# *Paranoia as a deficit in non-social belief updating*

Erin J. Reed[1,2], Stefan Uddenberg[3], Praveen Suthaharan[4], Christoph D. Mathys[5,6], Jane R. Taylor[4], Stephanie M. Groman[4], and Philip R. Corlett[4*]

[1]Interdepartmental Neuroscience Program, Yale School of Medicine, New Haven, CT, USA. [2]Yale MD-PhD Program, Yale School of Medicine, New Haven, CT, USA. [3]Department of Psychology, Princeton University, Princeton, NJ, USA. [4]Department of Psychiatry, Connecticut Mental Health Center, Yale University, New Haven, CT, USA. [5]Scuola Internazionale Superiore di Studi Avanzati (SISSA), Trieste, Italy. [6]Translational Neuromodeling Unit (TNU), Institute for Biomedical Engineering, University of Zurich and ETH Zurich, Zurich, Switzerland.
*email: philip.corlett@yale.edu

54 **Abstract**

55

56 Paranoia is the belief that harm is intended by others. It may arise from selective pressures to infer and avoid
57 social threats, particularly in ambiguous or changing circumstances. We propose that uncertainty may be
58 sufficient to elicit learning differences in paranoid individuals, without social threat. We used reversal learning
59 behavior and computational modeling to estimate belief updating across individuals with and without mental
60 illness, online participants, and rats chronically exposed to methamphetamine, an elicitor of paranoia in
61 humans. Paranoia is associated with a stronger prior on volatility, accompanied by elevated sensitivity to
62 perceived changes in the task environment. Methamphetamine exposure in rats recapitulates this impaired
63 uncertainty-driven belief updating and rigid anticipation of a volatile environment. Our work provides evidence
64 of fundamental, domain-general learning differences in paranoid individuals. This paradigm enables further
65 assessment of the interplay between uncertainty and belief-updating across individuals and species.
66
67

68    Paranoia is excessive concern that harm will occur due to deliberate actions of others[1]. It manifests along a

69    continuum of increasing severity[2-5]. Fleeting paranoid thoughts prevail in the general population[6]. A survey of

70    over 7,000 individuals found that nearly 20% believed people were against them at times in the past year;

71    approximately 8% felt people had intentionally acted to harm them[4]. At a national level, paranoia may fuel

72    divisive ideological intolerance. Historian Richard Hofstadter famously described catastrophizing, context

73    insensitive political discourse as the 'paranoid style':

74        *"*The paranoid spokesman sees the fate of conspiracy in apocalyptic terms—he traffics in the birth and

75        death of whole worlds, whole political orders, whole systems of human values. He is always manning the

76        barricades of civilization. *He constantly lives at a turning point* [emphasis added].*"*[7]

77

78    At its most severe, paranoia manifests as rigid beliefs known as delusions of persecution. These delusions

79    occur in nearly 90% of first episode psychosis patients[8]. Psychostimulants also elicit severe paranoid states.

80    Methamphetamine evokes new paranoid ideation particularly after repeated exposure  or escalating doses

81    (86% and 68%, respectively, in a survey of methamphetamine users)[9].

82

83    Paranoia has thus far defied explanation in mechanistic terms. Sophisticated Game Theory driven approaches

84    (such as the Dictator Game[10,11]) have largely re-described the phenomenon — people who are paranoid have

85    difficulties in laboratory tasks that require trust[12]. However, this is not driven by personal threat per se, but by

86    negative representations of others[10,11]. We posit that such representations are learned[13,14], via the same

87    fundamental learning mechanisms[15] that underwrite non-social learning in non-human species[16]. We

88    hypothesize that aberrations to these domain-general learning mechanisms underlie paranoia. One such

89    mechanism involves the judicious use of uncertainty to update beliefs: Expectations about the noisiness of the

90    environment constrain whether we update beliefs or dismiss surprises as probabilistic anomalies. The higher

91    the **expected uncertainty** (i.e., 'I expect variable outcomes'), the less surprising an atypical outcome may be,

92    and the less it drives belief updates ('this variation is normal'). **Unexpected uncertainty**, in contrast, describes

93    perceived change in the underlying statistics of the environment[17-19] (i.e. 'the world is changing'), which may

94    call for belief revision.

95

96    Since excessive unexpected uncertainty is a signal of change, it might drive the recategorization of allies as

97    enemies, which is a tenet of evolutionary theories of paranoia[12]. We tested the hypothesis that this drive to

98    flexibly recategorize associations extends to non-social, domain-general inferences. We dissected learning

99    mechanisms under **expected** and **unexpected uncertainty –** probabilistic variation and changes in underlying

100   task structure **(volatility).** Here, volatility is a property of the task. <u>Unexpected uncertainty is the perception of</u>

101   <u>that volatility</u>.  Participants completed a non-social, three-option learning task which challenged them to form

102   and revise associations between stimuli (colored card decks) and outcomes (points rewarded and lost), in

103   addition to their beliefs about the volatility of the task environment. They encountered **expected uncertainty**

104    as probabilistic win or loss feedback ('each option yields positive and negative outcomes, but in different

105    amounts,'), and **unexpected uncertainty** as reassignment of reward probabilities between options

106    ('sometimes the best option may change,' **reversal events**). Although reversal events elicit **unexpected**

107    **uncertainty** by driving re-evaluation of the options, participants increasingly anticipate reversals and develop

108    expectations about the stability of the task environment. We implemented an additional task manipulation: a

109    shift in the underlying probabilities themselves (**contingency transition, unsignaled to the participants**), that

110    effectively changes task volatility. Armed with the task structure and participants' choices, we applied **a**

111    **Hierarchical Gaussian Filter (HGF)**[20,21] model which allowed us to infer participants' initial beliefs (i.e., priors)

112    about task volatility, their readiness to learn about changes in the task volatility itself (meta-volatility learning

113    rate) and learning rates that captured their expected and unexpected uncertainty regarding the task**.**

114

115    We examined the behavioral and computational correlates of paranoia both in-person and in a large online

116    sample, spanning patients and healthy controls with varying degrees of paranoia. We also undertook a pre-

117    clinical replication in rodents exposed chronically to saline or methamphetamine to determine whether a drug

118    known to elicit paranoia in humans might induce similar perceptions of unexpected uncertainty, without

119    contingency transition[22]. We predicted that people with paranoia and rats administered methamphetamine

120    would exhibit stronger priors on volatility, facilitating aberrant learning through unexpected uncertainty. We

121    further hypothesized that this learning style would manifest as frequent and unnecessary choice switching

122    (increased choice stochasticity and 'win-switch' behavior) rather than increased sensitivity to negative

123    feedback (increased 'lose-switch' behavior / decreased 'lose-stay' behavior).

124

125    **Results**

126

127        We analyzed belief updating across three reversal-learning experiments (Fig. 1): an in laboratory pilot of

128    patients and healthy controls, stratified by stable, paranoid personality trait (Experiment 1); four online task

129    variants administered to participants via the Amazon Mechanical Turk (MTurk) marketplace (Experiment 2);

130    and a re-analysis of data from rats on chronic, escalating doses of methamphetamine, a translational model of

131    paranoia (Experiment 3)[22].

132

133    **Experiment 1.** First, we explored trans-diagnostic associations between paranoia and reversal-learning in-

134    person. Participants with and without psychiatric diagnoses (mood disorders: anxiety, depression, bipolar

135    disorder, n=8;  schizophrenia spectrum: schizophrenia or schizoaffective disorder, n=8; and healthy controls,

136    n=16), completed questionnaire versions of the *Structured Clinical Interview for DSM-IV Axis II Personality*

137    *Disorders* (SCID-II) screening assessment[23], Beck's Anxiety Inventory (BAI)[24], Beck's Depression Inventory

138    (BDI)[25], and demographic assessments (Table 1). Approximately two-thirds of participants endorsed three or

139    fewer items on the SCID-II paranoid personality subscale (median=1 item). Participants who endorsed four or

4

140 more items were classified as high paranoia (*n*=11), consistent with the diagnostic threshold for paranoid

141 personality disorder. Low paranoia (*n*=21) and high paranoia groups did not differ significantly by age, nor were

142 there significant group associations with gender, educational attainment, ethnicity, or race, although a larger

143 percentage of paranoid participants identified as racial minorities or "not specified" (Table 1). Diagnostic

144 category (i.e., healthy control, mood disorder, or schizophrenia spectrum) was significantly associated with

145 paranoia group membership, $\chi^2$ (2, *n*=32)=12.329, *P*=0.002, Cramer's V=0.621, as was psychiatric medication

146 usage, $\chi^2$ (1, *n*=32)=9.871, *P*=0.003, Cramer's V=0.555. These differences were due to the higher proportion

147 of healthy controls in the low paranoia group. As expected, paranoia, BAI, and BDI scores were significantly

148 elevated in the high paranoia group relative to low paranoia controls (Table 1; paranoia: mean difference

149 (MD)=0.536, CI=[0.455,0.618], *t*(30)=13.476, *P*=2.92E-14, Hedges' *g*=5.016; BAI: MD=0.585, CI=[0.239,

150 0.931], *t*(30)=3.453, *P*=0.002, Hedges' *g*=1.285, MD=-0.585; BDI: MD=0.427, CI=[0.078, 0.775],

151 *t*(11.854)=2.67, *P*=0.021, Hedges' *g*=1.255).

152

153 Participants completed a three-option reversal-learning task in which they chose between three decks of cards

154 with hidden reward probabilities (Fig. 1 a and b). They selected a deck on each turn and received positive or

155 negative feedback (+100 or -50 points, respectively). They were instructed to find the best deck with the caveat

156 that the best deck may change. Undisclosed to participants, reward probabilities switched among decks after

157 selection of the highest probability option in nine out of ten consecutive trials ("reversal events"). Thus, the task

158 was designed to elicit **expected uncertainty** (probabilistic reward associations) and **unexpected uncertainty**

159 (reversal events), requiring participants to distinguish probabilistic losses from change in the underlying deck

160 values. In addition, reward contingencies changed from 90%, 50%, and 10% chance of reward to 80%, 40%,

161 and 20% between the first and second halves of the task ("**contingency transition**"; block 1=80 trials, 90-50-

162 10%; block 2=80 trials, 80-40-20%, unsignaled to the participants). This transition altered the volatility of the

163 task environment, thereby making it more difficult to achieve reversals and often delaying their occurrence.

164 Successful achievement of reversals was contingent upon adapting stay-vs-switch strategies, thereby testing

165 subjects' abilities to update beliefs about the overall task volatility **("metavolatility learning").** High paranoia

166 subjects achieved fewer reversals (MD=-2.31, CI=[-4.504, -0.111,], *t*(30)=-2.145, *P*=0.04, Hedges' *g*=0.798),

167 but total points earned did not significantly differ, suggesting that there was no penalty for the different

168 behaviors expressed by the more paranoid subjects (Table 1). We predicted that paranoia would be associated

169 with unexpected uncertainty-driven belief updating.

170

171 **Experiment 2.** We aimed to replicate and extend our investigation of paranoia and reversal-learning in a larger

172 online sample. We administered three alternative task versions to control for the contingency transition (Fig.

173 1c). Version 1 (*n*=45 low paranoia, 20 high paranoia) provided a constant contingency of 90-50-10% reward

174 probabilities (Easy-Easy); version 2 (*n*=69 low paranoia, 18 high paranoia) provided a constant contingency of

175 80-40-20% (Hard-Hard); version 3 (*n*=56 low paranoia, 16 high paranoia) served to replicate Experiment 1 with

176    a contingency transition from 90-50-10% to 80-40-20% (Easy-Hard); version 4 ($n$=64 low paranoia, 19 high

177    paranoia) provided the reverse contingency transition, 80-40-20% to 90-50-10% (Hard-Easy). The stable

178    contingencies (versions 1 and 2) lacked contingency transitions. Versions 3 and 4 manipulated task volatility

179    mid-way, although the contingency transition was not signalled to participants. We predicted that high paranoia

180    participants would find versions 3 and 4 particularly challenging. Given that version 3 is easier to learn initially,

181    we expected participants to develop stronger priors and thus be more confounded by the contingency

182    transition, compared to version 4 participants.

183

184    Participants' demographic and mental health questionnaire responses did not differ significantly across task

185    version experiments (Table 2). Total points and reversals achieved suggest variations in task difficulty (Table

186    2, version effects: points earned, $F(3)$=232.88, $P$=4.16E-18, $\eta_p^2$=0.245; reversals achieved, $F(3)$=4.329,

187    $P$=0.005, $\eta_p^2$=0.042), but there was no significant association between task version and attrition rate (52.7%,

188    52.9%, 54.6%, and 53.1% attrition, respectively; $\chi^2(3)$=0.167, $P$=0.983, Cramer's V=0.015).

189

190    Across task versions, high paranoia participants endorsed higher BAI and BDI scores ($n$=73 high paranoia,

191    234 low paranoia; BAI: $F(1)$=38.752, $P$=1.63E-09, $\eta_p^2$=0.115; BDI: $F(1)$=74.528, $P$=3.62E-16, $\eta_p^2$=0.20; Table

192    2). Both correlated with paranoia (BAI: Pearson's $r$=0.450, $P$=1.09E-16, CI=[0.348, 0.55]; BDI: Pearson's

193    $r$=0.543, $P$=6.26E-25, CI=[0.448, 0.638]). Trial-by-trial reaction time did not differ significantly between low and

194    high paranoia (Table 2), but high paranoia participants earned fewer total points ($F(1)$=6.175, $P$=0.014,

195    $\eta_p^2$=0.020) and achieved fewer reversals ($F(1)$=5.762, p=0.017, $\eta_p^2$=0.019; Table 2). Deck choice

196    perseveration after negative feedback (lose-stay behavior) did not significantly differ by paranoia group, but

197    choice switching after positive feedback (win-switch behavior) was elevated in high paranoia (block 1:

198    $F(1)$=7.117, $P$=0.008, $\eta_p^2$=0.023; block 2: $F(1)$=9.918, $P$=0.002, $\eta_p^2$=0.032; Table 2).

199

200    **Experiment 3.** To translate across species, we performed a new analysis of published data from rats exposed

201    to chronic methamphetamine[22]. Rats chose between three operant chamber noseports with differing

202    probabilities of sucrose reward (70%, 30%, and 10%; Fig.1 d and e). Contingencies switched between the 70%

203    and 10% noseports after selection of the highest reinforced option in 21 out of 30 consecutive trials (Fig. 1e).

204    This task was most similar in structure to the first blocks of online versions 2 and 4. There was no increase in

205    unexpected volatility mid-way through the task. Rats were tested for 26 within-session reversal blocks (Pre-Rx,

206    $n$=10 per group), administered saline or methamphetamine according to a 23-day schedule mimicking the

207    escalating doses and frequencies of chronic human methamphetamine users[22], and tested once per week for

208    four weeks following completion of the drug regimen (Post-Rx; $n$=10 saline, 7 methamphetamine)[22]. Relative to

209    rats exposed to saline, those rats exposed to methamphetamine exhibited increased win-switch behavior,

210    similar to what we has observed in the high paranoia human participants, and additionally, unlike humans, they

211    perseverated after negative feedback[22].

212

213 **Computational modeling.** We employed hierarchical Gaussian filter (HGF) modeling to compare belief

214 updating across individuals with low and high paranoia, as well as across human participants and rats exposed

215 to methamphetamine (Table 3). We paired a three-level perceptual model with a softmax decision model

216 dependent upon third level volatility (Fig. 2a). We inverted the model from subject data (trial-by-trial choices

217 and feedback) to estimate parameters for each individual (Fig. 2b). Level 1 ($x_1$) characterizes trial-by-trial

218 perception of task feedback (win or loss in humans, reward or no reward in rats), Level 2 ($x_2$) distinguishes

219 stimulus-outcome associations (deck or noseport values), and Level 3 ($x_3$) renders perception of the overall

220 task volatility (i.e., frequency of reversal events, changes in the stimulus-outcome associations).

221

222 Belief trajectories were unique to each subject due to the probabilistic, performance-dependent nature of the

223 task, so we estimated initial beliefs (priors) for $x_2$ and $x_3$ ($\mathbf{\mu_2^0}$ and $\mathbf{\mu_3^0}$, respectively). We also estimated $\mathbf{\omega_2}$, the

224 tonic volatility of stimulus-outcome associations. Lower $\omega_2$ indicates that subjects are slower to adjust beliefs

225 about the value of each option; they maintain rigid beliefs about the underlying probabilities. The $\kappa$ parameter

226 captures the impact of phasic volatility on updating stimulus-outcome associations. In the setting of our

227 experiments, $\kappa$ approximates the influence of unexpected uncertainty. Higher $\kappa$ implies faster updating of

228 stimulus-outcome associations – that is, participants are more likely perceive volatility as reversal events.  Our

229 final parameter of interest, $\omega_3$, characterizes perception of 'meta-volatility,' such as changes in the frequency of

230 reversal events[26]. The lower $\omega_3$, the slower a subject is to adjust their volatility belief; they adhere more rigidly

231 to their volatility prior ($\mu_3^0$).

232

233 Priors did not differ between groups at $x_2$ (Table 3) but paranoid individuals and rats exposed to

234 methamphetamine exhibited elevated $\mu_3^0$, they expected greater task volatility (Fig. 2b, blue). In Experiment 1,

235 we observed an interaction between task block and paranoia group ($F(1)=5.344$, $P=0.028$, $\eta_p^2=0.151$; Table 1):

236 $\mu_3^0$ differed between high and low paranoia in both blocks (block 1, $F(1)=4.232$, $P=0.048$, $\eta_p^2=0.124$,

237 MD=0.658, CI=[0.005,1.312]; block 2, F(1)=7.497, $P=0.010$, $\eta_p^2=0.20$, MD=1.598, CI=[0.406, 2.789]), but only

238 low paranoia subjects significantly updated their priors between block 1 and block 2 ($F(30)=39.841$, $P=5.85E-$

239 07, $\eta_p^2=0.570$, MD=1.504, CI=[1.017, 1.99]). In Experiment 2, the analogous task design (version 3)

240 demonstrated significant effects of block ($F(1)=64.652$, $P=1.54E-11$, $\eta_p^2=0.480$, MD=1.303, CI=[0.980,1.627])

241 and paranoia ($F(1)=6.366$, $P=0.014$, $\eta_p^2=0.083$, MD=0.909, CI=[0.191, 1.628]; Table 1). Rats showed a similar

242 effect following methamphetamine exposure with a significant time (Pre-Rx, Post-Rx) by treatment

243 (methamphetamine, saline) interaction ($F(1)=5.159$, $P=0.038$, $\eta_p^2=0.256$; pre versus post methamphetamine

244 effect: $F(15)=12.186$, $P=0.003$, MD=1.265, CI=[-0.493, 2.037]; Pre-Rx mean [standard error]= -1.25 [0.56]

245 saline, -0.77 [0.80] methamphetamine; Post-Rx: $m$=-0.69 [0.74] saline, 0.58 [0.73] methamphetamine).

246 Random effects meta-analyses confirmed significant cross-experiment replication of elevated $\mu_3^0$ in human

247 participants with paranoia (in laboratory and online version 3; $MD_{META}$= 1.110, CI=[0.927, 1.292], $z_{META}$=11.929

248 , p=8.356E-33) and across humans with paranoia and rats exposed to methamphetamine (MD$_{META}$=2.090,

249 CI=[0.123, 4.056], $z_{META}$=2.083, p=0.037). Both paranoid humans and rats administered chronic

250 methamphetamine had strong beliefs that the task contingencies would change rapidly and unpredictably – in

251 other words, they expected frequent reversal events. Methamphetamine exposure made rats behave like

252 humans with high paranoia (Fig. 2b, Post-Rx condition, orange). This is particularly striking when compared to

253 human data from the first task block (before contingency transition), when task designs are most similar across

254 experiments.

255

256 Paranoid participants and methamphetamine exposed rats updated stimulus-outcome associations more

257 strongly in response to perceived volatility (e.g., correctly or incorrectly inferred reversals; Fig. 2b). $\kappa$ showed

258 significant paranoia group and block effects across the in laboratory experiment and online version 3 (Table 1;

259 paranoia effects, in laboratory: $F(1)$=7.599, $P$=0.010, $\eta_p^2$=0.202, MD=0.081, CI=[0.021, 0.140]; online version

260 3: $F(1)$=13.521, $P$=0.0005, $\eta_p^2$=0.162, MD=0.068, CI=[0.031-0.104]; MD$_{META}$ = 0.079, CI=[0.063, 0.095],

261 $z_{META}$=9.502 p=2.067E-21); see Table 3 for block effects). $\kappa$ increased from baseline in rats on

262 methamphetamine, yielding significant effects of treatment ($F(1)$=13.356, $P$=0.002, $\eta_p^2$=0.471, MD=0.045,

263 CI=[0.019, 0.072]) and time ($F(1)$=9.132, $P$=0.009, $\eta_p^2$=0.378, MD=0.041, CI=[0.012, 0.069]); however, the

264 interaction between time and treatment did not reach statistical significance (Table 3; Pre-Rx $m$=0.499 [0.015]

265 saline, 0.523 [0.040] methamphetamine; Post-Rx: $m$=0.518 [0.053] saline, 0.585 [0.029] methamphetamine).

266 Replication of group effects was significant across all three experiments (MD$_{META}$=2.063, CI=[0.341, 3.785],

267 $z_{META}$=2.348, p=0.019). Thus, learning was more strongly driven by unexpected uncertainty in high paranoia

268 participants and rats chronically administered methamphetamine; they were faster to interpret volatility as

269 reversal events than their low paranoia and saline exposed counterparts.

270

271　　Expected uncertainty ($\omega_2$) was decreased in paranoid participants and rats exposed to methamphetamine

272 (Fig. 2b). In laboratory and online (version 3), paranoid individuals were slower to update stimulus-outcome

273 associations in response to expected uncertainty(Table 1; $\omega_2$ paranoia effect, in laboratory: $F(1)$=4.186,

274 $P$=0.050, $\eta_p^2$=0.122, MD=-1.188, CI=[-2.375, -0.002]; online version 3: $F(1)$=8.7, $P$=0.004, $\eta_p^2$=0.111, MD=-

275 0.993, CI=[-1.665, -0.322]; MD$_{META}$=-1.154 , CI=[-1.455, -0.853], $z_{META}$=-7.521, p=5.450E-14). The effects of

276 methamphetamine exposure in rats were consistent (MD$_{META}$=-1.992 , CI=[-3.318, -0.665], $z_{META}$=-2.943,

277 p=0.003) yet more striking, with a strongly negative $\omega_2$ accounting for the more pronounced lose-stay behavior

278 or perseveration in rats (time by treatment interaction, $F(1)$=18.454, $P$=0.001, $\eta_p^2$=0.552; pre versus post

279 methamphetamine: $F(1)$=42.242, $P$=1.0E-5$^{22}$, $\eta_p^2$=0.738, MD=-1.604, CI=[-2.130, -1.078]; Pre-Rx $m$=0.198

280 [0.33] saline, -0.036 [0.42] methamphetamine; Post-Rx: $m$=-0.023 [0.56] saline, -1.640 [0.71]

281 methamphetamine). High paranoia humans and rats exposed to methamphetamine maintained rigid beliefs

282  about the underlying option probabilities relative to low paranoia and saline controls. This was associated with

283  perseverative behavior in the rats but not in humans.

284

285  Meta-volatility learning ($\omega_3$) was similarly decreased across paranoia and methamphetamine exposed groups

286  (in laboratory, online version 3, and rats: $MD_{META}$=-1.155, CI=[-2.139, -0.171], $z_{META}$=-2.3, p=0.021),

287  suggesting more reliance on expected task volatility (i.e., anticipated frequency of reversal events) than on

288  actual task feedback. In laboratory, we observed a block by paranoia group interaction (Table 1, $F$(1)=6.948,

289  $P$=0.010, $\eta_p^2$=0.188). Post-hoc tests differentiated first and second blocks for the low paranoia group only

290  ($F$(1)=26.640, $P$=1.5E-5, $\eta_p^2$=0.470, MD=-0.876, CI=[-1.222, -0.529]). The paranoia effect did not reach

291  statistical significance for online version 3 (block effect only, $F$(1)=14.932, $P$=0.0002, $\eta_p^2$=0.176, MD=-0.692,

292  CI=[-1.050, -0.335]; Table 3), but meta-analytic random effects analysis confirms a significant paranoia group

293  difference (in laboratory and online version 3: $MD_{META}$=-0.341, CI=[-0.522, -0.159], $z_{META}$=-3.68, p=0.0002).

294  Methamphetamine exposure rendered $\omega_3$ more negative in rats (time by treatment interaction, ($F$(1)=9.058,

295  $P$=0.009, $\eta_p^2$=0.376; pre versus post methamphetamine: $F$(1)=30.668, P=5.7E-5, $\eta_p^2$=0.672, MD=-1.210, CI=[-

296  1.676, -0.745]; Pre-Rx m=-0.692 [0.44] saline, -0.607 [0.51] methamphetamine; Post-Rx: $m$=-1.044 [0.44]

297  saline, -1.817 [0.32] methamphetamine). These data indicate that paranoia and methamphetamine are

298  associated with slower learning about changes in task volatility, suggesting greater reliance on volatility priors

299  than task feedback.

300

301  In summary, our modeling analyses suggest the following about paranoia in humans and methamphetamine

302  exposed animals: they expect the task to be volatile (high $\mu_3^0$), their expectations about task volatility are more

303  rigid (low $\omega_3$), and they confuse probabilistic errors and task volatility as a signal that the task has

304  fundamentally changed (high $\kappa$, low $\omega_2$).

305

306  We applied **False Discovery Rate (FDR) correction** for multiple comparisons of each model parameter[27]. $\kappa$

307  group effects survived corrections within each experiment (Table 4). In addition to $\kappa$, $\mu_3^0$ survived for

308  experiment 1; $\mu_3^0$ and $\omega_2$ survived in online version 3; and $\mu_3^0$, $\omega_2$, and $\omega_3$ survived in experiment 3 as group

309  effects. Such correction is not yet standard practice with this modeling approach[26,28,29] but we believe it should

310  be, and when effects survive correction we should increase our confidence in them.

311

312  **Paranoia effects across task versions.** To examine the relationship between beliefs about contingency

313  transition and paranoia within our HGF parameters, we performed split-plot, repeated measures ANOVAs

314  across all four task versions. Paranoia group effects were specific to versions of the task in which we explicitly

315  manipulated uncertainty via contingency transition which increased volatility (Fig. 3, Table 5, versions 3 and 4).

316  Specifically, we observed paranoia by version interactions for $\kappa$ ($F$(3)=4.178, $P$=0.006, $\eta_p^2$=0.040) and $\omega_2$

317    ($F(3)$=2.809, $P$=0.040, $\eta_p^2$=0.027; Table 2). Post-hoc tests confirmed that significant paranoia group effects

318    were restricted to version 3 ($\kappa$: $F(1)$=12.230, $P$=0.001, $\eta_p^2$=0.039, MD=0.068, CI=[0.03,0.106]; $\omega_2$: $F(1)$=8.734,

319    $P$=0.003, $\eta_p^2$=0.028, MD=-0.993, CI=[-1.655, -0.332]) and a trend for version 4 ($\omega_2$: $F(1)$=2.909, $P$=0.089,

320    $\eta_p^2$=0.010, MD=-0.528, CI=[-1.138, 0.081], Fig. 3a). $\mu_3^0$ also exhibited a paranoia by version trend (Table 2,

321    $F(3)$=2.329, $P$=0.075, $\eta_p^2$=0.023), largely driven by version 3 ($F(1)$=6.206, $P$=0.013, $\eta_p^2$=0.020, MD=0.909,

322    CI=[0.191, 1.628]; Fig. 3a). There were no significant paranoia effects or interactions for $\omega_3$ (Table 5). In sum,

323    our contingency shift manipulation – from easily discerned options to underlying probabilities that are closer

324    together – increased unexpected uncertainty the most, particularly in highly paranoid participants, compared to

325    the other task versions.

326

327

328 **Covariate analyses.** We completed three ANCOVAs for each HGF parameter derived from Experiment 2:

329 demographics (age, gender, ethnicity, and race); mental health factors (medication usage, diagnostic category,

330 BAI score, and BDI score); and metrics and correlates of global cognitive ability (educational attainment,

331 income, and cognitive reflection; Tables 6 and 7). For $\kappa$, our metric of unexpected uncertainty, the paranoia by

332 version interaction remained robust across all three ANCOVAs (demographics: $F(3)=3.753$, $P=0.011$,

333 $\eta_p^2=0.037$; mental health: $F(3)=4.417$, $P=0.005$, $\eta_p^2=0.049$; cognitive: $F(3)=4.304$, $P=0.005$ $\eta_p^2=0.043$). The

334 paranoia by version trend of $\mu_3^0$ diminished with inclusion of demographic, mental health, and cognitive

335 covariates (demographic: $F(3)=1.997$, $P=0.119$, $\eta_p^2=0.020$; mental health: $F(3)=1.942$, $P=0.123$, $\eta_p^2=0.022$;

336 cognitive: $F(3)=2.193$, $P=0.089$, $\eta_p^2=0.022$). The paranoia by version interaction for $\omega_2$ was robust to mental

337 health and cognitive factors ($F(3)=3.617$, $P=0.014$, $\eta_p^2=0.041$; $F(3)=3.017$, $P=0.030$, $\eta_p^2=0.030$). A paranoia

338 group effect and paranoia by version trend remained with inclusion of demographics ($\omega_2$, paranoia effect:

339 $F(1)=4.275$, $P=0.040$, $\eta_p^2=0.014$; interaction: $F(3)=2.507$, $P=0.059$, $\eta_p^2=0.025$). Thus $\kappa$ – participants'

340 perception of **unexpected uncertainty** – was the only parameter whose main effect of paranoia (higher $\kappa$ in

341 high paranoia participants) and paranoia-by-version interaction (higher $\kappa$ in high paranoia participants as a

342 function of increasing unexpected volatility in version 3) survived covariation for demographic, mental health

343 and cognitive covariates. We are most confident that high paranoia participants have higher **unexpected**

344 **uncertainty** which drives their excessive updating of stimulus-outcome associations.

345

346

347 **Relationships between Parameters and Paranoia**

348 We found a significant correlation between $\kappa$ and paranoia scores (Fig. 4). However, depression and anxiety

349 were also related to $\kappa$, and indeed, paranoia and depression correlate with one another, in our data and in

350 other studies[30]. In order to explore commonalities among the rating scales in the present data, we performed a

351 principle components analysis (Fig. 5), identifying three principle components. The first principle component

352 (PC 1) explained 82.3% of the variance in the scales and loaded similarly on anxiety, depression, and

353 paranoia. It correlated significantly with kappa (r=0.272, p=0.021). Depression, anxiety and paranoia all

354 contribute to PC1. We suggest that this finding is consistent with the idea that depression and anxiety

355 represent contexts in which paranoia can flourish and likewise, harboring a paranoid stance toward the world

356 can induce depression and anxiety.

357

358 **Multiple regression.** In order to make the case that our observations were most relevant to paranoia, we

359 examined the effects of paranoia, anxiety, and depression on $\kappa$ within the online version 3 dataset with multiple

360 regression. A significant regression equation was found (F(3,68)=3.681, p=0.016), with an $R^2$ of 0.140.

361 Participants' predicted $\kappa$ equaled 0.486 + 0.062 (PARANOIA)+0.012 (BDI) -0.006 (BAI). Paranoia was a

362   significant predictor of $\kappa$ ($\beta$=0.343, t=2.470, p=0.016, CI=[0.012, 0.113]) but depression and anxiety were not

363   (BDI: $\beta$=0.086, t=0.423, p=0.674, CI=[-0.043, 0.066]; BAI: $\beta$=-0.043, t=-0.218, p=0.828, CI=[-0.063, 0.050]).

364   Examination of correlation plots for $\kappa$ (Fig. 4) revealed a much stronger relationship when analyses were

365   restricted to individuals with paranoia scores greater than 0 (i.e., endorsement of at least one item); among

366   participants who denied all questionnaire items, a minority (seven out of 32) exhibited elevated $\kappa$. To account

367   for the possibility that some individuals with severe paranoia may avoid disclosing sensitive information, we

368   performed additional analyses of participants who endorsed one or more paranoia item. The correlation

369   between paranoia and $\kappa$ in the first block of the task increases from r=0.3, p=0.011, CI=[0.074, 0.497] (all

370   participants, n=72) to r=0.588, p=8.0E-5, CI=[0.335, 0.762] (participants with paranoia > 0, n=39). In this

371   subset, a significant regression equation was also found (F(3,35)=6.322, p=0.002), with an $R^2$ of 0.351 (Fig.

372   4.). Participants' predicted $\kappa$ was equal to 0.432 + 0.150 (PARANOIA)+0.013 (BDI) -0.004 (BAI). Paranoia was

373   a significant predictor of $\kappa$ ($\beta$=0.538, t=2.983, p=0.005, CI=[0.048, 0.252]) but depression and anxiety were not

374   (BDI: $\beta$=0.111, t=0.494, p=0.624, CI=[-0.041, 0.067]; BAI: $\beta$=-0.035, t=-0.163, p=0.872, CI=[-0.057, 0.049]).

375   Thus, paranoia predicts kappa across participants. Anxiety and depression do not predict kappa**.**

376

377   **Behavior and simulations.** Win-switching was the prominent behavioral feature of both paranoid participants

378   and rats exposed to methamphetamine (Table 1, Table 2[22]). Collapsed across blocks and task versions, our

379   Experiment 2 data demonstrated a main effect of paranoia group (Fig. 3b; $F$(1)=9.207, $P$=0.003, $\eta_p^2$=0.030,

380   MD=0.059, CI=[0.021, 0.097]; version trend: $F$(3)=2.263 $P$=0.081, $\eta_p^2$=0.022; low paranoia: $m$=0.06 [0.01], high

381   paranoia: $m$=0.12 [0.02]). To elucidate whether this behavior was stochastic or predictable (e.g., switching

382   back to a previously rewarding option), we calculated U-values[31], a metric of behavioral variability employed by

383   behavioral ecologists (increasingly an inspiration for human behavioral analysis[32]), particularly with regards to

384   predator-prey relationships[33]. When a predator is approaching a prey animal, the prey's best course of action is

385   to behave randomly, or in a *protean* fashion, in order to evade capture[33]. The more protean or stochastic the

386   behavior, the closer to the U-value is to 1. Across task blocks, paranoid participants exhibited elevated choice

387   stochasticity (paranoia by version interaction, $F$(3)=3.438, $P$=0.017, $\eta_p^2$=0.033; Table 2). Post-hoc tests

388   indicate that this stochasticity was specific to versions with contingency transition, suggesting a relationship to

389   unexpected uncertainty (Fig. 3b; version 3, $F$(1)=17.585, $P$=3.6E-5, $\eta_p^2$=0.056, MD=0.071, CI=[0.038, 0.104];

390   version 4, $F$(1)=6.397, $P$=0.012, $\eta_p^2$=0.021, MD=0.039, CI=[0.009, 0.07]). Our task manipulation, increasing

391   unexpected volatility, increases win-switching behavior and stochastic choice more in more paranoid

392   participants.

393

394   To test the propriety of our model, we simulated data for each subject in online version 3 and determined

395   whether or not key behavioral effects (Fig. 7a, Table 1, Table 9) were present. Using individually estimated

396   HGF parameters to generate ten simulations per participant, we recapitulated both elevated win-switch

397   behavior (paranoia effect, $F$(1)=15.394, $P$=2.01E-4, $\eta_p^2$=0.180, MD=0.186, CI=[0.091, 0.28]) and choice

398   stochasticity (U-value; paranoia effect, $F(1)=13.362$, *P*=0.0005, $\eta_p^2$=0.160, MD=0.065, CI=[0.030, 0.101]) in

399   simulated paranoid participants (Fig. 7b; simulated win-switch rate, low paranoia: *m*=0.24 [0.02], high paranoia:

400   *m*=0.43 [0.04]; simulated U-value, low paranoia: *m*=0.851 [0.008], high paranoia: *m*=0.916 [0.016]). Neither

401   real nor simulated data showed any significant relationship between lose-stay behavior and paranoia (Table 1,

402   Table 2, Table 9). To demonstrate the effects of parameters on task performance, we performed additional

403   simulations in which we doubled or halved a single parameter at a time from the baseline average of low

404   paranoia participants. These results confirmed the impact of $\kappa$, $\omega_2$, and $\omega_3$ on win-shift behavior (Fig. 4).

405   Parameter recovery revealed significant correlations for $\kappa$ and $\omega_2$ between original subject parameters and

406   those estimated from simulations (Fig. 6; $\omega$: r=0.702, p=2.52E-11, CI=[0.557, 0.805]; $\kappa$: r=0.305, p=0.011,

407   CI=[0.072, 0.506]).  Higher level parameters ($\omega_3$, $\mu_3^0$) were less consistently recovered, as noted in previous

408   publications[34]. Thus, the model we chose, with meta-volatility and three coupled layers of belief, successfully

409   simulates the key features of paranoid behavior: higher win-switching and stochastic choice.

410

411   **Alternate Models:** Our model is complex and other simpler reinforcement learning models might explain

412   behavior on this task. Given the win-switching behavior we sought to understand, we fit a model from Lefebvre

413   and colleagues that instantiated biased belief updating via differential weighting of positive and negative

414   prediction errors[35]. Fitting this model to online version 3, we saw no significant paranoia group differences in

415   learning rates for positive or negative prediction errors in parameters derived from all 180 trials (independent

416   samples t-test: $\alpha^+$, *t*(70)=-0.532, p=0.597; $\alpha^-$, *t*(70)=0.963, p=0.339), nor did we see any significant

417   block*paranoia or paranoia group effects by repeated measures ANOVA (block*paranoia: $\alpha^+$, *F*(1)=0.188,

418   p=0.732, $\alpha^-$, *F*(1)=0.378, p=0.540; paranoia group: $\alpha^+$, *F*(1)=0.243, p=0.623, $\alpha^-$, *F*(1)=1.292, p=0.260). See

419   Table 10.

420

421   We can also simplify within our hierarchical Gaussian Filter framework. The model we chose had three layers

422   of beliefs and the highest level seemed to capture most of the task and paranoia effects of interest (Fig. 8). To

423   confirm this suspicion, we removed the third layer, fitting an HGF model that had beliefs about outcomes and

424   deck values but no beliefs about volatility, no unexpected volatility learning rate, nor meta-volatility. This model

425   failed to capture the task effects or group differences in its parameters (see Table 10).

426

427   Therefore, a more complicated model, one that captures higher-level beliefs about contingency transitions or

428   learning when to learn, seems most appropriate, and indeed, that type of model was able to simulate the key

429   features of our data[36]. Future work will compare and contrast different potential computational models included,

430   but not limited to Bayesian Hidden State Markov Models[37], as well as switching[38] and volatile Kalman Filters[39].

431

432   **Clustering analysis.**

433  Given the apparent similarity in effects of paranoia and methamphetamine in humans and rats, respectively

434  (Fig. 2b), we searched for latent structure in our data using two-step cluster analysis[40]. This approach sorts

435  subjects into groups (clusters) on the basis of some experimenter-selected variables such as estimated model

436  parameters. The goal is to find distinct subsets in the data such that each cluster exhibits a cohesive pattern of

437  relationships between the variables. Whereas some clustering approaches require the experimenter to

438  predefine the expected number of clusters, two step-clustering determines both the optimal number of clusters

439  and the composition of each cluster. The greater the similarity (or homogeneity) within a group and the greater

440  the difference between groups, the better the clustering.

441

442  Considering that paranoia and methamphetamine exposure share a pattern of elevated $\mu_3^0$ and $\kappa$ accompanied

443  by decreased $\omega_2$ and $\omega_3$ (Table 8), we hypothesized that these four variables would yield a distinct cluster: a

444  'paranoid style' across species.  We analyzed $\mu_3^0$, $\kappa$, $\omega_2$, and $\omega_3$ estimates derived from the first block of

445  experiment 1 and online version 3 (pre-context change data, because rats do not experience a context shift)

446  with post-chronic exposure rat data (methamphetamine and saline). We identified two clusters with good

447  cohesion and separation, meaning that subjects sorted into two groups (each containing rodents and humans)

448  whose parameters travelled in such a way that their values were close to the centroid or mean of the cluster

449  they were in and as far as possible from the centroid of the other cluster (average silhouette coefficient=0.7;

450  cluster size ratio=2.46; Fig. 9a). All parameters contributed to clustering; $\kappa$ contributed most strongly (Fig. 9b).

451  Importantly, the cluster solution did not separate rats from humans (despite the differences in task structure,

452  incentives, manipulanda, and phylogeny). Relative to the overall distribution, Cluster 1 was characterized by

453  high $\kappa$ and $\mu_3^0$, and decreased $\omega_2$ and $\omega_3$. Cluster 1 membership was significantly associated with high

454  paranoia and methamphetamine exposure, $\chi^2$(1, $n$=121)=29.447, $P$=5.75E-8, Cramer's V=0.493 (Fig. 9c).

455  Notably, no participants in the low paranoia group with paranoia scores above zero were ascribed Cluster 1

456  membership. The cluster solution was robust to validation by split-half analysis (removing half of the

457  participants and repeating the clustering), removal of the rat subjects, and removal of human participants. In

458  each case, we identified two clusters with good cohesion and separation (**Split-half 1**, n=19 cluster 1, 42

459  cluster 2: silhouette coefficient = 0.6; **Split-half 2**, n = 17 cluster 1, 43 cluster 2: silhouette coefficient = 0.7; **No

460  Rat**, n=26 cluster 1, 78 cluster 2: silhouette coefficient = 0.7; **Rat Only**, n=6 cluster 1, 11 cluster 2: silhouette

461  coefficient = 0.7).  In summary, paranoid participants and methamphetamine-exposed rats cluster together

462  (high $\mu_3^0$, high $\kappa$, low $\omega_2$, and low $\omega_3$), suggesting that these parameters share an underlying generative

463  process and that paranoia and methamphetamine have similar effects on reversal-learning.

464

465  **Discussion**

466  During non-social probabilistic reversal-learning, paranoid individuals and rats chronically exposed to

467  methamphetamine have higher initial expectations of task volatility ($\mu_3^0$). In other words, they start the task

468  anticipating more changes in stimulus-outcome associations, and they switch choices readily and excessively

14

469   in anticipation of reversal events. By relying more on their expectations of volatility than on actual experience

470   (exemplified by switching even after positive feedback), they are slower to learn about changes in task

471   volatility. This manifests as decreased meta-volatility learning ($\omega_3$) and failure to significantly adjust $\mu_3^0$ after

472   contingency transitions. More paranoid individuals are similarly slower to adjust expected deck values (lower

473   $\omega_2$) but faster to attribute volatility to reversal events (elevated $\kappa$), perceiving change (**unexpected**

474   **uncertainty**) instead of normal statistical variation (**expected uncertainty**). They sit at Hofstadter's 'turning

475   point', constantly expecting change but never learning appropriately from it.

476

477   In the reversal learning literature, choice switching after positive feedback has garnered less attention than

478   perseverative behavior and sensitivity to negative feedback[41,42]. Individuals with depression and schizophrenia

479   seemingly perseverate less than healthy controls, but this has formerly been attributed to increased sensitivity

480   to negative feedback[42,43]. However, elevated win-switch tendencies have been reported in youths with bipolar

481   disorder, major depressive disorder, and anxiety disorder[44]. A prior study in people with schizophrenia

482   described excessive win-switch behavior that correlated with the severity of delusional beliefs and

483   hallucinations[42]. Likewise, an elevated prior on environmental volatility ($\mu_3^0$) and higher sensitivity to this

484   volatility ($\kappa$) have been observed in HGF analyses of 2-choice probabilistic reversal-learning in medicated and

485   unmedicated patients with schizophrenia[45]. These authors did not explore paranoia specifically.

486

487   We assessed paranoia across the continuum of health and mental illness, provided three choice options, and

488   explicitly manipulated unexpected volatility across task versions. The version that shifted from an easier to a

489   more difficult contingency context (version 3) was associated with paranoia group effects on $\mu_3^0$, $\kappa$, and $\omega_2$, and

490   a meta-analytic effect on $\omega_3$. Furthermore, this contingency transition – an exposure to truly unexpected

491   volatility – rendered low paranoia controls more similar to their paranoid counterparts by decreasing their meta-

492   volatility learning ($\omega_3$). Paranoid participants responded to contingency transitions in version 3 and version 4 by

493   switching stochastically. These findings suggest a continuum of behavioral responses to volatility, moving from

494   optimal learning to diminished feedback sensitivity (i.e, decreased $\omega_3$ in low paranoia participants) and from

495   diminished feedback sensitivity (lower $\omega_3$ and increased win-switching in high paranoia participants) toward

496   complete dissociation from experienced feedback (stochastic switching).

497

498   Unexpected uncertainty, the perception of change in the probabilities of the environment — particularly

499   "unsignaled context switches"[17] which increase unexpected volatility — is thought to promote abandonment of

500   old associations and new learning. However, our results suggest that this response might vary according to a

501   hierarchy of belief. Paranoid participants were quick to abandon "best deck" associations and explore

502   alternative options (i.e., $x_2$ beliefs), but in turn they relied more on their higher-level beliefs about the task

503   volatility ($x_3$ beliefs) and less on sensory feedback (lower metavolatility learning). Our analysis of covariates

504   warrants specific focus on $\kappa$, the sensitivity to unexpected volatility. Other parameter-paranoia associations did

505  not endure after controlling for demographic factors (age, gender, ethnicity, and race), although we see their

506  derangement in our rodent study as well as in the significant meta-analytic effects across our experiments.

507  Furthermore, these demographic factors are themselves strong predictors of paranoia[46-48]. It is notable too that

508  $\kappa$ was the most powerful discriminator of the two clusters of human and animal participants. We conclude that

509  elevated $\kappa$ - belief updating tethered to unexpected volatility - is the parameter change most robustly

510  associated with paranoia. Doubling $\kappa$ in our simulations induced significantly more win-switching.

511

512  Multiple neurobiological manipulations may induce such win-switching behavior. Lesions of the mediodorsal

513  thalamus in non-human primates[49] or neurons projecting from the amygdala to orbitofrontal cortex in rats[50]

514  engender win-switching. Unexpected uncertainty, and the $\kappa$ parameter of the HGF in particular[51], are thought

515  to be signaled via the locus coeruleus and noradrenaline[17-19,52]. This mechanism is thought to modulate

516  switching versus staying behaviors[53-56], as well as responses to stress[57-59] and subliminal fear cues[60] to

517  coordinate fight-or-flight responses[59]. The dual role of the locus coeruleus in recognizing and responding to

518  threats as well as unexpected uncertainty suggests that dysfunction could produce both paranoia and the

519  inferential abnormalities we observed. Methamphetamine may induce similar dysfunction[61-63]. Acute moderate

520  doses increase pre-synaptic catecholamine release, particularly noradrenaline[64], and induce exploratory

521  locomotive effects modulated through adrenoceptors on dopamine neurons[62].

522

523  Excessive release of noradrenaline from the locus coeruleus into the anterior cingulate cortex drives feedback

524  insensitivity and stochastic switching behavior in rats completing a three-option counter prediction task[52].

525  Evolutionarily, departure from predictable, rational actions might offer an adaptive mechanism for escape from

526  intractable threat. As a protean defense mechanism, behavioral stochasticity impedes predators' abilities to

527  create accurate, actionable countermeasures[33,65,66]. If driven by excessive unexpected uncertainty,

528  underwritten by noradrenaline, protean defense may represent a heavily conserved, continuous common

529  mechanism underlying vigilance and false alarms[67-69], arousal-linked attentional biases[56] and selective

530  processing of social threats. However, protean behaviors are not necessarily adaptive. Pathological

531  insensitivity to feedback and reliance on internal beliefs over evidence constitute a "break from reality" – in

532  other words, psychosis.

533

534  Efference copy models of motor control[70] have been evoked to explain psychotic symptoms[71-78]. Aberrant

535  mismatches between expected and experienced sensory consequences of actions, weighted  by their

536  uncertainty[70], can lead to the misattribution of one's movements to an external agent[71-78]. Since we model

537  others' intentions with reference to our model of ourselves[79],  volatile experiences of ones' body and actions

538  will lead to uncertain and ultimately more threatening inferences about others[79]. This would be entirely

539  consistent with the present observations.

540

541 When confronted with intractable unexpected uncertainty our participants rely on higher-level beliefs about the

542 task environment. When humans experience non-social volatility, (For example through threats to their sense

543 of control[80] or exposure to surprising non-social stimuli[81,82]), they appeal to the influence of powerful enemies,

544 even when those enemies' influence is not obviously linked to the volatility[83]. Our account places the locus of

545 paranoia at the level of the individual. Here, our account departs from evolutionary accounts of paranoia

546 grounded in coalitional threat[12]; persecutors are not scapegoats that increase group cohesion. Rather, when

547 paranoid, we have a ready explanation for hazards. With a well-defined persecutor in mind, a volatile world

548 may be perceived to have less randomly distributed risk[83]. However, paranoia might become a self-fulfilling

549 prophecy, engendering more volatility and negative social interactions. This aspect may be captured in our

550 task through win-switch behavior. By failing to incorporate positive feedback from the best option, paranoid

551 individuals sample sub-optimal options which delivers misleading positive feedback.

552

553 There are some important limitations to our conclusions. Compared with humans, rats are relatively asocial.

554 But they are not completely asocial. In our experiment they were housed in pairs, and, more broadly, they

555 evince social affiliative interactions with other rats[84-86]. A further limitation centers on the comparability of our

556 experimental designs. In humans our comparisons were both within (contingency transition) and between groups

557 (low versus high paranoia). In rats, the model was also mixed with some between (saline versus

558 methamphetamine) and some within-subject (pre versus post chronic treatment) comparisons. We should be

559 clear that there was no contingency context transition in the rat study. However, just as that transition made low

560 paranoia humans behave like high paranoia, chronic methamphetamine exposure made rats behave on a stable

561 contingency much like high paranoia humans - even in the absence of contingency transition. The comparable

562 results across species, despite these differences, warrant the inference that our basic, relatively asocial,

563 approach provides a robust tool for computational dissection of learning mechanisms.

564

565 Social interactions play a rich and undeniable role in paranoia, but translational, domain-general approaches

566 may ultimately facilitate biological insights into paranoia, psychosis and delusions[87,88]. Whilst we contend that

567 our task is relatively free of social features (certainly compared to others[11]), the possibility remains that the

568 elevated U-values in our participants are reflective of attempts (and perhaps failures) to predict our intentions

569 as experimenters. Indeed, this is a possibility raised previously with regards to simple conditioned behaviors in

570 experimental animals. Even during Pavlovian conditioning, animals may attempt to infer a generative model of

571 the task environment, which might, ultimately, include the experimenter arranging the contingencies[89,90]. It is

572 possible that all instances of human cognitive testing involve an element of inference by the participant with

573 regards to the intentions of the experimenter, whether or not the task at hand is explicitly social, and indeed, all

574 cognitive functions may be aimed at or modulated by such inferences[91].

575

576  In summary, a strong belief in the volatility of the world necessitates hypervigilance and a facility with change.

577  However, in paranoia, that belief in the volatility of the world is itself resistant to change, making it difficult to

578  reassure, teach, or change the minds of people who are paranoid. They remain "on guard," adhering to

579  expectations over evidence. By using a non-social task, we have shown that this paranoid style is not restricted

580  to the social domain, and that it can be modeled in relatively asocial animals. Additionally, our domain-general

581  approach reaffirms the merit of establishing expectations of a stable, predictable environment to promote

582  recovery from paranoia-associated illness[92]. We note with interest the apparent relationship between

583  conspiratorial ideation and societal crisis situations (terrorist attacks, plane crashes, natural disasters or war)

584  throughout history, with peaks around the great fire of Rome (AD 64), the industrial revolution, the beginning of

585  the cold war, 9/11, and contemporary financial crises[93]. In today's world of escalating uncertainty and volatilty –

586  particularly environmental climate change and viral pandemics – our findings suggest that the paranoid style of

587  inference may prove particularly maladaptive for coordinating collaborative solutions.

588

589  **Methods**

590

591  Experiments were conducted at Yale University and the Connecticut Mental Health Center (New Haven, CT) in

592  strict accordance with Yale University's Human Investigation Committee and Institutional Animal Care and Use

593  Committee. Informed consent was provided by all research participants.

594

595  **Experiment 1**. English-speaking participants aged 18 to 65 ($n$=34) were recruited from the greater New Haven

596  area through public fliers and mental health provider referrals. Exclusion criteria included history of cognitive or

597  neurologic disorder (e.g., dementia), intellectual impairment, or epilepsy; current substance dependence or

598  intoxication; cognition-impairing medications or doses (e.g. opiates, high dose benzodiazepines); history of

599  special education; and color blindness. Participants were classified as healthy controls ($n$=18), schizophrenia

600  spectrum patients (schizophrenia or schizoaffective disorder; $n$=8), and mood disorder patients (depression,

601  bipolar disorder, generalized anxiety disorder, post-traumatic stress disorder; $n$=8) on the basis of clinician

602  referrals and/or self-report. Participants were compensated $10 for enrolment with an additional $10 upon

603  completion. Two healthy controls were excluded from analyses due to failure to complete the questionnaires

604  and suspected substance use, respectively.

605

606  **Experiment 2**. 332 participants were recruited online via Amazon Mechanical Turk (MTurk). The study

607  advertisement was accessible to MTurk workers with a 90% or higher HIT approval rate located within the

608  United States. To discourage bot submissions and verify human participation, we required participants to

609  answer open-ended free response questions; submit unique, separate completion codes for the behavioral

610  task and questionnaires; and enter MTurk IDs into specific boxes within the questionnaires. All submissions

611  were reviewed for completion code accuracy, completeness of responses (i.e., declining no more than 30% of

612 questionnaire items), quality of free response items (e.g., length, appropriate grammar and content), and use

613 of virtual private servers (VPS) to submit multiple responses and/or conceal non-US locations (Dennis VPS

614 paper, 2018). Upon approval, workers were compensated $6. Those who scored in the top 25% on the card

615 game (reversal-learning task) earned a $2 bonus. We rejected or excluded 19 submissions that geolocation

616 services (https://www.iplocation.net/) identified as originating outside of the United States or from suspected

617 server farms, 4 submissions for failure to manually enter MTurk ID codes, and 2 submissions for insufficient

618 questionnaire completion. Submissions with grossly incorrect completion codes were rejected without further

619 review.

620

621 **Experiment 3.** Subject information, behavioral data acquisition, and behavioral analyses were described

622 previously [22]. Long Evans rats (Charles River; $n$=20) ranged from 7 to 9 weeks of age. Rats were exposed to

623 escalating doses and frequency of saline ($n$=10) or methamphetamine ($n$=10, 3 withdrawn during dosing),

624 imitating patterns of human methamphetamine users[94,95]. Prior to dosing (Pre-Rx), rats completed 26 within-

625 session reversal sessions, including up to 8 reversals per session. Post-dosing (Post-Rx), rats completed one

626 test session per week for four weeks. Computational model parameters were estimated from each session and

627 averaged across treatment conditions to yield one Pre-Rx and Post-Rx set of parameters per rat.

628

629 **Behavioral task.** Participants completed a 3-option probabilistic reversal-learning paradigm. Three decks of

630 cards were displayed on a computer monitor for 160 trials. Participants selected a deck on each trial by

631 pressing the predesignated key. We advised participants that each deck contained winning and losing cards

632 (+100 and -50 points), but in different amounts. We also stated that the best deck may change. Participants

633 were instructed to find the best deck and earn as many points as possible. Probabilities switched between

634 decks when the highest probability deck was selected in 9 out of 10 consecutive trials (performance-dependent

635 reversal). Every 40 trials the participant was provided a break, following which probabilities automatically

636 reassigned (performance-independent reversal).

637

638 In Experiment 1, the task was presented via Eprime® 2.0 software (Psychology Software Tools, Sharpsburg,

639 PA). Participants were limited to a 3-second response window, after which the trial would time out and record a

640 null response. A fixation cross appeared during variable inter-trial intervals (jittering). Task pacing remained

641 independent of response time. In block 1 (trials 1-80) the reward probabilities (contingency) of the three decks

642 were 90%, 50%, and 10% (90-50-10%). Without cue or warning (i.e. unsignaled to the participants) the

643 contingency transitioned to 80%, 40%, and 20% (80-40-20%) upon initiation of block 2 (trials 81-160).

644

645 In Experiment 2, the task was administered via web browser link from the MTurk marketplace. We changed the

646 task timing to self-paced and eliminated null trials and inter-trial jittering. A progress tracker was provided every

647 40 trials. Workers were randomly assigned to one of four task versions, using restricted block randomization to

19

648  ensure comparable numbers of high paranoia participants across task versions. Version 1 had a constant

649  contingency of 90-50-10%. Version 4 maintained a constant contingency of 80-40-20%. Version 3 replicated

650  the 90-50-10% (block 1) to 80-40-20% (block 2) context transition of Experiment 1. Version 4 presented the

651  reversed contingency transition, 80-40-20% (block 1) to 90-50-10% (block 2). We analyzed attrition rates

652  across the four versions.

653

654  **Questionnaires.** Following task completion, questionnaires were administered via the Qualtrics® survey

655  platform (Qualtrics Labs, Inc., Provo, UT). Items included demographic information (age, gender, educational

656  attainment, ethnicity, and race) and mental health questions (past or present diagnosis, medication use,

657  *Structured Clinical Interview for DSM-IV Axis II Personality Disorders* (SCID-II)[23], Beck's Anxiety Inventory

658  (BAI)[24], Beck's Depression Inventory (BDI)[25]. We removed the single suicidality question from the BDI for

659  Experiment 2. Experiment 2 included additional items: income, three cognitive reflection questions (Table 7),

660  and three free response items ('What do you think the card game was testing?', 'Did you use any particular

661  strategy or strategies? If yes, please describe', and 'Did you find yourself switching strategies over the course

662  of the game?'). We quantified trait-level paranoia using the paranoid personality subscale of the SCID-II, and

663  we included an ideas of reference item from the schizotypy subscale ('When you are out in public and see

664  people talking, do you often feel that they are talking about you?') This item, along with other SCID-II items,

665  has previously been included as a metric of paranoia in the general population[5,96]. Participants who endorsed 4

666  or more paranoid personality items (i.e., the cut-off for the top third identified in Experiment 1) were classified

667  as 'high paranoia.' Each participant's SCID-II, BAI, and BDI scores were normalized by total scale items

668  answered. Response rates were higher than 90% for all questionnaire items and scales (Table 11).

669

670  **Behavioral analysis.** We analyzed tendencies to choose alternative decks after positive feedback (win-switch)

671  and select the same deck after negative feedback (lose-stay). Win-switch rates were calculated as the number

672  of trials in which the participant switched after positive feedback divided by the number of trials in which they

673  received positive feedback. Lose-stay rates were calculated as number of trials in which a participant persisted

674  after negative feedback divided by total negative feedback trials. In Experiment 1, we excluded post-null trials

675  from these analyses. To further characterize switching behavior, we calculated U-values, a measure of choice

676  stochasticity:

677  $$U-value = -\Sigma_{i=1}^{\beta} \frac{\log(\alpha_i) \; x \; \alpha_i}{\log(\beta)} \qquad (1)$$

678  where $\beta$ is the number of possible choice options (i.e., card decks or noseports) and $\alpha$ equals the relative

679  frequency of choice option $i$[31]. To avoid any choice counterbalancing effects across reversals, choice

680  frequencies were determined by the underlying probabilities of the decks rather than their physical attributes

681  (e.g., deck position or color). Additional behavioral analyses included trials to first reversal, trials to post-

682  reversal recovery, and trials to post-reversal switch. The latter two were restricted to the first reversal in the first

683   block. Trials post-reversal were counted from the first-negative feedback trial following the true reversal event.

684   Recovery was defined as switching to the best deck and staying for at least one additional trial.

685

686   **Computational modeling**

687   **Materials.** The Hierarchical Gaussian Filter (HGF) toolbox v5.3.1 is freely available for download in the TAPAS

688   package at https://translationalneuromodeling.github.io/tapas[20,21]. We installed and ran the package in

689   MATLAB and Statistics Toolbox Release 2016a (MathWorks ®, Natick, MA).

690   **Perceptual parameter estimation.** In the human reversal-learning experiments, we estimated perceptual

691   parameters individually for the first and second halves of the task (i.e., blocks 1 and 2). Each participant's

692   choices (i.e., deck 1, 2, or 3) and outcomes (win or loss) were entered as separate column vectors with rows

693   corresponding to trials. Wins were encoded as '1', losses as '0', and choices as '1', '2', or '3'. We selected the

694   autoregressive 3-level HGF multi-arm bandit configuration for our perceptual model and paired it with the

695   softmax-mu03 decision model.

696

697   Rat reversal-learning data was entered similarly, with choices designated as '1', '2', or '3' and reward presence

698   or absence noted as '1' and '0', respectively. Perceptual parameters were estimated as a single block per

699   session and averaged across Pre-Rx or Post-Rx sessions for each subject. Since the contingency remained

700   70-30-10%, we used the default start point values of $\mu_2$ and $\mu_3$, as in block 1 estimations for the human

701   reversal-learning experiments).

702

703   **Simulations.** We performed ten simulations per participant (online version 3) to determine whether our

704   parameter estimates and model successfully captured behavioral differences between groups (e.g., win-switch

705   rates). Each simulation required the participant's actual data (i.e., the column vectors 'outcomes' and 'choices')

706   and the corresponding set of derived perceptual parameters. On each trial, a new choice was simulated

707   conditional on the actual inputs in previous trials.

708

709   To illustrate the effects of each parameter on task behavior we doubled or halved one parameter at a time, by

710   establishing a baseline set of perceptual parameters containing the average values from the low paranoia

711   participants (online version 3). We then ran 10 simulations per subject for each of the following conditions:

712   baseline, $2\kappa$, $0.5\kappa$, $2\mu_3^0$, $0.5\mu_3^0$, $2\omega_3$, $0.5\omega_3$, $2\omega_2$, $0.5\omega_2$, and the average perceptual parameters ($\kappa$, $\mu_3^0$, $\omega_3$,

713   and $\omega_2$) from Post-Rx methamphetamine rats. The $2\omega_2$ condition yielded parameters in a region where model

714   assumptions were violated (negative posterior precision error message) and was excluded from further

715   analysis. Win-shift and lose-stay rates were calculated from each simulation as follows, and then averaged for

716   each condition:

717

$$Win\text{-}switch\ rate = \frac{Number\ of\ trials\ in\ which\ choice\ switched\ after\ positive\ feedback}{Total\ positive\ feedback\ trials}$$

21

718

719

$$Lose\text{-}stay\ rate = \frac{Number\ of\ trials\ in\ which\ choice\ repeated\ after\ negative\ feedback}{Total\ negative\ feedback\ trials}$$

720

721 For each participant, we divided rates derived from each condition by the baseline rates to determine relative

722 win-switch and lose-stay rates. We compared each relative rate to the baseline condition (i.e., 1.0) with paired-

723 samples t-tests using Bonferroni-corrected p-values.

724

725 **Parameter recovery.** We performed perceptual parameter estimation (see above) on 10 simulations per

726 subject using first block data from online version 3. These simulations were generated from each subject's

727 corresponding perceptual parameters. We averaged recovered parameters across simulations and low versus

728 high paranoia (Fig. 7).

729

730 **Alternative models.** We employed a Q-learning model with separate parameter weights for positive and

731 negative prediction errors to determine whether differential weighting might contribute to paranoia group

732 effects. This model has been described previously[35]. We also evaluated whether a simpler two-level HGF

733 model might suffice to capture paranoia group differences. To sever the third level from the model, we fixed the

734 log- $\kappa$ parameter at negative infinity (i.e., by additionally setting the variance to zero), and similarly fixed the

735 values of $\mu_3$, $\omega_3$, $\omega_2$, $\varphi_3$ at the values previously assigned in the configuration file. Parameter estimation was

736 performed as described above, with a softmax decision model.

737

738 **Statistics.** Unless otherwise specified, statistical analyses and effect size calculations were performed in IBM

739 SPSS Statistics, Version 25 (IBM Corp., Armonk, NY), with an alpha of 0.05. Box-plots were created with the

740 web tool BoxPlotR[97]. Model parameters were corrected for multiple comparisons using the Benjamini

741 Hochberg (False Discovery Rate) method. Bonferroni corrected results were largely consistent (Table 4)

742

743 To compare questionnaire item means between two groups (Table 1, low versus high paranoia), we conducted

744 independent samples t-tests. To compare questionnaire item means across paranoia groups and task versions

745 (Table 2), we employed univariate analyses. Associations between characteristic frequencies and subject

746 group or task version were evaluated by Chi-Square Exact tests (two groups) or Monte Carlo tests (more than

747 2 groups). Pearson correlations established the associations between paranoia and BDI scores, BAI scores,

748 win-switch rates, and $\kappa$. We selected two-tailed p-values where applicable and assumed normality. Multiple

749 regressions were conducted with $\kappa$ estimates from the first task block (dependent variable) and paranoia, BAI,

750 and BDI scores from online version 3.

751

752 To compare HGF parameter estimates and behavioral patterns (win-switch, U-value, lose-stay) across block,

753 paranoia group (Experiment 1, Experiment 2 version 3), and/or task version (Experiment 2), we employed

754 repeated measures and split-plot ANOVAs (i.e., block designated within-subject factor, paranoia group and

755 task version as between subject). We similarly evaluated Experiment 3 parameter estimates for treatment by

756 time interactions. For Experiment 2, we performed ANCOVAs for $\mu_3^0$, $\kappa$, $\omega_2$, and $\omega_3$ to evaluate three sets of

757 covariates: (1) demographics (age, gender, ethnicity, and race); (2) mental health factors (medication usage,

758 diagnostic category, BAI score, and BDI score); (3) and metrics and correlates of global cognitive function

759 (educational attainment, income, and cognitive reflection). Unless otherwise stated, post-hoc tests were

760 conducted as least significant difference (LSD)-corrected estimated marginal means.

761

762 Meta-analyses were conducted using random effects models with the R Metafor package[98]. Mean differences

763 were assessed for low versus high paranoia groups in the in-laboratory experiment and online version 3.

764 Standardized mean differences (methamphetamine or high paranoia versus saline or low paranoia) were

765 employed to account for the differences in task design between animal and human studies.

766

767 The 2-step clustering analysis approach was selected to automatically determine optimal cluster count and

768 cluster group assignment. Clustering variables included paranoia-relevant parameter estimates ($\mu_3^0$, $\kappa$, $\omega_2$, and

769 $\omega_3$) from Experiment 1 (block 1); online, version 3 (block 1), and rats (Post-Rx) as continuous variables with a

770 Log-likelihood distance measure, maximum cluster count of 15, and Schwarz's Bayesian Criterion (BIC)

771 clustering criterion. We validated our clustering solution by sorting the data into two halves and running

772 separate cluster analyses. We also compared cluster solutions derived exclusively from rat data versus human

773 data. A Chi-Square test determined the significance of the association between cluster membership and group

774 (methamphetamine/high paranoia versus saline/low paranoia). See Fig. 10.

775

776 **Data availability**

777 Data are available on ModelDB[99] (http://modeldb.yale.edu/258631) with accession code **p2c8q74m**.

778

779 **Code availability**

780 Code for the HGF toolbox v5.3.1 is freely available at  https://translationalneuromodeling.github.io/tapas/.
781

782 **References**

783

784 1 Freeman, D. & Garety, P. A. Comments on the content of persecutory delusions: does the definition
785 need clarification? *Br J Clin Psychol* **39 ( Pt 4)**, 407-414 (2000).
786 2 Freeman, D. *et al.* Psychological investigation of the structure of paranoia in a non-clinical population.
787 *The British journal of psychiatry : the journal of mental science* **186**, 427-435,
788 doi:10.1192/bjp.186.5.427 (2005).

789  3    Freeman, D., Pugh, K., Vorontsova, N., Antley, A. & Slater, M. Testing the continuum of delusional
790       beliefs: an experimental study using virtual reality. *J Abnorm Psychol* **119**, 83-92,
791       doi:10.1037/a0017514 (2010).
792  4    Freeman, D. *et al.* Concomitants of paranoia in the general population. *Psychol Med* **41**, 923-936,
793       doi:10.1017/S0033291710001546 (2011).
794  5    Bebbington, P. E. *et al.* The structure of paranoia in the general population. *The British journal of*
795       *psychiatry : the journal of mental science* **202**, 419-427, doi:10.1192/bjp.bp.112.119032 (2013).
796  6    Freeman, D. Delusions in the nonclinical population. *Curr Psychiatry Rep* **8**, 191-204 (2006).
797  7    Hofstadter, R. The Paranoid Style in American Politics. *Harper's Magazine*, 77-86 (1964).
798  8    Freeman, D. Suspicious minds: the psychology of persecutory delusions. *Clin Psychol Rev* **27**, 425-457,
799       doi:10.1016/j.cpr.2006.10.004 (2007).
800  9    Leamon, M. H. *et al.* Methamphetamine and paranoia: the methamphetamine experience
801       questionnaire. *Am J Addict* **19**, 155-168, doi:10.1111/j.1521-0391.2009.00014.x (2010).
802  10   Raihani, N. J. & Bell, V. Conflict and cooperation in paranoia: a large-scale behavioural experiment.
803       *Psychol Med* **48**, 1523-1531, doi:10.1017/S0033291717003075 (2018).
804  11   Raihani, N. J. & Bell, V. Paranoia and the social representation of others: a large-scale game theory
805       approach. *Sci Rep* **7**, 4544, doi:10.1038/s41598-017-04805-3 (2017).
806  12   Raihani, N. J. & Bell, V. An evolutionary perspective on paranoia. *Nat Hum Behav* **3**, 114-121,
807       doi:10.1038/s41562-018-0495-0 (2019).
808  13   Fineberg, S. K., Steinfeld, M., Brewer, J. A. & Corlett, P. R. A Computational Account of Borderline
809       Personality Disorder: Impaired Predictive Learning about Self and Others Through Bodily Simulation.
810       *Front Psychiatry* **5**, 111, doi:10.3389/fpsyt.2014.00111 (2014).
811  14   Behrens, T. E., Hunt, L. T., Woolrich, M. W. & Rushworth, M. F. Associative learning of social value.
812       *Nature* **456**, 245-249, doi:10.1038/nature07538 (2008).
813  15   Cramer, R. E. *et al.* Human agency and associative learning: Pavlovian principles govern social process
814       in causal relationship detection. *Q J Exp Psychol B* **55**, 241-266, doi:10.1080/02724990143000289
815       (2002).
816  16   Heyes, C. & Pearce, J. M. Not-so-social learning strategies. *Proceedings. Biological sciences / The Royal*
817       *Society* **282**, doi:10.1098/rspb.2014.1709 (2015).
818  17   Yu, A. J. & Dayan, P. Uncertainty, neuromodulation, and attention. *Neuron* **46**, 681-692,
819       doi:10.1016/j.neuron.2005.04.026 (2005).
820  18   Payzan-LeNestour, E. & Bossaerts, P. Risk, unexpected uncertainty, and estimation uncertainty:
821       Bayesian learning in unstable settings. *PLoS Comput Biol* **7**, e1001048,
822       doi:10.1371/journal.pcbi.1001048 (2011).
823  19   Payzan-LeNestour, E., Dunne, S., Bossaerts, P. & O'Doherty, J. P. The neural representation of
824       unexpected uncertainty during value-based decision making. *Neuron* **79**, 191-201,
825       doi:10.1016/j.neuron.2013.04.037 (2013).
826  20   Mathys, C., Daunizeau, J., Friston, K. J. & Stephan, K. E. A bayesian foundation for individual learning
827       under uncertainty. *Frontiers in human neuroscience* **5**, 39, doi:10.3389/fnhum.2011.00039 (2011).
828  21   Mathys, C. D. *et al.* Uncertainty in perception and the Hierarchical Gaussian Filter. *Front Hum Neurosci*
829       **8**, 825, doi:10.3389/fnhum.2014.00825 (2014).
830  22   Groman, S. M., Rich, K. M., Smith, N. J., Lee, D. & Taylor, J. R. Chronic Exposure to Methamphetamine
831       Disrupts Reinforcement-Based Decision Making in Rats. *Neuropsychopharmacology* **43**, 770-780,
832       doi:10.1038/npp.2017.159 (2018).
833  23   Ryder, A. G., Costa, P. T. & Bagby, R. M. Evaluation of the SCID-II personality disorder traits for DSM-IV:
834       coherence, discrimination, relations with general personality traits, and functional impairment. *J Pers*
835       *Disord* **21**, 626-637, doi:10.1521/pedi.2007.21.6.626 (2007).

836  24  Beck, A. T., Epstein, N., Brown, G. & Steer, R. A. An inventory for measuring clinical anxiety:
837      psychometric properties. *J Consult Clin Psychol* **56**, 893-897 (1988).
838  25  Beck, A. T., Ward, C. H., Mendelson, M., Mock, J. & Erbaugh, J. An inventory for measuring depression.
839      *Archives of general psychiatry* **4**, 561-571 (1961).
840  26  Lawson, R. P., Mathys, C. & Rees, G. Adults with autism overestimate the volatility of the sensory
841      environment. *Nat Neurosci* **20**, 1293-1299, doi:10.1038/nn.4615 (2017).
842  27  Hochberg, Y. & Benjamini, Y. More powerful procedures for multiple significance testing. *Stat Med* **9**,
843      811-818, doi:10.1002/sim.4780090710 (1990).
844  28  Powers, A. R., Mathys, C. & Corlett, P. R. Pavlovian conditioning-induced hallucinations result from
845      overweighting of perceptual priors. *Science* **357**, 596-600, doi:10.1126/science.aan3458 (2017).
846  29  Sevgi, M., Diaconescu, A. O., Tittgemeyer, M. & Schilbach, L. Social Bayes: Using Bayesian Modeling to
847      Study Autistic Trait-Related Differences in Social Cognition. *Biological psychiatry* **80**, 112-119,
848      doi:10.1016/j.biopsych.2015.11.025 (2016).
849  30  Na, E. J. *et al.* Paranoid Ideation Without Psychosis Is Associated With Depression, Anxiety, and Suicide
850      Attempts in General Population. *J Nerv Ment Dis* **207**, 826-831, doi:10.1097/NMD.0000000000001050
851      (2019).
852  31  Kong, X., McEwan, J.S., Bizo, L.A., Foster, T.M. An Analysis of U-Value as a Measure of Variability.
853      *Psychological Rec* **67**, 581-586 (2017).
854  32  Fung, B. J., Qi, S., Hassabis, D., Daw, N. & Mobbs, D. Slow escape decisions are swayed by trait anxiety.
855      *Nat Hum Behav* **3**, 702-708, doi:10.1038/s41562-019-0595-5 (2019).
856  33  Humphries, D. A. & Driver, P. M. Protean defence by prey animals. *Oecologia* **5**, 285-302,
857      doi:10.1007/BF00815496 (1970).
858  34  Broker, F., Marshall, L., Bestmann, S. & Dayan, P. Forget-me-some: General versus special purpose
859      models in a hierarchical probabilistic task. *PLoS One* **13**, e0205974, doi:10.1371/journal.pone.0205974
860      (2018).
861  35  Lefebvre, G., Nioche, A., Bourgeois-Gironde, S. & Palminteri, S. Contrasting temporal difference and
862      opportunity cost reinforcement learning in an empirical money-emergence paradigm. *Proc Natl Acad*
863      *Sci U S A* **115**, E11446-E11454, doi:10.1073/pnas.1813197115 (2018).
864  36  Palminteri, S., Wyart, V. & Koechlin, E. The Importance of Falsification in Computational Cognitive
865      Modeling. *Trends Cogn Sci* **21**, 425-433, doi:10.1016/j.tics.2017.03.011 (2017).
866  37  Hampton, A. N., Bossaerts, P. & O'Doherty, J. P. The role of the ventromedial prefrontal cortex in
867      abstract state-based inference during decision making in humans. *J Neurosci* **26**, 8360-8367,
868      doi:10.1523/JNEUROSCI.1010-06.2006 (2006).
869  38  Gershman, S. J., Radulescu, A., Norman, K. A. & Niv, Y. Statistical computations underlying the
870      dynamics of memory updating. *PLoS Comput Biol* **10**, e1003939, doi:10.1371/journal.pcbi.1003939
871      (2014).
872  39  Piray, P., Daw, N.D. A simple model for learning in volatile environments. *Bioarxiv* (2020).
873  40  Tkaczynski, A. in *Segmentation in Social Marketing*  (ed Rundle-Thiele S. Dietrich T., Kubacki K)  (2017).
874  41  Izquierdo, A., Brigman, J. L., Radke, A. K., Rudebeck, P. H. & Holmes, A. The neural basis of reversal
875      learning: An updated perspective. *Neuroscience* **345**, 12-26, doi:10.1016/j.neuroscience.2016.03.021
876      (2017).
877  42  Waltz, J. A. The neural underpinnings of cognitive flexibility and their disruption in psychotic illness.
878      *Neuroscience* **345**, 203-217, doi:10.1016/j.neuroscience.2016.06.005 (2017).
879  43  Robinson, O. J., Cools, R., Carlisi, C. O., Sahakian, B. J. & Drevets, W. C. Ventral striatum response during
880      reward and punishment reversal learning in unmedicated major depressive disorder. *Am J Psychiatry*
881      **169**, 152-159, doi:10.1176/appi.ajp.2011.11010137 (2012).

882    44    Dickstein, D. P. *et al.* Impaired probabilistic reversal learning in youths with mood and anxiety
883          disorders. *Psychol Med* **40**, 1089-1100, doi:10.1017/S0033291709991462 (2010).
884    45    Deserno, L. Overestimating environmental volatility increases switching behavior and is linked to
885          activation of dorsolateral prefrontal cortex in schizophrenia. *Bioarxiv* (2018).
886    46    Holt, A. E. & Albert, M. L. Cognitive neuroscience of delusions in aging. *Neuropsychiatr Dis Treat* **2**, 181-
887          189, doi:10.2147/nedt.2006.2.2.181 (2006).
888    47    Iacovino, J. M., Jackson, J. J. & Oltmanns, T. F. The relative impact of socioeconomic status and
889          childhood trauma on Black-White differences in paranoid personality disorder symptoms. *J Abnorm
890          Psychol* **123**, 225-230, doi:10.1037/a0035258 (2014).
891    48    Mahoney, J. J., 3rd, Hawkins, R. Y., De La Garza, R., 2nd, Kalechstein, A. D. & Newton, T. F. Relationship
892          between gender and psychotic symptoms in cocaine-dependent and methamphetamine-dependent
893          participants. *Gend Med* **7**, 414-421, doi:10.1016/j.genm.2010.09.003 (2010).
894    49    Chakraborty, S., Kolling, N., Walton, M. E. & Mitchell, A. S. Critical role for the mediodorsal thalamus in
895          permitting rapid reward-guided updating in stochastic reward environments. *Elife* **5**,
896          doi:10.7554/eLife.13588 (2016).
897    50    Groman, S. M. *et al.* Orbitofrontal Circuits Control Multiple Reinforcement-Learning Processes. *Neuron*
898          **103**, 734-746 e733, doi:10.1016/j.neuron.2019.05.042 (2019).
899    51    Marshall, L. *et al.* Pharmacological Fingerprints of Contextual Uncertainty. *PLoS Biol* **14**, e1002575,
900          doi:10.1371/journal.pbio.1002575 (2016).
901    52    Tervo, D. G. *et al.* Behavioral variability through stochastic choice and its gating by anterior cingulate
902          cortex. *Cell* **159**, 21-32, doi:10.1016/j.cell.2014.08.037 (2014).
903    53    Kane, G. A. *et al.* Increased locus coeruleus tonic activity causes disengagement from a patch-foraging
904          task. *Cogn Affect Behav Neurosci* **17**, 1073-1083, doi:10.3758/s13415-017-0531-y (2017).
905    54    Aston-Jones, G. & Cohen, J. D. An integrative theory of locus coeruleus-norepinephrine function:
906          adaptive gain and optimal performance. *Annu Rev Neurosci* **28**, 403-450,
907          doi:10.1146/annurev.neuro.28.061604.135709 (2005).
908    55    Aston-Jones, G., Rajkowski, J. & Cohen, J. Role of locus coeruleus in attention and behavioral flexibility.
909          *Biol Psychiatry* **46**, 1309-1320 (1999).
910    56    Eldar, E., Cohen, J. D. & Niv, Y. The effects of neural gain on attention and learning. *Nat Neurosci* **16**,
911          1146-1153, doi:10.1038/nn.3428 (2013).
912    57    Borodovitsyna, O., Flamini, M. D. & Chandler, D. J. Acute Stress Persistently Alters Locus Coeruleus
913          Function and Anxiety-like Behavior in Adolescent Rats. *Neuroscience* **373**, 7-19,
914          doi:10.1016/j.neuroscience.2018.01.020 (2018).
915    58    McCall, J. G. *et al.* CRH Engagement of the Locus Coeruleus Noradrenergic System Mediates Stress-
916          Induced Anxiety. *Neuron* **87**, 605-620, doi:10.1016/j.neuron.2015.07.002 (2015).
917    59    Atzori, M. *et al.* Locus Ceruleus Norepinephrine Release: A Central Regulator of CNS Spatio-Temporal
918          Activation? *Front Synaptic Neurosci* **8**, 25, doi:10.3389/fnsyn.2016.00025 (2016).
919    60    Liddell, B. J. *et al.* A direct brainstem-amygdala-cortical 'alarm' system for subliminal signals of fear.
920          *Neuroimage* **24**, 235-243, doi:10.1016/j.neuroimage.2004.08.016 (2005).
921    61    Ferrucci, M. *et al.* The Effects of Amphetamine and Methamphetamine on the Release of
922          Norepinephrine, Dopamine and Acetylcholine From the Brainstem Reticular Formation. *Front
923          Neuroanat* **13**, 48, doi:10.3389/fnana.2019.00048 (2019).
924    62    Ferrucci, M., Giorgi, F. S., Bartalucci, A., Busceti, C. L. & Fornai, F. The effects of locus coeruleus and
925          norepinephrine in methamphetamine toxicity. *Curr Neuropharmacol* **11**, 80-94,
926          doi:10.2174/157015913804999522 (2013).
927    63    Ferrucci, M., Pasquali, L., Paparelli, A., Ruggieri, S. & Fornai, F. Pathways of methamphetamine toxicity.
928          *Ann N Y Acad Sci* **1139**, 177-185, doi:10.1196/annals.1432.013 (2008).

929  64    Rothman, R. B. *et al.* Amphetamine-type central nervous system stimulants release norepinephrine
930        more potently than they release dopamine and serotonin. *Synapse* **39**, 32-41, doi:10.1002/1098-
931        2396(20010101)39:1<32::AID-SYN5>3.0.CO;2-3 (2001).
932  65    Richardson, G., Dickinson, P., Burman, O. H. P. & Pike, T. W. Unpredictable movement as an anti-
933        predator strategy. *Proc Biol Sci* **285**, doi:10.1098/rspb.2018.1112 (2018).
934  66    Humphries, D. A. & Driver, P. M. Erratic display as a device against predators. *Science* **156**, 1767-1768
935        (1967).
936  67    Aston-Jones, G., Rajkowski, J., Kubiak, P. & Alexinsky, T. Locus coeruleus neurons in monkey are
937        selectively activated by attended cues in a vigilance task. *J Neurosci* **14**, 4467-4480 (1994).
938  68    Rajkowski, J., Kubiak, P. & Aston-Jones, G. Locus coeruleus activity in monkey: phasic and tonic changes
939        are associated with altered vigilance. *Brain Res Bull* **35**, 607-616 (1994).
940  69    Usher, M., Cohen, J. D., Servan-Schreiber, D., Rajkowski, J. & Aston-Jones, G. The role of locus
941        coeruleus in the regulation of cognitive performance. *Science* **283**, 549-554 (1999).
942  70    Wolpert, D. M. & Ghahramani, Z. Computational principles of movement neuroscience. *Nat Neurosci* **3
943        Suppl**, 1212-1217, doi:10.1038/81497 (2000).
944  71    Blakemore, S. J., Wolpert, D. & Frith, C. Why can't you tickle yourself? *Neuroreport* **11**, R11-16 (2000).
945  72    Blakemore, S. J., Wolpert, D. M. & Frith, C. D. Central cancellation of self-produced tickle sensation. *Nat
946        Neurosci* **1**, 635-640, doi:10.1038/2870 (1998).
947  73    Blakemore, S. J., Wolpert, D. M. & Frith, C. D. The cerebellum contributes to somatosensory cortical
948        activity during self-produced tactile stimulation. *Neuroimage* **10**, 448-459,
949        doi:10.1006/nimg.1999.0478 (1999).
950  74    Blakemore, S. J., Wolpert, D. M. & Frith, C. D. Abnormalities in the awareness of action. *Trends Cogn Sci*
951        **6**, 237-242 (2002).
952  75    Frith, C. D., Blakemore, S. & Wolpert, D. M. Explaining the symptoms of schizophrenia: abnormalities in
953        the awareness of action. *Brain Res Brain Res Rev* **31**, 357-363 (2000).
954  76    Frith, C. D., Blakemore, S. J. & Wolpert, D. M. Abnormalities in the awareness and control of action.
955        *Philos Trans R Soc Lond B Biol Sci* **355**, 1771-1788, doi:10.1098/rstb.2000.0734 (2000).
956  77    Shergill, S. S., Samson, G., Bays, P. M., Frith, C. D. & Wolpert, D. M. Evidence for sensory prediction
957        deficits in schizophrenia. *Am J Psychiatry* **162**, 2384-2386, doi:10.1176/appi.ajp.162.12.2384 (2005).
958  78    Shergill, S. S. *et al.* Functional magnetic resonance imaging of impaired sensory prediction in
959        schizophrenia. *JAMA psychiatry* **71**, 28-35, doi:10.1001/jamapsychiatry.2013.2974 (2014).
960  79    Friston, K. & Frith, C. A Duet for one. *Conscious Cogn* **36**, 390-405, doi:10.1016/j.concog.2014.12.003
961        (2015).
962  80    Whitson, J. A. & Galinsky, A. D. Lacking control increases illusory pattern perception. *Science* **322**, 115-
963        117, doi:10.1126/science.1159845 (2008).
964  81    Proulx, T., Inzlicht, M. & Harmon-Jones, E. Understanding all inconsistency compensation as a palliative
965        response to violated expectations. *Trends Cogn Sci* **16**, 285-291, doi:10.1016/j.tics.2012.04.002 (2012).
966  82    Heine, S. J., Proulx, T. & Vohs, K. D. The meaning maintenance model: on the coherence of social
967        motivations. *Pers Soc Psychol Rev* **10**, 88-110, doi:10.1207/s15327957pspr1002_1 (2006).
968  83    Sullivan, D., Landau, M. J. & Rothschild, Z. K. An existential function of enemyship: evidence that
969        people attribute influence to personal and political enemies to compensate for threats to control.
970        *Journal of personality and social psychology* **98**, 434-449, doi:10.1037/a0017457 (2010).
971  84    Donaldson, T. N. *et al.* Social Order: Using The Sequential Structure of Social Interaction to Discriminate
972        Abnormal Social Behavior in the Rat. *Learn Motiv* **61**, 41-51, doi:10.1016/j.lmot.2017.03.003 (2018).
973  85    Kondrakiewicz, K., Kostecki, M., Szadzinska, W. & Knapska, E. Ecological validity of social interaction
974        tests in rats and mice. *Genes Brain Behav* **18**, e12525, doi:10.1111/gbb.12525 (2019).

86  Urbach, Y. K., Bode, F. J., Nguyen, H. P., Riess, O. & von Horsten, S. Neurobehavioral tests in rat models of degenerative brain diseases. *Methods Mol Biol* **597**, 333-356, doi:10.1007/978-1-60327-389-3_24 (2010).

87  Corlett, P. R., Fletcher, P.C. Computational Psychiatry: A Rosetta Stone linking the brain to mental illness. *Lancet Psychiatry* (2014).

88  Feeney, E. J., Groman, S. M., Taylor, J. R. & Corlett, P. R. Explaining Delusions: Reducing Uncertainty Through Basic and Computational Neuroscience. *Schizophr Bull*, doi:10.1093/schbul/sbw194 (2017).

89  Gershman, S. J. & Niv, Y. Exploring a latent cause theory of classical conditioning. *Learn Behav* **40**, 255-268, doi:10.3758/s13420-012-0080-8 (2012).

90  Gershman, S. J. & Niv, Y. Learning latent structure: carving nature at its joints. *Curr Opin Neurobiol* **20**, 251-256, doi:10.1016/j.conb.2010.02.008 (2010).

91  Turner, J. C., Oakes, P.J., Haslam, S.A., McGarty, C. Self and Collective: Cognition and Social Context. *Personality and Social Psychology B* **20**, 454-463 (1994).

92  Powers, A. R., 3rd, Bien, C. & Corlett, P. R. Aligning Computational Psychiatry With the Hearing Voices Movement: Hearing Their Voices. *JAMA psychiatry* **75**, 640-641, doi:10.1001/jamapsychiatry.2018.0509 (2018).

93  van Prooijen, J. W. & Douglas, K. M. Conspiracy theories as part of history: The role of societal crisis situations. *Mem Stud* **10**, 323-333, doi:10.1177/1750698017701615 (2017).

94  Segal, D. S., Kuczenski, R., O'Neil, M. L., Melega, W. P. & Cho, A. K. Escalating dose methamphetamine pretreatment alters the behavioral and neurochemical profiles associated with exposure to a high-dose methamphetamine binge. *Neuropsychopharmacology* **28**, 1730-1740, doi:10.1038/sj.npp.1300247 (2003).

95  Han, E., Paulus, M. P., Wittmann, M., Chung, H. & Song, J. M. Hair analysis and self-report of methamphetamine use by methamphetamine dependent individuals. *J Chromatogr B Analyt Technol Biomed Life Sci* **879**, 541-547, doi:10.1016/j.jchromb.2011.01.002 (2011).

96  Bell, V. & O'Driscoll, C. The network structure of paranoia in the general population. *Soc Psychiatry Psychiatr Epidemiol* **53**, 737-744, doi:10.1007/s00127-018-1487-0 (2018).

97  Spitzer, M., Wildenhain, J., Rappsilber, J. & Tyers, M. BoxPlotR: a web tool for generation of box plots. *Nat Methods* **11**, 121-122, doi:10.1038/nmeth.2811 (2014).

98  Viechtbauer, W. Conducting meta-analyses in R with the metafor package. *Journal of statistical software* **36** (2010).

99  McDougal, R. A. *et al.* Twenty years of ModelDB and beyond: building essential modeling tools for the future of neuroscience. *J Comput Neurosci* **42**, 1-10, doi:10.1007/s10827-016-0623-7 (2017).

**Acknowledgements**

**Competing interests**

The authors declare no competing interests.

032 **Figure Legends**

033 **Fig. 1. Probabilistic reversal learning task. a**, Human paradigm: participants choose between three decks of cards with
034 different colored backs (Blue, Red, and Green) with different, unknown probabilities of reward and loss. **b,** Reward
035 contingency schedule for in laboratory experiment (Reward probabilities associated with the different colored decks,
036 Blue, Red, Green, across trials and blocks). On trial 81, the probability context shifted from 90%, 50%, and 10% (dark
037 grey) to 80%, 40%, and 20% without warning (light grey). **c**, Reward contingency schedules for online experiment. **d,** Rat
038 paradigm: subjects choose between three noseports (A - Blue, B – Red, C- Green, for illustrative puposes) with different
039 probabilities of sucrose pellet reward. **e,** Reward contingency schedule for rat experiment[22] (Probabilities of reward
040 associated with the different noseports, A - Blue, B – Red, C- Green). Performance dependent reversals occur after a
041 certain number of choices of the high reward deck. Performance independent reversals occur regardless of participant
042 behavior.

043

044 **Fig. 2. Hierarchical Gaussian Filter (HGF) model parameters. a,** 3-level HGF perceptual model (blue) with a softmax
045 decision model (green). **Level 1 ($x_1$):** trial-by-trial perception of win or loss feedback. **Level 2 ($x_2$):** stimulus-outcome
046 associations (i.e., deck values). **Level 3 ($x_3$):** perception of the overall reward contingency context. The impact of phasic
047 volatility upon $x_2$ is captured by $\kappa$ (i.e., coupling). Tonic volatility modulates $x_3$ and $x_2$ via $\omega_3$ and $\omega_2$, respectively. $\mu_3^0$ is the
048 initial value of the third level volatility belief. **b,** HGF model parameter estimates from each of our three studies (in
049 laboratory, online, rat - columns), $\omega_3, \mu_3^0, \kappa,$ and $\omega_2$, displayed hierarchically, in rows, in parallel with the position of the
050 particular parameter in the model depiction in **a**. Parameters replicate across high paranoia groups in the in-laboratory
051 experiment (*n*=21 low paranoia [gray], 11 high paranoia [orange]; dark bars are initial task blocks, lighter bars follow the
052 contingency transition); the analogous online task (version 3, *n*=56 low paranoia [gray], 16 high paranoia [orange]; dark
053 bars are initial task blocks, lighter bars follow the contingency transition); and rats exposed to chronic, escalating saline
054 or methamphetamine (*n*=10 per group, Pre-Rx [dark gray]; Post-Rx, *n*=10 saline [light gray], 7 methamphetamine
055 [orange]). Center lines depict medians; box limits indicate the 25th and 75th percentiles; whiskers extend 1.5 times the
056 interquartile range from the 25th and 75th percentiles, outliers are represented by dots; crosses represent sample
057 means; data points are plotted as open circles. *$P \leq 0.05$, **$P \leq 0.01$, ***$P \leq 0.001$.

058

059 **Fig. 3. Paranoia effects across task versions. a, Estimated model parameters** derived from participant choices in
060 response to the tasks. Low paranoia is shown in gray, high paranoia is shown in orange. $\mu_3^0, \kappa,$ and $\omega_2$ are shown in
061 separate panels (top, middle, and bottom panels, respectively; y-axes). X-axes depict each separate online task version
062 from Experiment 2 (version 1: Easy-Easy, version 2: Hard-Hard, version 3: Easy-Hard, version 4: Hard-Easy) . **b, Behavior.**
063 **Win-switch rate** (top)**:** paranoid participants switched between decks more frequently after positive feedback. Rates are
064 collapsed across all task versions and blocks (paranoia group effect; *n*=234 low paranoia [gray], 73 high paranoia
065 [orange]). **U-value** (bottom): a measure of choice stochasticity, calculated for low (gray) and high (orange) paranoia
066 participants and collapsed across task blocks. U-values are shown separately for each online task version (1 through 4, as
067 in part **a)**. In versions 3 and 4 only (the versions containing unsignaled contingency transitions), paranoid participants
068 showed higher U-values, suggesting increasingly stochastic switching rather than perseverative returns to a previously
069 rewarding option. Center lines show the medians; box limits indicate the 25th and 75th percentiles; whiskers extend 1.5
070 times the interquartile range from the 25th and 75th percentiles, outliers are represented by dots; crosses represent
071 sample means; data points are plotted as open circles. *P*-values correspond to estimated marginal means post-hoc
072 comparisons: *$P \leq 0.05$, **$P \leq 0.01$, ***$P \leq 0.001$.

073

074 **Fig. 4. Correlations between $\kappa$ and symptoms, with and without paranoia scores of zero.** Paranoia (SCID-II, top),
075 depression (BDI, middle), and anxiety (BAI, bottom). **a,** Among all 72 subjects from online version 3, $\kappa$ correlates with
076 paranoia (r=0.30, p=0.011, top) and depression (r=0.250, p=0.034, middle), but not anxiety (r=0.210, p=0.077, bottom).
077 **b,** Among participants who endorse at least one paranoia item (SCID-II paranoia > 0, n=39), $\kappa$ correlates with paranoia
078 (r=0.588, p=8.1E-5, top), depression (r=0.427, p=0.007, middle), and anxiety (r=0.367, p=0.021, bottom). All correlations
079 are two-tailed.

080

081 **Fig. 5. Dimensionality reduction analysis.** Principal component analysis (PCA) was performed on behavioral data to
082 explain the relationship between $\kappa$ and the rating scales - paranoia (SCID), depression (BDI) and anxiety (BAI). **a,** Scree
083 plot of PCA illustrates percent of variance for each component explained by SCID, BDI and BAI. **b,** Principal component 1
084 (PC1) plotted against $\kappa$ values. $\kappa$ correlates with PC1 (r=0.272, p=0.021).
085

086 **Fig. 6. Parameter effects on simulated task performance.** We simulated behavior from low paranoia participants (online
087 Version 3, n=54) to evaluate the effects of $\kappa$, $\mu_3^0$, $\omega_2$, and $\omega_3$ on win-shift and lose-stay rates. Estimated perceptual
088 parameters were averaged across subjects to create a single set of baseline parameters. Additional parameter sets were
089 created by doubling or halving one parameter at a time (e.g., 2 $\kappa$ or 0.5 $\kappa$), while the others were held constant (n.b., 2
090 $\omega_2$ violated model assumptions and was excluded from analysis). We also included the average parameter values of rats
091 exposed to methamphetamine (Meth). Ten simulations were run per subject for each condition (i.e., parameter set).
092 Win-shift and lose-stay rates were calculated, then averaged across simulations and subjects. Rates from each condition
093 were divided by the baseline condition rate to generate relative win-shift and lose-stay rates. We compared relative
094 rates for each condition to the baseline (relative rate of 1, depicted as the dotted line; paired t-tests, Bonferroni-
095 corrected p-values). Of note, baseline parameters were positive for $\kappa$ and $\omega_2$, and negative for $\mu_3^0$ and $\omega_3$.
096 Consequently, the doubled (2x) condition makes $\mu_3^0$ and $\omega_3$ more negative (lower). (n=54). Box-plots: center lines show
097 the medians; box limits indicate the 25th and 75th percentiles; whiskers extend 1.5 times the interquartile range from
098 the 25th and 75th percentiles, outliers are represented by dots; crosses represent sample means; data points are
099 plotted as open circles; *$P \leq 0.05$, **$P \leq 0.01$, ***$P \leq 0.001$.
100

101 **Fig. 7. Parameter recovery. a, Actual subject trajectory:** this is an example choice trajectory from one participant (top).
102 The layers correspond to the three layers of belief in the HGF model (depicted in Figure 2a). Focusing on the low-level
103 beliefs (yellow box): The purple line represents the subject's estimated first-level belief about the value of choosing deck
104 1; blue, their belief about the value of choosing deck 2; and red, their belief about the value of choosing deck 3.
105 **Simulated subject trajectory** represents the estimated beliefs from choices simulated from estimated perceptual
106 parameters from that participant (middle), and **Recovered subject trajectory** represents what happens when we re-
107 estimate beliefs from the simulated choices (bottom). Crucially, Simulated trajectories closely align with real trajectories
108 (the increases and decreased in estimated beliefs about the values of each deck [purple, blue, red lines] align with each
109 other across actual, simulated and recovered trajectories), although trial-by-trial choices (colored dots and arrow)
110 occasionally differ. Outcomes (1 or 0; black dots and arrows) remain the same. **b, Actual versus Recovered:** these data
111 represent the belief parameters estimated from the participant's responses (**Actual**) compared to those estimated from
112 the choices simulated with the participant's perceptual parameters (**Recovered**). Actual and Recovered values
113 significantly correlate for $\omega_2$ (r=0.702, p=2.52E-11) and $\kappa$ (r=0.305, p=0.011) but not $\omega_3$ (r=0.172, p=0.16) or $\mu_3^0$
114 (r=0.186, p=0.13). Box plots: gray indicates low paranoia, orange designates high paranoia; center lines depict medians;
115 box limits indicate the 25th and 75th percentiles; whiskers extend 1.5 times the interquartile range from the 25th and
116 75th percentiles, outliers are represented by dots; crosses represent sample means; data points are plotted as open
117 circles. Online version 3 dataset.
118

119 **Fig. 8. Behavioral data and simulations. a,** Plots of in laboratory and online behavioral metrics. Win-switch rate
120 (switching after positive feedback), U-value (behavioral stochasticity) and Lose-stay rate (perseverating after a loss). Low
121 paranoia participants are shown in gray, High paranoia in orange. Win-switch rates and U-values are collapsed across
122 blocks. For Lose-stay rates, darker colors are block 1 data and lighter colors are block 2 data. Behavioral switching
123 patterns replicate across in laboratory and online version 3 experiments. Perseveration after negative feedback (lose-
124 stay behavior) did not significantly differ between paranoia groups or task block. **b,** Simulated data generated from HGF
125 perceptual parameters (version 3). Win-switch rate, U-value and Lose-stay rate of the simulated data are depicted. The
126 model simulated data replicate the win-switch and U-value behavioral differences between high and low paranoia
127 participants presented in panel **a**. Like the real participants, there was no difference in lose-stay rates in the simulated
128 data. Center lines show the medians; box limits indicate the 25th and 75th percentiles; whiskers extend 1.5 times the
129 interquartile range from the 25th and 75th percentiles, outliers are represented by dots; crosses represent sample
130 means; data points are plotted as open circles.*$P \leq 0.05$, **$P \leq 0.01$, ***$P \leq 0.001$. Plots of participant behavioral metrics
131 (**a**) are presented side by side with simulated data (**b**).

31

132

133 **Fig. 9. Cluster analysis of HGF parameters.** Two-step cluster analysis of model parameters ($\omega_3$, $\mu_3^0$, $\kappa$, $\omega_2$) across rat and
134 human data sets (rat, post-Rx; in laboratory and online version 3, block 1). Automated clustering yielded an optimal two
135 clusters with good cohesion and separation (average silhouette coefficient=0.7; cluster size ratio=2.46). **a, Density plots**
136 for $\mu_3^0$, $\kappa$, $\omega_2$, and $\omega_3$ (light pink) depict cluster-specific distributions for each parameter (red). Unlike frequency
137 histograms (that depict the number of data points in bins), density plots employ smoothing to prioritize distribution
138 shape and are not restricted by bin size. Beneath each density plot, box-plots of overall median, 25[th] quartile, and
139 75[th] quartile for each parameter are aligned (pink), with cluster medians and quartiles superimposed (red).  Relative to
140 the overall distribution, Cluster 1 ($n$=35) medians are elevated for $\mu_3^0$ and $\kappa$, decreased for $\omega_2$ and $\omega_3$. Cluster 2 ($n$=86)
141 falls within each overall distribution. **b, Predictor importance** of included parameters. Consistent with the color scheme
142 in Fig 2a, Uncertainty weighting parameters ($\kappa$, $\omega_2$, $\omega_3$ ) are depicted in purple and $\mu_3^0$ the prior is in blue  **c, Distribution**
143 **of cluster identities within groups**. Black bars signify the proportion of group members assigned to Cluster 1 and gray
144 bars represent the proportion of group members assigned to Cluster 2. Cluster 1 membership is significantly associated
145 with paranoia and methamphetamine groups ($\chi^2$(1, $n$=121)=29.447, $P$=5.75E-8).
146

147

148

149

150

151

**Table 1. In laboratory vs. online version 3**

| | In laboratory | | | | Online version 3 | | | |
|---|---|---|---|---|---|---|---|---|
| | Low paranoia (n=21) | High paranoia (n=11) | Statistic | p-value | Low paranoia (n=56) | High paranoia (n=16) | Statistic | p-value |
| **Demographics** | | | | | | | | |
| Age (years) | 36.0 [3.2] | 38.9 [3.9] | -0.531 (27)† | 0.6 | 38.6 [1.6] | 32.9 [1.7] | 2.441 (41.842)† | **0.019¶** |
| Gender | | | 0.006 (1)‡ | 1§ | | | .780 (1)‡ | 0.410 |
| % Female | 71.4% | 72.7% | n/a | n/a | 50.0% | 62.5% | n/a | n/a |
| % Male | 28.6% | 27.3% | n/a | n/a | 50.0% | 37.5% | n/a | n/a |
| % Other or not specified | 0% | 0% | n/a | n/a | 0% | 0% | n/a | n/a |
| Education | | | 4.972 (6)‡ | 0.638§ | | | 5.351 (6)‡ | 0.549§ |
| % High school degree or equivalent | 19.0% | 45.5% | n/a | n/a | 16.1% | 6.3% | n/a | n/a |
| % Some college or university, no degree | 14.3% | 0% | n/a | n/a | 17.9% | 25.0% | n/a | n/a |
| % Associate degree | 9.5% | 9.1% | n/a | n/a | 12.5% | 12.5% | n/a | n/a |
| % Bachelor's degree | 23.8% | 27.3% | n/a | n/a | 35.7% | 56.3% | n/a | n/a |
| % Master's degree | 9.5% | 0% | n/a | n/a | 14.3% | 0% | n/a | n/a |
| % Doctorate or professional degree | 4.8% | 0% | n/a | n/a | 1.8% | 0% | n/a | n/a |
| % Completed some postgraduate | 0% | 0% | n/a | n/a | 1.8% | 0% | n/a | n/a |
| % Other / not specified | 19.0% | 18.2% | n/a | n/a | 0% | 0% | n/a | n/a |
| Ethnicity | | | .134 (1)‡ | 1§ | | | .117 (1)‡ | 1§ |
| % Hispanic, Latino, or Spanish origin | 23.8% | 18.2% | n/a | n/a | 8.9% | 6.3% | n/a | n/a |
| % Not of Hispanic, Latino, or Spanish origin | 76.2% | 81.8% | n/a | n/a | 91.1% | 93.8% | n/a | n/a |
| Race | | | 6.250 (4)‡ | 0.186§ | | | 5.368 (4)‡ | 0.229§ |
| % White | 61.9% | 36.4% | n/a | n/a | 85.7% | 75.0% | n/a | n/a |
| % Black or African American | 19.0% | 36.4% | n/a | n/a | 0% | 12.5% | n/a | n/a |
| % Asian | 14.3% | 9.1% | n/a | n/a | 3.6% | 6.3% | n/a | n/a |
| % American Indian or Alaska Native | 4.8% | 0% | n/a | n/a | 1.8% | 6.3% | n/a | n/a |
| % Multiracial | 0% | 0% | n/a | n/a | 3.6% | 0% | n/a | n/a |
| % Other / not specified | 0% | 18.2% | n/a | n/a | 5.4% | 0% | n/a | n/a |
| **Mental health** | | | | | | | | |
| Psychiatric diagnosis | | | 12.329 (2)‡ | **0.002§** | | | 7.850 (3)‡ | **0.039§** |
| % No psychiatric diagnosis | 71.4% | 9.1% | adj. residuals | **0.004** | 71.4% | 50.0% | adj. residuals | 0.465 |
| % Schizophrenia spectrum | 19.0% | 36.4% | adj. residuals | 0.546 | 0% | 6.3% | adj. residuals | 0.307 |
| % Mood disorder | 9.5% | 54.5% | adj. residuals | 0.020ª | 21.4% | 43.8% | adj. residuals | 0.356 |
| % Not specified | 0% | 0% | adj. residuals | n/a | 7.1% | 0% | adj. residuals | 0.751 |
| % Medicated | 23.8% | 81.8% | 9.871 (1)‡ | **0.003§** | 7.1% | 31.3% | 8.730 (2)‡ | **0.023§** |
| Beck's Anxiety Inventory | 0.27 [0.08] | 0.85 [0.17] | -3.453 (30)† | **0.002** | 0.24 [0.04] | 0.90 [0.20] | -3.303 (16.179)† | **0.004¶** |
| Beck's Depression Inventory | 0.23 [0.05] | 0.66 [0.15] | -2.67 (11.854)† | **0.021¶** | 0.25 [0.04] | 1.03 [0.19] | -3.951 (16.659)† | **0.001¶** |
| SCID Paranoia Personality Score | 0.09 [0.02] | 0.63 [0.04] | -13.476 (30)† | **2.92E-14** | 0.1 [0.02] | 0.72 [0.04] | -16.551 (70)† | **6.712E-26** |
| **Reversal learning performance** | | | | | | | | |
| Total points earned | 7061.9 [286.9] | 6290.9 [372.2] | 1.608 (30)† | 0.118 | 7533.0 [143.8] | 6503.1 [340.6] | 3.177 (70)† | **0.002** |
| Total reversals achieved | 4.8 [0.7] | 2.5 [0.8] | 2.145 (30)† | **0.04** | 6.3 [0.3] | 4.9 [0.8] | 1.758 (20.14)† | 0.094§ |
| % Achieving reversals | 90.5% | 72.7% | 1.407 (1)‡ | 0.327§ | 100% | 87.5% | 7.200 (1)‡ | **0.047§** |
| Trials to switch | 1.68 [0.22] | 1.43 [0.20] | 0.671 (24)† | 0.509 | 2.1 [0.2] | 2.6 [0.6] | -1.088 (64)† | 0.280 |
| Trials to recovery | 3.75 [0.51] | 4 [0.93] | -0.285 (21)† | 0.779 | 2.9 [0.3] | 4.9 [0.8] | -2.694 (60)† | **0.009** |
| Win-switch rate, block 1 (90-50-10) | 0.08 [0.03] | 0.24 [0.09] | -1.742 (12.379)† | 0.106¶ | 0.04 [0.01] | 0.13 [0.05] | -1.906 (15.762)† | 0.075§ |
| Win-switch rate, block 2 (80-40-20) | 0.07 [0.04] | 0.21 [0.1] | -1.601 (30)† | 0.12 | 0.02 [0.01] | 0.12 [0.05] | -2.02 (15.915)† | 0.061§ |
| Lose-stay rate, block 1 (90-50-10) | 0.19 [0.03] | 0.13 [0.06] | 0.919 (30)† | 0.365 | 0.30 [0.03] | 0.39 [0.06] | -1.425 (70)† | 0.158 |
| Lose-stay rate, block 2 (80-40-20) | 0.26 [0.05] | 0.12 [0.05] | 1.817 (30)† | 0.079 | 0.33 [0.03] | 0.37 [0.06] | -0.554 (70)† | 0.581 |
| Null trials | 8.5 [2.8] | 10.4 [3.7] | -0.391 (30)† | 0.699 | n/a | n/a | n/a | n/a |

152
153
154
155
156
157
158
159
160
161
162
163

Columns display means [standard error] or percentage of participants within the described category, test-statistics, and p-values.
†Independent samples t-test: t-value (df). Two-tailed *P*-values reported.
‡Chi square coefficient (df).
§Fisher's exact test, exact significance (2-sided).
¶Equal variances not assumed.
ªNot significant (Bonferonni correction).
††Data presented in Fig. 8; repeated measures ANOVA, paranoia group trend or effect: *F*(df), *P*; estimated marginal means and standard error.
‡‡Data presented in Fig. 2; repeated measures ANOVA, *F*(df), *P*. In laboratory: paranoia x block interactions for $\omega_2$, $\mu_3^0$; paranoia group effects for $\kappa$, $\omega_2$. Version 3: paranoia group effects reported. See Table 3 for complete ANOVA results.

**Table 2. Online experiment**

| | Version 1 | | Version 2 | | Version 3 | | Version 4 | | Version Effect | | Paranoia Effect | | Interaction | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Low paranoia (n=45) | High paranoia (n=20) | Low paranoia (n=69) | High paranoia (n=18) | Low paranoia (n=56) | High paranoia (n=16) | Low paranoia (n=64) | High paranoia (n=19) | Statistic | P-value | Statistic | P-value | Statistic | P-value |
| **Demographics** | | | | | | | | | | | | | | |
| Age (years) | 36.5 [1.5] | 35.4 [2.4] | 36.2 [1.4] | 39.5 [2.8] | 38.6 [1.6] | 32.9 [1.7] | 37.6 [1.3] | 30.7 [1.6] | 1.12 (3)†† | 0.342 | 3.20 (1)†† | 0.075 | 2.62 (3)†† | 0.051 |
| Gender | | | | | | | | | 7.29 (6)‡ | 0.238§ | 1.37 (2)‡ | 0.503§ | n/a | n/a |
| *% Female* | 44.4% | 45.0% | 47.8% | 50.0% | 50.0% | 62.5% | 57.8% | 73.7% | n/a | n/a | n/a | n/a | n/a | n/a |
| *% Male* | 55.6% | 55.0% | 50.7% | 50.0% | 50.0% | 37.5% | 42.2% | 26.3% | n/a | n/a | n/a | n/a | n/a | n/a |
| *% Other or not specified* | 0.0% | 0.0% | 1.4% | 0.0% | 0.0% | 0.0% | 0.0% | 0.0% | n/a | n/a | n/a | n/a | n/a | n/a |
| Education | | | | | | | | | 15.94 (21)‡ | 0.812§|| | 7.33 (7)‡ | 0.4§ | n/a | n/a |
| *% High school degree or equivalent* | 17.8% | 20.0% | 13.0% | 16.7% | 16.1% | 6.3% | 25.0% | 10.5% | n/a | n/a | n/a | n/a | n/a | n/a |
| *% Some college or university, no degree* | 22.2% | 30.0% | 24.6% | 22.2% | 17.9% | 25.0% | 25.0% | 26.3% | n/a | n/a | n/a | n/a | n/a | n/a |
| *% Associate degree* | 13.3% | 15.0% | 17.4% | 22.2% | 12.5% | 12.5% | 9.4% | 21.1% | n/a | n/a | n/a | n/a | n/a | n/a |
| *% Bachelor's degree* | 33.3% | 35.0% | 40.6% | 22.2% | 35.7% | 56.3% | 28.1% | 31.6% | n/a | n/a | n/a | n/a | n/a | n/a |
| *% Master's degree* | 8.9% | 0.0% | 2.9% | 0.0% | 14.3% | 0.0% | 7.8% | 10.5% | n/a | n/a | n/a | n/a | n/a | n/a |
| *% Doctorate or professional degree* | 4.4% | 0.0% | 0.0% | 5.6% | 1.8% | 0.0% | 1.6% | 0.0% | n/a | n/a | n/a | n/a | n/a | n/a |
| *% Completed some postgraduate* | 0.0% | 0.0% | 1.4% | 5.6% | 1.8% | 0.0% | 3.1% | 0.0% | n/a | n/a | n/a | n/a | n/a | n/a |
| *% Other / not specified* | 0.0% | 0.0% | 0.0% | 5.6% | 0.0% | 0.0% | 0.0% | 0.0% | n/a | n/a | n/a | n/a | n/a | n/a |
| Income | | | | | | | | | 15.0 (18)‡ | .671|| | 1.18 (6)‡ | 0.981§ | n/a | n/a |
| *Less than $20,000* | 24.4% | 25.0% | 24.6% | 33.3% | 17.9% | 37.5% | 23.4% | 15.8% | n/a | n/a | n/a | n/a | n/a | n/a |
| *$20,000 to $34,999* | 40.0% | 25.0% | 20.3% | 22.2% | 33.9% | 31.3% | 28.1% | 31.6% | n/a | n/a | n/a | n/a | n/a | n/a |
| *$35,000 to $49,999* | 15.6% | 15.0% | 18.8% | 16.7% | 12.5% | 6.3% | 18.8% | 15.8% | n/a | n/a | n/a | n/a | n/a | n/a |
| *$50,000 to $74,999* | 13.3% | 35.0% | 20.3% | 5.6% | 21.4% | 12.5% | 18.8% | 21.1% | n/a | n/a | n/a | n/a | n/a | n/a |
| *$75,000 to $99,999* | 4.4% | 0.0% | 7.2% | 11.1% | 8.9% | 6.3% | 7.8% | 15.8% | n/a | n/a | n/a | n/a | n/a | n/a |
| *Over $100,000* | 0.0% | 0.0% | 5.8% | 5.6% | 3.6% | 6.3% | 1.6% | 0.0% | n/a | n/a | n/a | n/a | n/a | n/a |
| *Not specified* | 2.2% | 0.0% | 2.9% | 5.6% | 1.8% | 0.0% | 1.6% | 0.0% | n/a | n/a | n/a | n/a | n/a | n/a |
| Cognitive Reflection | | | | | | | | | 11.92 (9)‡ | 0.223§|| | 7.00 (3)‡ | 0.071§ | n/a | n/a |
| Ethnicity | | | | | | | | | 5.16 (3)‡ | 0.157§ | 3.72 (1)‡ | 0.069§ | n/a | n/a |
| *% Hispanic, Latino, or Spanish origin* | 4.4% | 15.0% | 1.4% | 0.0% | 8.9% | 6.3% | 1.6% | 15.8% | n/a | n/a | n/a | n/a | n/a | n/a |
| *% Not of Hispanic, Latino, or Spanish origin* | 95.6% | 85.0% | 98.6% | 100.0% | 91.1% | 93.8% | 98.4% | 84.2% | n/a | n/a | n/a | n/a | n/a | n/a |
| Race | | | | | | | | | 19.56 (15)‡ | .173|| | 9.63 (5)‡ | 0.084§ | n/a | n/a |
| *% White* | 82.2% | 75.0% | 84.1% | 88.9% | 85.7% | 75.0% | 85.9% | 73.7% | n/a | n/a | n/a | n/a | n/a | n/a |
| *% Black or African American* | 6.7% | 15.0% | 5.8% | 11.1% | 0.0% | 12.5% | 4.7% | 10.5% | n/a | n/a | n/a | n/a | n/a | n/a |
| *% Asian* | 8.9% | 10.0% | 7.2% | 0.0% | 3.6% | 6.3% | 7.8% | 0.0% | n/a | n/a | n/a | n/a | n/a | n/a |
| *% American Indian or Alaska Native* | 0.0% | 0.0% | 0.0% | 0.0% | 1.8% | 6.3% | 0.0% | 0.0% | n/a | n/a | n/a | n/a | n/a | n/a |
| *% Multiracial* | 2.2% | 0.0% | 1.4% | 0.0% | 3.6% | 0.0% | 1.6% | 15.8% | n/a | n/a | n/a | n/a | n/a | n/a |
| *% Other / not specified* | 0.0% | 0.0% | 1.4% | 0.0% | 5.4% | 0.0% | 0.0% | 0.0% | n/a | n/a | n/a | n/a | n/a | n/a |
| **Mental health** | | | | | | | | | | | | | | |
| Psychiatric diagnosis | | | | | | | | | 10.78 (9)‡ | 0.292§ | 2.96 (3)‡ | 0.361§ | n/a | n/a |
| *% No psychiatric diagnosis* | 73.3% | 80.0% | 60.9% | 55.6% | 71.4% | 50.0% | 65.6% | 42.1% | n/a | n/a | n/a | n/a | n/a | n/a |
| *% Schizophrenia spectrum* | 2.2% | 0.0% | 0.0% | 0.0% | 0.0% | 6.3% | 0.0% | 0.0% | n/a | n/a | n/a | n/a | n/a | n/a |
| *% Mood disorder* | 13.3% | 15.0% | 27.5% | 22.2% | 21.4% | 43.8% | 26.6% | 31.6% | n/a | n/a | n/a | n/a | n/a | n/a |
| *% Not specified* | 11.1% | 5.0% | 11.6% | 22.2% | 7.1% | 0.0% | 7.8% | 26.3% | n/a | n/a | n/a | n/a | n/a | n/a |
| *% Medicated* | 8.9% | 10.0% | 13.0% | 22.2% | 7.1% | 31.3% | 14.1% | 10.5% | 4.164 (2)‡ | 0.121§ | n/a | n/a | n/a | n/a |
| Beck's Anxiety Inventory | 0.34 [0.06] | 0.52 [0.14] | 0.31 [0.04] | 0.6 [0.13] | 0.24 [0.04] | 0.90 [0.20] | 0.33 [0.06] | 0.79 [0.18] | 1.24 (3)†† | 0.294 | 38.75 (1)†† | **1.63E-09** | 2.58 (3)†† | 0.054 |
| Beck's Depression Inventory | 0.36 [0.07] | 0.86 [0.15] | 0.32 [0.05] | 0.79 [0.13] | 0.25 [0.04] | 1.03 [0.19] | 0.38 [0.07] | 1.06 [0.20] | 1.02 (3)†† | 0.382 | 74.53(1)†† | **3.62E-16** | 1.09 (3)†† | 0.354 |
| SCID Paranoia Personality Score | 0.11 [0.02] | 0.67 [0.04] | 0.11 [0.02] | 0.61 [0.03] | 0.1 [0.02] | 0.72 [0.04] | 0.11 [0.02] | 0.65 [0.03] | 1.30 (3)†† | 0.276 | 879.38 (1)†† | **4.81E-91** | 2.02 (3)†† | 0.111 |
| **Reversal learning performance** | | | | | | | | | | | | | | |
| Total points earned | 8656.7 [182.9] | 8372.5 [405.2] | 6045.7 [135.7] | 6266.7 [288.0] | 7533.0 [143.8] | 6503.1 [340.6] | 7171.1 [175.6] | 6510.5 [403.6] | 32.29 (3)†† | **4.16E-18** | 6.18 (1)†† | **0.0135** | 2.26 (3)†† | 0.082 |
| Total reversals achieved | 7.2 [0.3] | 6.5 [0.5] | 5.5 [0.3] | 5.7 [0.5] | 6.3 [0.3] | 4.9 [0.8] | 5.9 [0.3] | 4.8 [0.6] | 4.33 (3)†† | **0.005** | 5.76 (1)†† | **0.017** | 1.10 (3)†† | 0.349 |
| % Achieving reversals | 100.0% | 100.0% | 98.6% | 94.4% | 100.0% | 87.5% | 96.9% | 94.7% | 2.26 (3)‡ | 0.598§ | 4.40 (1)‡ | 0.058§ | n/a | n/a |
| Win-switch rate, block 1 (90-50-10) | 0.09 [0.03] | 0.09 [0.04] | 0.07 [0.01] | 0.11 [0.05] | 0.04 [0.01] | 0.13 [0.05] | 0.1 [0.03] | 0.21 [0.06] | 2.28 (3)†† | 0.079 | 7.12 (1)†† | **0.008** | 1.15 (3)†† | 0.329 |
| Win-switch rate, block 2 (80-40-20) | 0.05 [0.02] | 0.08 [0.03] | 0.04 [0.01] | 0.05 [0.04] | 0.02 [0.01] | 0.12 [0.05] | 0.06 [0.02] | 0.15 [0.05] | 2.07 (3)†† | 0.105 | 9.92 (1)†† | **0.002** | 1.17 (3)†† | 0.32 |
| Lose-stay rate, block 1 (90-50-10) | 0.27 [0.03] | 0.34 [0.05] | 0.37 [0.03] | 0.34 [0.04] | 0.3 [0.03] | 0.39 [0.06] | 0.32 [0.03] | 0.34 [0.06] | 0.56 (3)†† | 0.641 | 1.83 (1)†† | 0.177 | 0.75 (3)†† | 0.521 |
| Lose-stay rate, block 2 (80-40-20) | 0.28 [0.03] | 0.23 [0.05] | 0.4 [0.03] | 0.32 [0.05] | 0.33 [0.03] | 0.37 [0.06] | 0.29 [0.03] | 0.33 [0.06] | 2.47 (3)†† | 0.062 | 0.18 (1)†† | 0.674 | 0.83 (3)†† | 0.476 |
| Reaction time, block 1 | 433.6 [28.8] | 789.3 [282.7] | 548.1 [77.8] | 365.6 [26.4] | 448 [60.1] | 442.1 [59.5] | 557.2 [108.2] | 530 [130.2] | 0.79 (3)†† | 0.499 | 0.16 (1)†† | 0.689 | 1.73 (3)†† | 0.161 |
| Reaction time, block 2 | 370.7 [23.3] | 494.3 [88.6] | 465.3 [61.6] | 331.4 [22.9] | 391.7 [52.3] | 555.9 [121.2] | 385.4 [29.2] | 504.1 [82.7] | 0.39 (3)†† | 0.757 | 1.92 (1)†† | 0.167 | 1.95 (3)†† | 0.122 |
| U-value‡‡ | 0.798 [0.009] | 0.81 [0.01] | 0.868 [0.007] | 0.871 [0.01] | 0.824 [0.008] | 0.894 [0.02] | 0.837 [0.007] | 0.877 [0.01] | 13.61 (3) | **2.42E-08** | 15.28 (1) | **0.0001** | 3.44 (3) | **0.017** |
| **Model parameters‡‡** | | | | | | | | | | | | | | |
| $\omega_0$ | -0.537 [0.12] | -0.736 [0.17] | -1.04 [0.93] | -0.821 [0.18] | -0.663 [0.10] | -0.898 [0.19] | -0.912 [0.10] | -0.993 [0.18] | 2.06 (3) | 0.105 | 0.50 (1) | 0.481 | 1.01 (3) | 0.391 |
| $\mu_1^0$ | -1.001 [0.19] | -0.721 [0.29] | -0.402 [0.16] | -0.804 [0.30] | -1.089 [0.17] | -0.180 [0.32] | -0.401 [0.16] | -0.067 [0.30] | 2.32 (3) | 0.075 | 45.08 (1) | 0.108 | 2.33 (3) | 0.075 |
| $\kappa$ | 0.480 [0.010] | 0.490 [0.015] | 0.528 [0.008] | 0.503 [0.016] | 0.470 [0.009] | 0.538 [0.017] | 0.525 [0.009] | 0.543 [0.016] | 5.06 (3) | **0.002** | 3.60 (1) | 0.059 | 4.18 (3) | **0.006** |
| $\omega_\beta$ | 1.102 [0.177] | 1.017 [0.265] | 0.330 [0.143] | 0.590 [0.280] | 1.246 [0.158] | 0.252 [0.296] | 0.603 [0.148] | 0.074 [0.272] | 4.16 (3) | **0.007** | 4.44 (1) | **0.036** | 2.81 (3) | **0.04** |

Version columns display means [standard error] or percentage of participants within the described category.
††Univariate analysis, F(df).
‡Exact test, chi-square coefficient (df).
§Exact significance (2-sided).
||Monte Carlo significance (2-sided).
‡‡Data presented in Fig. 3; repeated measures ANOVA, *F*(df), *P*. Mean values collapsed across blocks.

**Table 3. ANOVA Results for HGF Parameters**

| | Block Effect [†] | | Group Effect[‡] | | Interaction Effect | |
|---|---|---|---|---|---|---|
| | Statistic[§] | p-value | Statistic[§] | p-value | Statistic[§] | p-value |
| **Experiment 1** | | | | | | |
| $\omega_3$ | 11.672 (1) | **0.002** | 1.294 (1) | 0.264 | 6.948 (1) | **0.013** |
| $\mu_3^0$ | 25.904 (1) | **1.809E-5** | 7.063 (1) | **0.012** | 5.344 (1) | **0.028** |
| $\kappa$ | 7.768 (1) | **0.009** | 7.599 (1) | **0.010** | 0.003 (1) | 0.960 |
| $\omega_2$ | 2.182 (1) | 0.150 | 4.186 (1) | **0.050** | 0.058 (1) | 0.811 |
| $\mu_2^0$ | 4.831 (1) | **0.036** | 1.261 (1) | 0.270 | 0.370 (1) | 0.547 |
| **BIC** | 0.061 (1) | 0.807 | 8.801 (1) | **0.006** | 1.7 (1) | 0.202 |
| **Experiment 2, Version 3** | | | | | | |
| $\omega_3$ | 14.932 (1) | **0.0002** | 1.128 (1) | 0.292 | 1.406 (1) | 0.240 |
| $\mu_3^0$ | 64.651 (1) | **1.54E-11** | 6.366 (1) | **0.014** | 0.003 (1) | 0.959 |
| $\kappa$ | 15.53 (1) | **0.0002** | 13.521 (1) | **0.0005** | 0.011 (1) | 0.916 |
| $\omega_2$ | 0.027 (1) | 0.869 | 8.70 (1) | **0.004** | 0.090 (1) | 0.765 |
| $\mu_2^0$ | 11.432 (1) | **0.001** | 0.030 (1) | 0.864 | 0.203 (1) | 0.653 |
| **BIC** | 1.110E-5 (1) | 0.997 | 16.336 (1) | **0.0001** | 1.678 (1) | 0.199 |
| **Experiment 3: Rats** | | | | | | |
| $\omega_3$ | 30.086 (1) | **6.2785E-5** | 4.579 (1) | **0.049** | 9.058 (1) | **0.009** |
| $\mu_3^0$ | 31.416 (1) | **5.0188E-5** | 8.454 (1) | **0.011** | 5.159 (1) | **0.038** |
| $\kappa$ | 9.132 (1) | **0.009** | 13.356 (1) | **0.002** | 2.644 (1) | 0.125 |
| $\omega_2$ | 32.192 (1) | **4.4173E-5** | 22.344 (1) | **0.0003** | 18.454 (1) | **0.001** |
| $\mu_2^0$ | 5.226 (1) | **0.037** | 0.368 (1) | 0.553 | 2.087 (1) | 0.169 |
| **BIC** | 5.052 (1) | **0.040** | 1.890 (1) | 0.189 | 0.331 (1) | 0.573 |

[†] Block refers to first versus second half in human studies, Pre-Rx vs Post-Rx in rat studies.

[‡] Group refers to low versus high paranoia in humans, saline versus methamphetamine in rats

**Table 4. Corrections for Multiple Comparisons**

| | Group Effect [†] | | | | Interaction Effect [‡] | | | |
|---|---|---|---|---|---|---|---|---|
| | Survives Bonferroni?[§] | Survives FDR? | Critical Value | Benjamini-Hochberg p-value | Survives Bonferroni?[§] | Survives FDR? | Critical Value | Benjamini-Hochberg p-value |
| **Experiment 1** | | | | | | | | |
| $\omega_3$ | N/A | N/A | 0.05 | 0.264 | No | No | 0.0125 | 0.052 |
| $\mu_3^{\ 0}$ | Yes | Yes | 0.025 | 0.024 | No | No | 0.025 | 0.056 |
| $\kappa$ | Yes | Yes | 0.0125 | 0.04 | N/A | N/A | 0.05 | 0.96 |
| $\omega_2$ | No | No | 0.0375 | 0.0667 | N/A | N/A | 0.0375 | 1.081 |
| **Experiment 2, Version 3** | | | | | | | | |
| $\omega_3$ | N/A | N/A | 0.05 | 0.292 | N/A | N/A | 0.0125 | 0.96 |
| $\mu_3^{\ 0}$ | No | Yes | 3.75E-02 | 0.0187 | N/A | N/A | 0.05 | 0.959 |
| $\kappa$ | Yes | Yes | 0.0125 | 0.002 | N/A | N/A | 0.0375 | 1.221 |
| $\omega_2$ | Yes | Yes | 0.025 | 0.008 | N/A | N/A | 0.025 | 1.53 |
| **Experiment 3: Rats** | | | | | | | | |
| $\omega_3$ | No | Yes | 5.00E-02 | 0.049 | Yes | Yes | 0.025 | 0.018 |
| $\mu_3^{\ 0}$ | Yes | Yes | 3.75E-02 | 0.0147 | No | No | 0.0375 | 0.0507 |
| $\kappa$ | Yes | Yes | 0.025 | 0.004 | N/A | N/A | 0.05 | 0.125 |
| $\omega_2$ | Yes | Yes | 0.0125 | 0.0012 | Yes | Yes | 0.0125 | 0.004 |

N/A denotes to p-values that were not significant before corrections.

[†] Low versus high paranoia in humans, saline versus methamphetamine in rats.

[‡] Group by time (i.e., first versus second half in human studies, Pre-Rx vs Post-Rx in rat studies).

[§] p-value < 0.0125

**Table 5. Experiment 2 Effects Across Block, Paranoia Group, and Task Version**

| | Block | | Group | | Version | | Block*Group* Version | | Group*Version | | Block*Group | | Block*Version | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | F (df)[†] | p | F (df)[†] | p | F (df)[†] | p | F (df)[†] | p | F (df)[†] | p | F (df)[†] | p | F (df)[†] | p |
| $\omega_3$ | 3.722 (1) | 0.055 | 0.499 (1) | 0.481 | 2.061 (3) | 0.105 | 0.415 (3) | 0.742 | 1.005 (3) | 0.391 | 0.145 (1) | 0.704 | 7.0155 (3) | **1.42E-4** |
| $\mu_3^0$ | 288.1 (1) | **1.01E-45** | 2.604 (1) | 0.108 | 2.321 (3) | 0.075 | 0.261 (3) | 0.853 | 2.329 (3) | 0.075 | 0.281 (1) | 0.597 | 0.061 (3) | 0.98 |
| $\kappa$ | 120.9 (1) | **7.65E-24** | 3.602 (1) | 0.059 | 5.06 (3) | **0.002** | 0.08 (3) | 0.971 | 4.178 (3) | **0.006** | 1.028 (1) | 0.312 | 2.559 (3) | 0.055 |
| $\omega_2$ | 35.3 (1) | **7.92E-9** | 4.435 (1) | **0.036** | 4.155 (3) | **0.007** | 0.166 (3) | 0.919 | 2.809 (3) | **0.04** | 2.387 (1) | 0.123 | 8.697 (3) | **1.5E-5** |
| $\mu_2^0$ | 71.3 (1) | **1.33E-15** | 0.242 (1) | 0.623 | 0.616 (3) | 0.605 | 1.081 (3) | 0.358 | 0.412 (3) | 0.744 | 0.057 (1) | 0.812 | 1.505 (3) | 0.213 |
| **BIC** | 56.6 (1) | **6.23E-13** | 8.073 (1) | **0.005** | 5.385 (3) | **0.001** | 0.262 (3) | 0.853 | 4.927 (3) | **0.002** | 0.451 (1) | 0.502 | 11.905 (3) | **2.19E-07** |

[†] F-statistic (degrees of freedom); split-plot ANOVA (i.e., repeated measures with two between-subjects factors).

178
179
180
181
182
183

## Table 6. Experiment 2 ANCOVAs

| Effect | | $\omega 3$ | | | $\mu 30$ | | | $\kappa$ | | | $\omega 2$ | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | df | F | p-value | df | F | p-value | df | F | p-value | df | F | p-value |
| **Demographics (age, gender, ethnicity, and race)** | | | | | | | | | | | | |
| Block | 1 | 0.328 | 0.568 | 1 | 10.835 | **0.001** | 1 | 3.425 | 0.066 | 1 | 2.711 | 0.101 |
| Block * Age | 1 | 0.659 | 0.418 | 1 | 2.035 | 0.155 | 1 | 2.195 | 0.14 | 1 | 0.212 | 0.646 |
| Block * Gender | 1 | 0.363 | 0.547 | 1 | 0.105 | 0.746 | 1 | 4.042 | **0.046** | 1 | 0.096 | 0.757 |
| Block * Ethnicity | 1 | 0.016 | 0.901 | 1 | 0.042 | 0.837 | 1 | 0.268 | 0.605 | 1 | 0.024 | 0.876 |
| Block * Race | 1 | 3.244 | 0.073 | 1 | 0.279 | 0.598 | 1 | 0.082 | 0.775 | 1 | 1.386 | 0.24 |
| Block * Paranoia Group | 1 | 0.001 | 0.969 | 1 | 0.162 | 0.687 | 1 | 0.738 | 0.391 | 1 | 1.189 | 0.277 |
| Block * Version | 3 | 7.61 | **7.25E-05** | 3 | 0.561 | 0.641 | 3 | 2.568 | 0.055 | 3 | 8.613 | **1.97E-05** |
| Block * Paranoia Group * Version | 3 | 0.451 | 0.717 | 3 | 0.135 | 0.939 | 3 | 0.119 | 0.949 | 3 | 0.1 | 0.96 |
| Age | 1 | 3.054 | 0.082 | 1 | 2.974 | 0.086 | 1 | 2.101 | 0.149 | 1 | 2.339 | 0.128 |
| Gender | 1 | 0.438 | 0.509 | 1 | 0.02 | 0.886 | 1 | 0.005 | 0.941 | 1 | 0.014 | 0.905 |
| Ethnicity | 1 | 0.029 | 0.865 | 1 | 0.059 | 0.808 | 1 | 0.087 | 0.768 | 1 | 0.221 | 0.639 |
| Race | 1 | 0.072 | 0.789 | 1 | 2.218 | 0.138 | 1 | 0.373 | 0.542 | 1 | 0.333 | 0.564 |
| Paranoia Group | 1 | 4.71E-04 | 0.983 | 1 | 0.741 | 0.39 | 1 | 1.795 | 0.182 | 1 | 3.302 | 0.071 |
| Version | 3 | 1.845 | 0.14 | 3 | 1.914 | 0.128 | 3 | 4.975 | **0.002** | 3 | 3.786 | **0.011** |
| Paranoia Group * Version | 3 | 0.935 | 0.424 | 3 | 1.911 | 0.129 | 3 | 3.599 | **0.014** | 3 | 1.919 | 0.127 |
| **Mental health factors (medication usage, diagnostic category, BAI score, and BDI score)** | | | | | | | | | | | | |
| Block | 1 | 3.333 | 0.069 | 1 | 95.753 | **3.12E-19** | 1 | 25.498 | **8.78E-07** | 1 | 8.341 | **0.004** |
| Block * BAI | 1 | 0.26 | 0.611 | 1 | 1.532 | 0.217 | 1 | 2.852 | 0.093 | 1 | 0.394 | 0.531 |
| Block * BDI | 1 | 0.009 | 0.926 | 1 | 0.208 | 0.649 | 1 | 6.55 | **0.011** | 1 | 0.597 | 0.441 |
| Block * Medication Usage | 1 | 0.027 | 0.87 | 1 | 1.288 | 0.258 | 1 | 0.691 | 0.407 | 1 | 0.871 | 0.352 |
| Block * Diagnostic Category | 1 | 1.366 | 0.244 | 1 | 1.785 | 0.183 | 1 | 0.063 | 0.803 | 1 | 0.208 | 0.649 |
| Block * Paranoia Group | 1 | 0.068 | 0.795 | 1 | 0.298 | 0.586 | 1 | 0.298 | 0.586 | 1 | 0.007 | 0.935 |
| Block * Version | 3 | 5.872 | **0.001** | 3 | 0.531 | 0.662 | 3 | 0.906 | 0.439 | 3 | 6.16 | **0.0005** |
| Block * Paranoia Group * Version | 3 | 1.024 | 0.383 | 3 | 0.869 | 0.458 | 3 | 0.266 | 0.85 | 3 | 0.095 | 0.963 |
| BAI | 1 | 1.108 | 0.294 | 1 | 0.012 | 0.913 | 1 | 0.954 | 0.33 | 1 | 0.921 | 0.338 |
| BDI | 1 | 0.037 | 0.848 | 1 | 0.574 | 0.449 | 1 | 1.343 | 0.248 | 1 | 2.372 | 0.125 |
| Medication Usage | 1 | 0.327 | 0.568 | 1 | 0.058 | 0.81 | 1 | 0.002 | 0.966 | 1 | 0.467 | 0.495 |
| Diagnostic Category | 1 | 4.252 | **0.04** | 1 | 0.004 | 0.949 | 1 | 1.443 | 0.231 | 1 | 1.743 | 0.188 |
| Paranoia Group | 1 | 0.057 | 0.811 | 1 | 0.233 | 0.63 | 1 | 1.032 | 0.311 | 1 | 1.695 | 0.194 |
| Version | 3 | 3.183 | **0.025** | 3 | 2.73 | **0.045** | 3 | 5.274 | **0.002** | 3 | 4.468 | **0.004** |
| Paranoia Group * Version | 3 | 0.311 | 0.818 | 3 | 2.307 | 0.077 | 3 | 4.556 | **0.004** | 3 | 3.397 | **0.019** |
| **Global cognitive ability (educational attainment, income, and cognitive reflection)** | | | | | | | | | | | | |
| Block | 1 | 1.19E-04 | 0.991 | 1 | 51.264 | **7.60E-12** | 1 | 28.675 | **1.83E-07** | 1 | 18.388 | **2.51E-05** |
| Block * Education | 1 | 0.603 | 0.438 | 1 | 0.001 | 0.975 | 1 | 0.033 | 0.856 | 1 | 0.258 | 0.612 |
| Block * Income | 1 | 1.211 | 0.272 | 1 | 2.874 | 0.091 | 1 | 3.483 | 0.063 | 1 | 2.421 | 0.121 |
| Block * Cognitive Reflection | 1 | 1.83 | 0.177 | 1 | 0.709 | 0.401 | 1 | 1.221 | 0.27 | 1 | 4.667 | **0.032** |

## Paranoia & Belief Updating

| | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Block * Paranoia Group | 1 | 0.005 | 0.946 | 1 | 0.359 | 0.55 | 1 | 0.263 | 0.608 | 1 | 0.885 | 0.348 |
| Block * Version | 3 | 8.861 | **1.27E-05** | 3 | 0.182 | 0.909 | 3 | 2.325 | 0.075 | 3 | 8.815 | **1.35E-05** |
| Block * Paranoia Group * Version | 3 | 0.826 | 0.48 | 3 | 0.478 | 0.698 | 3 | 0.15 | 0.929 | 3 | 0.3 | 0.825 |
| Education | 1 | 0.111 | 0.739 | 1 | 0.578 | 0.448 | 1 | 1.395 | 0.239 | 1 | 0.608 | 0.436 |
| Income | 1 | 2.763 | 0.098 | 1 | 1.382 | 0.241 | 1 | 0.055 | 0.814 | 1 | 1.035 | 0.31 |
| Cognitive Reflection | 1 | 0.164 | 0.686 | 1 | 12.807 | **0.0004** | 1 | 0.224 | 0.636 | 1 | 0.807 | 0.37 |
| Paranoia Group | 1 | 0.069 | 0.793 | 1 | 0.555 | 0.457 | 1 | 2.477 | 0.117 | 1 | 4.715 | **0.031** |
| Version | 3 | 2.104 | 0.1 | 3 | 2.55 | 0.056 | 3 | 5.53 | **0.001** | 3 | 3.799 | **0.011** |
| Paranoia Group * Version | 3 | 1.288 | 0.279 | 3 | 2.568 | 0.055 | 3 | 4.469 | **0.004** | 3 | 2.793 | **0.041** |

184
185
186
187

### Table 7. Modified Cognitive Reflection Questionnaire Items

| Item | Prompt |
|---|---|
| 1 | A folder and a paper clip cost $1.10 in total. The folder costs $1.00 more than the paper clip. How much does the paper clip cost? |
| 2 | If it takes 5 clerks 5 minutes to review 5 applications, how long would it take 100 clerks to review 100 applications? |
| 3 | In a garden, there is a cluster of weeds. Every day, the cluster doubles in size. If it takes 48 days for the cluster to cover the entire garden, how long would it take for the cluster to cover half of the garden? |

188
189
190

191

192

193

194

195
196

**Table 8. Summary of Paranoia / Methamphetamine Effects on Belief-Updating**

|  | In lab | Online | Rats |
|---|---|---|---|
| $\omega_3$ | $\downarrow^\dagger$ | $\downarrow$ | $\downarrow$ |
| $\mu_3{}^0$ | $\uparrow$ | $\uparrow^{\ddagger\S}$ | $\uparrow$ |
| $\kappa$ | $\uparrow$ | $\uparrow^\ddagger$ | $\uparrow$ |
| $\omega_2$ | $\downarrow$ | $\downarrow^{\ddagger\P}$ | $\downarrow$ |
| $\mu_2{}^0$ | - | - | - |

↑↓   Non-significant increase/decrease in high paranoia or meth, relative to low paranoia or saline

↑↓   Trend-level increase/decrease in high paranoia or meth, relative to low paranoia or saline

↑↓   Significantly higher/lower in high paranoia or meth, relative to low paranoia or saline

- -   No significant findings or trends

[†]Baseline trend; parameter decreases in second block for low but not high paranoia

[‡]Version 3 only

[§]Trend-level significance disappears with inclusion of demographic covariates

[¶] Significance reduced to trend with inclusion of demographic covariates

197

198
199

**Table 9. Simulations and behavior**

| Effect | df | Win-switch Rate F | p-value | df | U-value F | p-value | df | Lose-stay Rate F | p-value |
|---|---|---|---|---|---|---|---|---|---|
| **Experiment 1** | | | | | | | | | |
| Block | 1 | 1.465 | 0.236 | 1 | 16.999 | **0.0003** | 1 | 1.334 | 0.257 |
| Block*Paranoia Group | 1 | 0.602 | 0.444 | 1 | 2.393 | 0.132 | 1 | 2.575 | 0.119 |
| Paranoia Group | 1 | 3.579 | 0.068 | 1 | 3.312 | 0.079 | 1 | 2.283 | 0.141 |
| **Experiment 2, Version 3** | | | | | | | | | |
| Block | 1 | 0.935 | 0.337 | 1 | 10.153 | **0.002** | 1 | 0.122 | 0.728 |
| Block*Paranoia Group | 1 | 0.001 | 0.982 | 1 | 0.003 | 0.958 | 1 | 1.93 | 0.169 |
| Paranoia Group | 1 | 12.698 | **0.001** | 1 | 19.209 | **4.03E-05** | 1 | 1.095 | 0.299 |
| **Simulations**[†] | | | | | | | | | |
| Block | 1 | 0.176 | 0.676 | 1 | 3.335 | 0.072 | 1 | 5.073 | **0.027** |
| Block*Paranoia Group | 1 | 2.039 | 0.158 | 1 | 2.624 | 0.11 | 1 | 0.036 | 0.85 |
| Paranoia Group | 1 | 15.394 | **0.0002** | 1 | 13.362 | **0.0005** | 1 | 0.042 | 0.839 |

200 [†]Simulated data from experiment 2, Version 3
201

Table 10. Alternative models fail to capture paranoia group differences

| | Low Paranoia (n=56)[†] | | | High Paranoia (n=16)[†] | | | Paranoia Group Effect[‡] | | Paranoia x Block Effect[‡] | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Mean | SEM | 95% CI | Mean | SEM | 95% CI | F(df) | P | F(df) | P |
| **Q-learning with learning rates for positive and negative prediction errors** | | | | | | | | | | |
| *Positive prediction error (α+)* | | | | | | | | | | |
| 1st half | 0.463 | 0.038 | [0.388, 0.538] | 0.475 | 0.071 | [0.335, 0.616] | 0.243 (1) | 0.623 | 0.118 (1) | 0.732 |
| 2nd half | 0.476 | 0.039 | [0.398, 0.555] | 0.535 | 0.074 | [0.379, 0.672] | | | | |
| *Negative prediction error (α-)* | | | | | | | | | | |
| 1st half | 0.421 | 0.022 | [0.377, 0.464] | 0.365 | 0.041 | [0.284, 0.446] | 1.292 (1) | 0.260 | 0.320 (1) | 0.573 |
| 2nd half | 0.386 | 0.021 | [0.344, 0.427] | 0.364 | 0.039 | [0.285, 0.442] | | | | |
| *Inverse temperature ($\beta$)* | | | | | | | | | | |
| 1st half | 271 | 74.0 | [126, 416] | 147 | 133 | [-114, 408] | 1.626 (1) | 0.207 | 0.043 (1) | 0.837 |
| 2nd half | 316 | 82.3 | [155, 477] | 145 | 132 | [-114, 403] | | | | |
| **2-level HGF with softmax decision model** | | | | | | | | | | |
| $\mu_2$ | | | | | | | | | | |
| 1st half | -0.059 | 0.081 | [-0.218, 0.100] | -0.303 | 0.157 | [-0.611, 0.005] | 3.039 (1) | 0.086 | 0.385 (1) | 0.537 |
| 2nd half | -0.244 | 0.082 | [-0.405, -0.082] | -0.566 | 0.155 | [-0.869, -0.262] | | | | |
| *Inverse temperature ($\beta$)* | | | | | | | | | | |
| 1st half | 131 | 30.6 | [71.3, 191] | 35.3 | 6.20 | [23.2, 47.5] | 2.665 (1) | 0.107 | 0.250 (1) | 0.619 |
| 2nd half | 119 | 30.6 | [58.7, 179] | 52.1 | 12.1 | [28.3, 75.9] | | | | |

[†]Online version 3 data
[‡]Repeated measures ANOVA

**Table 11. Questionnaire item completion (% responses)**

| Questionnaire / subscale | Experiment 1 | Experiment 2 |
|---|---|---|
| **Age** | 90.6% | 99.7% |
| **Gender** | 100.0% | 100.0% |
| **Ethnicity** | 100.0% | 100.0% |
| **Race** | 100.0% | 100.0% |
| **Education** | 100.0% | 99.7% |
| **Meds** | 100.0% | 90.6% |
| **Dx** | 100.0% | 94.1% |
| **Income** | N/A | 98.0% |
| **SCID-II Paranoia - all items** | 96.9% | 94.1% |
| SCID-II Paranoia - 1 item missing | 3.1% | 5.5% |
| SCID-II Paranoia - 3 items missing | 0.0% | 0.3% |
| **Cognitive reflection - all items** | N/A | 97.7% |
| **Beck's Anxiety Inventory (BAI) - all items** | 90.6% | 96.7% |
| BAI - 1 item missing | 3.1% | 2.9% |
| BAI - 2 items missing | 6.3% | 0.3% |
| **Beck's Depression Inventory (BDI) - all items** | 100.0% | 99.0% |
| BDI - 1 item missing | 0.0% | 1.0% |

211
212
213
214
215
216
217
218
219
220

**a**

**Performance-independent**
— Deck A
— Deck B
— Deck C

**Performance-dependent**
- - Deck A
- - Deck B
···· Deck C

**b**

**c**

Version 1

Version 2

Version 3

Version 4

**d**

**Performance-dependent**
- - Noseport A
- - Noseport B
···· Noseport C

**e**

**a**

3-level HGF model

$\omega_3$ Metavolatility

Level 3

$X_3^{(t-1)}$ → $X_3^{(t)}$

Contingency context (Phasic volatility)

$\kappa$ Phasic volatility coupling

Perceptual model

$\omega_2$ Tonic volatility (Level 2)

Level 2

$X_2^{(t-1)}$ → $X_2^{(t)}$

Stimulus-outcome associations

Level 1

$X_1^{(t-1)}$ → $X_1^{(t)}$

Win or loss feedback

Stay or Switch

Softmax, $\beta = \exp(-\mu_3^{(t)})$

Response model

**b**

In laboratory

Version 3 (online)

Rat

Low paranoia   High paranoia

Low paranoia   High paranoia

Saline   Meth

$\omega_3$

$\mu_3^0$

$\kappa$

$\omega_2$

Low paranoia, block 1

Low paranoia, block 2

High paranoia, block 1

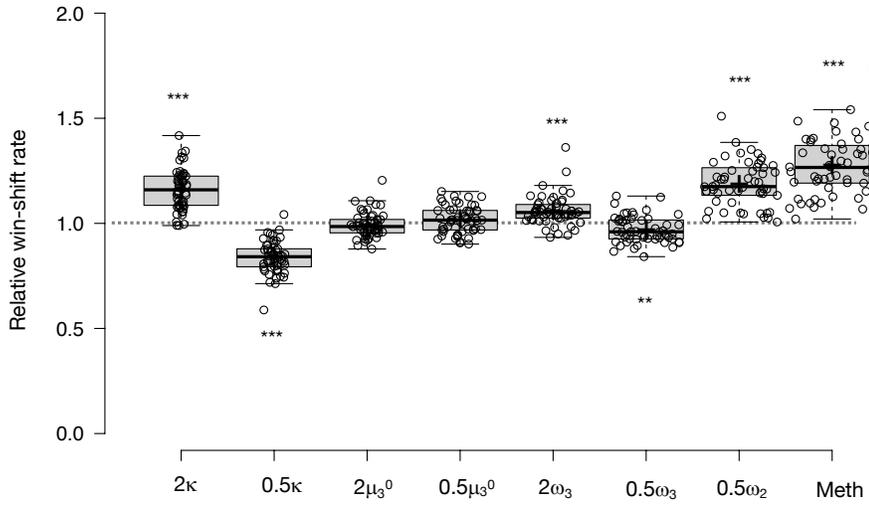High paranoia, block 2
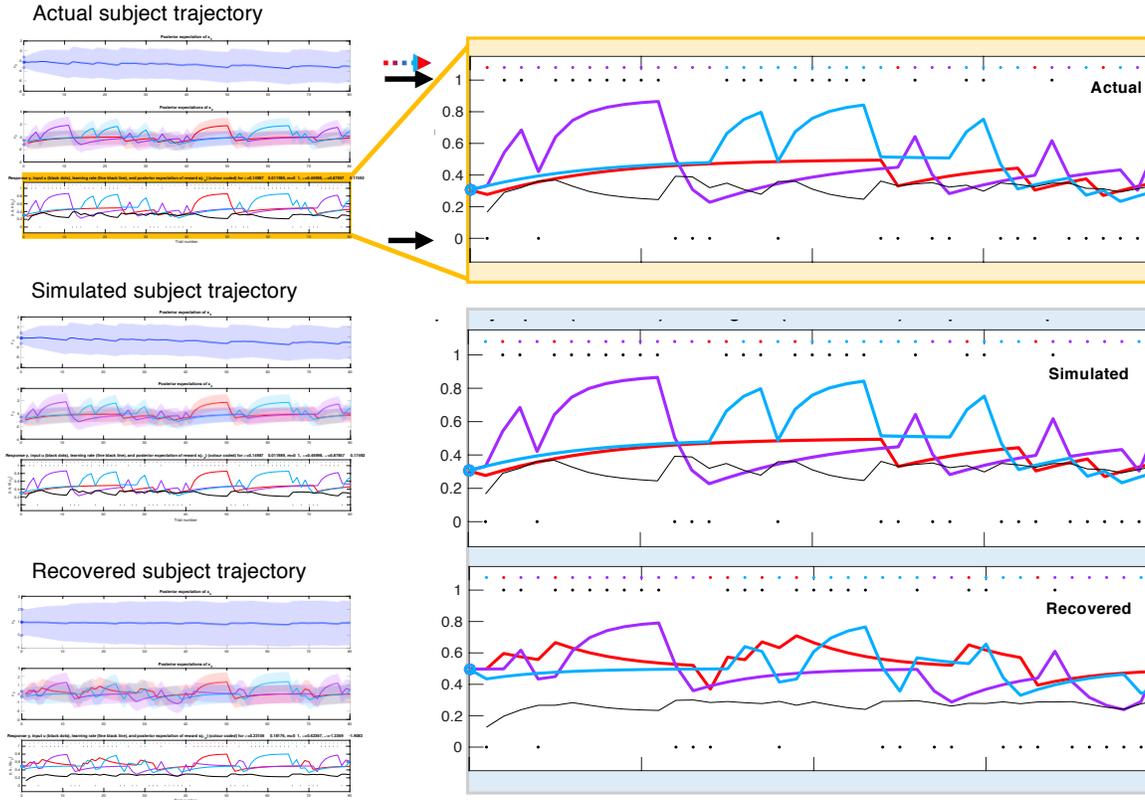
Pre-Rx

Post-Rx, saline

Post-Rx, methamphetamine

**a** Model parameters

**b** Behavior

Low paranoia

High paranoia

**a**



y = 0.1631x + 0.4323
R² = 0.3434

κ block 1 vs Paranoia

y = 0.1631x + 0.4323
R² = 0.3434

κ block 1 vs BDI

y = 0.1631x + 0.4323
R² = 0.3434

κ block 1 vs BAI

**b**

y = 0.1631x + 0.4323
R² = 0.3434

κ block 1 vs Paranoia

y = 0.1631x + 0.4323
R² = 0.3434

κ block 1 vs BDI

y = 0.1631x + 0.4323
R² = 0.3434

κ block 1 vs BAI

**a**



**b**

**a**

Actual subject trajectory

Simulated subject trajectory

Recovered subject trajectory

**b**

$\omega_3$

$\kappa$

$\mu_3{}^0$

$\omega_2$

a

**In laboratory**

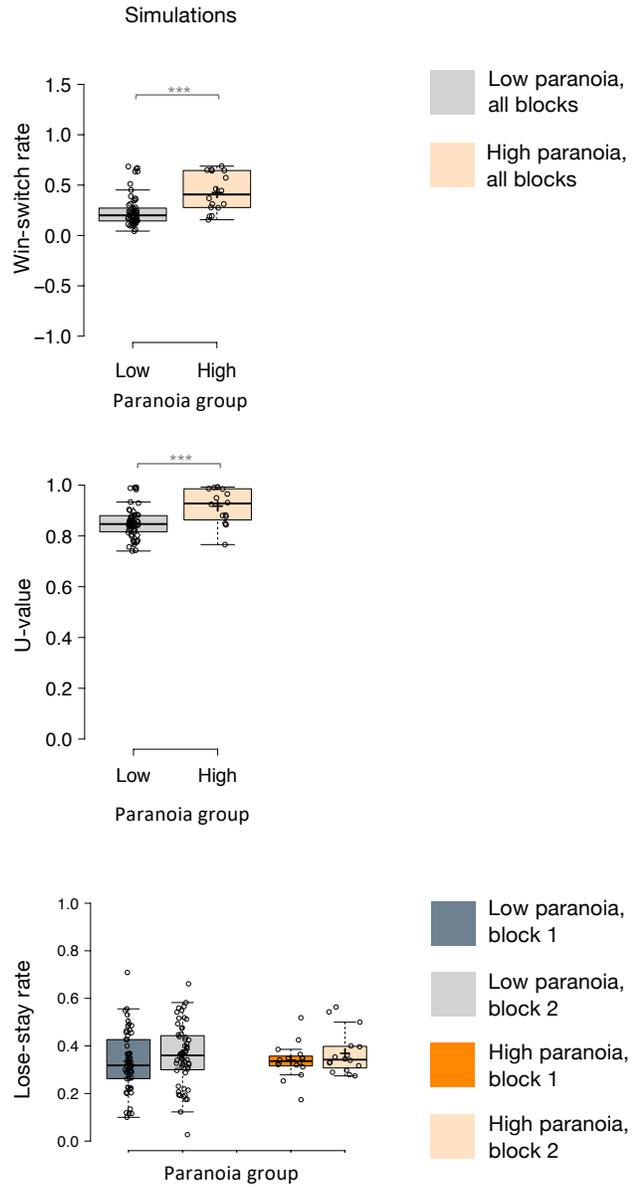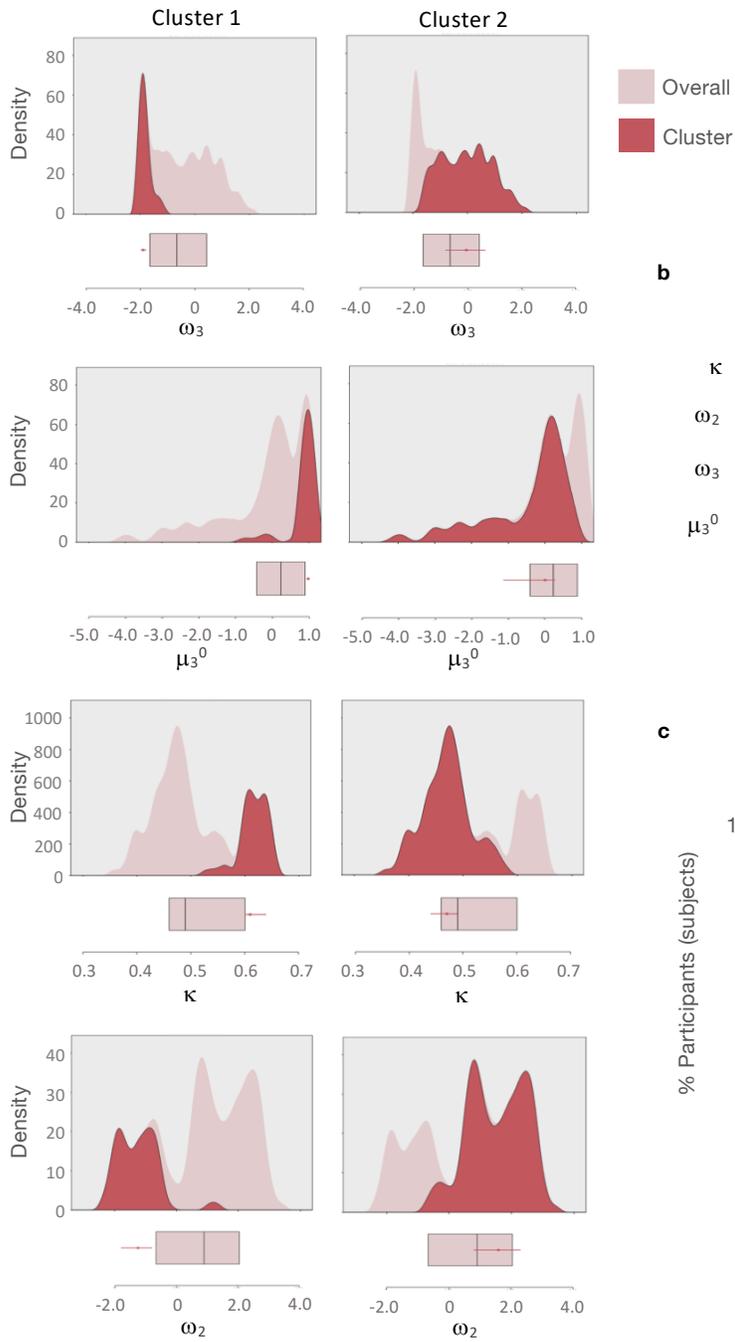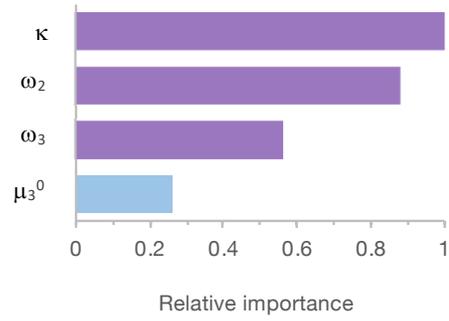**Version 3 (online)**

b

**Simulations**

**a** Cluster analysis cell distribution

**b** Predictor importance

**c** Cluster group membership