very consistently (Gold et al., 2012), which may reflect a lack of certainty about its value. We lacked sufficient power to detect differences between Scz with exclusively high positive or negative symptoms, however.

Limitations

Each of our datasets contains some limitations of the beads task that are addressed by the other. Dataset 1 did not include a memory aid or measure working memory, but Dataset 2 did both, and Dataset 2 also matched IQ across groups much better than Dataset 1; Dataset 2 used a Likert scale for responding and so could potentially exaggerate small changes in belief updating, but Dataset 1 used a continuous measure; Dataset 2 only tested stable outpatients, but Dataset 1 tested more unwell inpatients and retested them once they were better. The main limitation common to both datasets is that all subjects with psychotic diagnoses were taking antipsychotic medication when tested. Although the correlation between ν and medication dose was almost significant in Dataset 1, this relationship seems likely to be driven by illness severity rather than medication itself. Dopamine 2 receptor antagonists seem to both reduce overconfidence in probabilistic reasoning (Andreou et al., 2014) and also reduce motor response variability (Galea et al., 2013) and so, if anything, likely reduce our group differences.

In conclusion, we have shown that Scz subjects in two independent beads task datasets have consistent differences in two parameters of a belief updating model that attempts to reproduce consequences of attractor network instability. This study was designed to link patterns of inferences to model parameters that (do or do not) mimic the effects of abnormal attractor states on belief updating. The HGF itself does not contain attractor states, and no relation between its parameters and NMDAR function has hitherto been tested. More detailed spiking network modeling, pharmacological (or other NMDAR) manipulations, and imaging are required in the future to understand how neuromodulatory function in both pyramidal cells and inhibitory interneurons contributes to real attractor dynamics and probabilistic inference, and to seek empirical evidence for a correspondence between the stability of network states and the stability of its inferences (especially in Scz). This work underscores the importance of relating psychological biases to their underlying computational mechanisms, and thence (in future) to the constraints (e.g., the hypofunction of NMDARs) that neurobiology imposes on these mechanisms.

References

- Abi-Saab WM, D’Souza DC, Moghaddam B, Krystal JH (1998) The NMDA antagonist model for schizophrenia: promise and pitfalls. *Pharmacopsychiatry* 31 [Suppl 2]:104–109. [CrossRef Medline](#)
- Adams RA, Huys QJ, Roiser JP (2016) Computational psychiatry: towards a mathematically informed understanding of mental illness. *J Neurol Neurosurg Psychiatry* 87:53–63. [CrossRef Medline](#)
- Ammons RB, Ammons CH (1962) The quick test (QT): provisional manual. *Psychol Rep* 11:111–161.
- Andreou C, Moritz S, Veith K, Veckenstedt R, Naber D (2014) Dopaminergic modulation of probabilistic reasoning and overconfidence in errors: a double-blind study. *Schizophr Bull* 40:558–565. [CrossRef Medline](#)
- Averbeck BB, Evans S, Chouhan V, Bristow E, Shergill SS (2011) Probabilistic learning and inference in schizophrenia. *Schizophr Res* 127:115–122. [CrossRef Medline](#)
- Baker S, Konova A, Daw N, Horga G (2018) T216. Deficient Belief Updating Explains Abnormal Information Seeking Associated With Delusions in Schizophrenia. *Biol Psychiatry* 83:S212.
- Beal MJ (2003) Variational algorithms for approximate Bayesian inference. <https://cse.buffalo.edu/faculty/mbeal/papers/beal03.pdf>. Accessed March 26, 2012.
- Brunel N, Wang XJ (2001) Effects of neuromodulation in a cortical network model of object working memory dominated by recurrent inhibition. *J Comput Neurosci* 11:63–85. [CrossRef Medline](#)
- Collins AG, Brown JK, Gold JM, Waltz JA, Frank MJ (2014) Working memory contributions to reinforcement learning impairments in schizophrenia. *J Neurosci* 34:13747–13756. [CrossRef Medline](#)
- Collins AG, Albrecht MA, Waltz JA, Gold JM, Frank MJ (2017) Interactions among working memory, reinforcement learning, and effort in value-based choice: a new paradigm and selective deficits in schizophrenia. *Biol Psychiatry* 82:431–439. [CrossRef Medline](#)
- Diaconescu AO, Mathys C, Weber LA, Daunizeau J, Kasper L, Lomakina EI, Fehr E, Stephan KE (2014) Inferring on the intentions of others by hierarchical Bayesian learning. *PLoS Comput Biol* 10:e1003810. [CrossRef Medline](#)
- Dudley R, Taylor P, Wickham S, Hutton P (2016) Psychosis, delusions and the “jumping to conclusions” reasoning bias: a systematic review and meta-analysis. *Schizophr Bull* 42:652–665. [CrossRef Medline](#)
- Durstewitz D, Seamans JK (2008) The dual-state theory of prefrontal cortex dopamine function with relevance to catechol-o-methyltransferase genotypes and schizophrenia. *Biol Psychiatry* 64:739–749. [CrossRef Medline](#)
- Fear CF, Healy D (1997) Probabilistic reasoning in obsessive-compulsive and delusional disorders. *Psychol Med* 27:199–208. [CrossRef Medline](#)
- Foulds GA, Bedford A (1975) Hierarchy of classes of personal illness. *Psychol Med* 5:181–192. [CrossRef Medline](#)
- Friston K (2005) A theory of cortical responses. *Philos Trans R Soc Lond B Biol Sci* 360:815–836. [CrossRef Medline](#)
- Galea JM, Ruge D, Buijink A, Bestmann S, Rothwell JC (2013) Punishment-

- induced behavioral and neurophysiological variability reveals dopamine-dependent selection of kinematic movement parameters. *J Neurosci* 33:3981–3988. [CrossRef Medline](#)
- Garety PA, Hemsley DR, Wessely S (1991) Reasoning in deluded schizophrenic and paranoid patients: biases in performance on a probabilistic inference task. *J Nerv Ment Dis* 179:194–201. [CrossRef Medline](#)
- Gepperth A, Lefort M (2016) Learning to be attractive: probabilistic computation with dynamic attractor networks. In: 2016 Joint IEEE International Conference on Development and Learning and Epigenetic Robotics (ICDL-EpiRob), pp 270–277.
- Gold JM, Waltz JA, Matveeva TM, Kasanova Z, Strauss GP, Herbener ES, Collins AG, Frank MJ (2012) Negative symptoms and the failure to represent the expected reward value of actions: behavioral and computational modeling evidence. *Arch Gen Psychiatry* 69:129–138. [CrossRef Medline](#)
- Hamm JP, Peterka DS, Gogos JA, Yuste R (2017) Altered cortical ensembles in mouse models of schizophrenia. *Neuron* 94:153–167.e8. [CrossRef Medline](#)
- Hauser TU, Iannaccone R, Ball J, Mathys C, Brandeis D, Walitza S, Brem S (2014) Role of the medial prefrontal cortex in impaired decision making in juvenile attention-deficit/hyperactivity disorder. *JAMA Psychiatry* 71:1165–1173. [CrossRef Medline](#)
- Hopfield JJ (1982) Neural networks and physical systems with emergent collective computational abilities. *Proc Natl Acad Sci U S A* 79:2554–2558. [CrossRef Medline](#)
- Iglesias S, Mathys C, Brodersen KH, Kasper L, Piccirelli M, den Ouden HE, Stephan KE (2013) Hierarchical prediction errors in midbrain and basal forebrain during sensory learning. *Neuron* 80:519–530. [CrossRef Medline](#)
- Jardri R, Duverne S, Litvinova AS, Denève S (2017) Experimental evidence for circular inference in schizophrenia. *Nat Commun* 8:14218. [CrossRef Medline](#)
- Javitt DC, Zukin SR, Heresco-Levy U, Umbricht D (2012) Has an angel shown the way? Etiological and therapeutic implications of the PCP/NMDA model of schizophrenia. *Schizophr Bull* 38:958–966. [CrossRef Medline](#)
- Kay SR, Fiszbein A, Opler LA (1987) The positive and negative syndrome scale (PANSS) for schizophrenia. *Schizophr Bull* 13:261–276. [CrossRef Medline](#)
- Krawiecka M, Goldberg D, Vaughan M (1977) A standardized psychiatric assessment scale for rating chronic psychotic patients. *Acta Psychiatr Scand* 55:299–308. [CrossRef Medline](#)
- Lam NH, Borduqui T, Hallak J, Roque AC, Anticevic A, Krystal JH, Wang XJ, Murray JD (2017) Effects of altered excitation-inhibition balance on decision making in a cortical circuit model. [bioRxiv:100347](#).
- Langdon R, Ward PB, Coltheart M (2010) Reasoning anomalies associated with delusions in schizophrenia. *Schizophr Bull* 36:321–330. [CrossRef Medline](#)
- Marshall L, Mathys C, Ruge D, de Berker AO, Dayan P, Stephan KE, Bestmann S (2016) Pharmacological fingerprints of contextual uncertainty. *PLoS Biol* 14:e1002575. [CrossRef Medline](#)
- Mathys C, Daunizeau J, Friston KJ, Stephan KE (2011) A Bayesian foundation for individual learning under uncertainty. *Front Hum Neurosci* 5:39. [CrossRef Medline](#)
- Miller GA, Chapman JP (2001) Misunderstanding analysis of covariance. *J Abnorm Psychol* 110:40–48. [CrossRef Medline](#)
- Moritz S, Woodward TS (2005) Jumping to conclusions in delusional and non-delusional schizophrenic patients. *Br J Clin Psychol* 44:193–207. [CrossRef Medline](#)
- Moutoussis M, Bentall RP, El-Deredy W, Dayan P (2011) Bayesian modeling of jumping-to-conclusions bias in delusional patients. *Cogn Neuro-psychiatry* 16:422–447. [CrossRef Medline](#)
- Murray JD, Anticevic A, Gancsos M, Ichinose M, Corlett PR, Krystal JH, Wang XJ (2014) Linking microcircuit dysfunction to cognitive impairment: effects of disinhibition associated with schizophrenia in a cortical working memory model. *Cereb Cortex* 24:859–872. [CrossRef Medline](#)
- Nelson HE (1982) National Adult Reading Test (NART): for the Assessment of Premorbid Intelligence in Patients with Dementia: Test Manual. NFER-Nelson.
- Pantelis C, Barber FZ, Barnes TR, Nelson HE, Owen AM, Robbins TW (1999) Comparison of set-shifting ability in patients with chronic schizophrenia and frontal lobe damage. *Schizophr Res* 37:251–270. [CrossRef Medline](#)
- Peters E, Garety P (2006) Cognitive functioning in delusions: a longitudinal analysis. *Behav Res Ther* 44:481–514. [CrossRef Medline](#)
- Peters ER, Joseph SA, Garety PA (1999) Measurement of delusional ideation in the normal population: introducing the PDI (Peters et al. Delusions Inventory). *Schizophr Bull* 25:553–576. [CrossRef Medline](#)
- Powers AR, Mathys C, Corlett PR (2017) Pavlovian conditioning-induced hallucinations result from overweighting of perceptual priors. *Science* 357:596–600. [CrossRef Medline](#)
- Raine A (1991) The SPQ: a scale for the assessment of schizotypal personality based on DSM-III-R criteria. *Schizophr Bull* 17:555–564. [CrossRef Medline](#)
- Rao RP, Ballard DH (1999) Predictive coding in the visual cortex: a functional interpretation of some extra-classical receptive-field effects. *Nat Neurosci* 2:79–87. [CrossRef Medline](#)
- Redish AD, Jensen S, Johnson A, Kurth-Nelson Z (2007) Reconciling reinforcement learning models with behavioral extinction and renewal: implications for addiction, relapse, and problem gambling. *Psychol Rev* 114:784–805.
- Rigoux L, Stephan KE, Friston KJ, Daunizeau J (2014) Bayesian model selection for group studies—revisited. *Neuroimage* 84:971–985. [CrossRef Medline](#)
- Rolls ET, Loh M, Deco G, Winterer G (2008) Computational models of schizophrenia and dopamine modulation in the prefrontal cortex. *Nat Rev Neurosci* 9:696–709. [CrossRef Medline](#)
- Schlagenhauf F, Huys QJ, Deserno L, Rapp MA, Beck A, Heinze HJ, Dolan R, Heinz A (2014) Striatal dysfunction during reversal learning in unmedicated schizophrenia patients. *Neuroimage* 89:171–180. [CrossRef Medline](#)
- Speechley WJ, Whitman JC, Woodward TS (2010) The contribution of hypersaliency to the “jumping to conclusions” bias associated with delusions in schizophrenia. *J Psychiatry Neurosci* 35:7–17. [CrossRef Medline](#)
- Standage D, You H, Wang DH, Dorris MC (2013) Trading speed and accuracy by coding time: a coupled-circuit cortical model. *PLoS Comput Biol* 9:e1003021. [CrossRef Medline](#)
- Stephan KE, Penny WD, Daunizeau J, Moran RJ, Friston KJ (2009) Bayesian model selection for group studies. *Neuroimage* 46:1004–1017. [CrossRef Medline](#)
- Stuke H, Stuke H, Weinhhammer VA, Schmack K (2017) Psychotic experiences and overhasty inferences are related to maladaptive learning. *PLoS Comput Biol* 13:e1005328. [CrossRef Medline](#)
- Sutton R (1992) Gain adaptation beats least squares? <http://citeseer.ist.psu.edu/viewdoc/summary?doi=10.1.1.55.9218>. Accessed Jan. 26, 2018.
- Vinckier F, Gaillard R, Palminteri S, Rigoux L, Salvador A, Fornito A, Adapa R, Krebs MO, Pessiglione M, Fletcher PC (2016) Confidence and psychosis: a neuro-computational account of contingency learning disruption by NMDA blockade. *Mol Psychiatry* 21:946–955. [CrossRef Medline](#)
- Vossel S, Mathys C, Daunizeau J, Bauer M, Driver J, Friston KJ, Stephan KE (2014) Spatial attention, precision, and Bayesian inference: a study of saccadic response speed. *Cereb Cortex* 24:1436–1450. [CrossRef Medline](#)
- Waltz JA, Kasanova Z, Ross TJ, Salmeron BJ, McMahon RP, Gold JM, Stein EA (2013) The roles of reward, default, and executive control networks in set-shifting impairments in schizophrenia. *PLoS One* 8:e57257. [CrossRef Medline](#)
- Wang XJ (2013) The prefrontal cortex as a quintessential “cognitive-type” neural circuit: working memory and decision making. <https://nyuscholars.nyu.edu/en/publications/the-prefrontal-cortex-as-a-quintessential-cognitive-type-neural-c>. Accessed January 31, 2018.
- Wechsler D (1997) WAIS-III: administration and scoring manual: Wechsler Adult Intelligence Scale. San Antonio: Psychological Corporation.
- Woodward TS, Moritz S, Cuttler C, Whitman JC (2006) The contribution of a cognitive bias against disconfirmatory evidence (BADE) to delusions in schizophrenia. *J Clin Exp Neuropsychol* 28:605–617. [CrossRef Medline](#)
- Young HF, Bentall RP (1997) Probabilistic reasoning in deluded, depressed and normal subjects: effects of task difficulty and meaningful versus non-meaningful material. *Psychol Med* 27:455–465. [CrossRef Medline](#)