

1 **Transdiagnostic factors differentially shape choices and reaction**
2 **times in social evaluative learning**

3

4 Aamir Sohail^{*a,b,c}, Aleksandra Magda^{^a}, Muna Said^{^a}, Laura Phillips^{^a}, Yifei Lu^a,
5 Georgina Reynolds^a, Alice Comer-Fenn^a, Jessie Russell^a, Huafeng Lu^d, Tally Sims^a,
6 Anne C. Saulin^{a,b}, Lei Zhang^{*a,b,c,e}

7

8 ^aSchool of Psychology, University of Birmingham, Birmingham, B15 2TT, UK

9 ^bCentre for Human Brain Health, University of Birmingham, Birmingham, B15 2TT, UK

10 ^cInstitute for Mental Health, University of Birmingham, Birmingham, B15 2TT, UK

11 ^dSchool of Architecture, Built Environment, Computing and Engineering, Birmingham
12 City University, Birmingham, B4 7BD, UK

13 ^eCentre for Developmental Science, University of Birmingham, Birmingham, B15 2TT,
14 UK

15

16 ORCID ID

17 Aamir Sohail: 0009-0000-6584-4579

18 Anne C. Saulin: 0000-0002-7200-468X

19 Lei Zhang: 0000-0002-9586-595X

20

21 *Correspondence should be addressed to:

22 Aamir Sohail: axs2210@bham.ac.uk

23 Lei Zhang: l.zhang.13@bham.ac.uk

24 Centre for Human Brain Health, University of Birmingham, Birmingham, B15 2TT, UK

25

26 [^]These authors contributed equally to this work.

27 **Abstract**

28 **Background**

29 Learning from feedback provided by others is important for navigating complex social
30 environments. Aberrant processing of such social feedback is implicated in the
31 psychopathology of several mental health disorders including social anxiety and depression.
32 Recent psychiatric research increasingly adopts transdiagnostic approaches that examine
33 mechanisms and constructs cutting across traditional diagnostic categories.

34 **Methods**

35 Here, we applied reinforcement learning models together with transdiagnostic measures of
36 psychopathology to formally assess social evaluation learning. Human participants (n = 193)
37 completed the Social Evaluation Learning Task, learning whether computer personas 'liked,'
38 'disliked,' or were 'neutral' toward them. Participants also completed six questionnaires
39 assessing socially relevant psychiatric traits.

40 **Results**

41 Computational modeling revealed that a model with separate learning parameters for positive
42 and negative feedback best accounted for the learning behavior. Exploratory factor analysis
43 identified three transdiagnostic factors. A factor reflecting social avoidance predicted
44 enhanced learning from negative feedback and a reduced positive learning bias. Socially
45 avoidant individuals with heightened negative learning rates additionally showed faster
46 reaction times specifically in positive social contexts. A second factor linking with emotional
47 insensitivity for others predicted a lower choice accuracy in both positively and negatively
48 valanced conditions and was associated with faster reaction times overall. No significant
49 associations emerged for a third factor reflecting depressive and mood-related symptoms.

50 **Conclusions**

51 These results demonstrate how transdiagnostic trait dimensions shape social learning
52 mechanisms through specific behavioral and computational processes.

53

54 **Keywords**

55 social evaluation; social learning; reinforcement learning; transdiagnostic; computational
56 psychiatry

57 **Introduction**

58 Learning from social feedback is a vital cognitive process enabling individuals to adjust their
59 behavior based on social cues, such as praise and criticism (Bauer, 1967). This mechanism
60 is crucial for navigating complex social environments, maintaining relationships, and shaping
61 one's self-concept (Boyd et al., 2011; Button et al., 2012; Elder et al., 2022; FeldmanHall &
62 Chang, 2018; Kendal & Watson, 2023; Olsson et al., 2020; Peters et al., 2024; Schröder et
63 al., 2025). In the form of social evaluative feedback, individuals update beliefs about the self
64 and others (Behrens et al., 2008; Diaconescu et al., 2014) from external judgements or
65 feedback which can differ in valence and expectancy (Peters et al., 2024). Appropriately
66 learning from social feedback is integral to mental and physical well-being, with core self-
67 appraisals being significantly influenced by external evaluations (Lundgren, 2004).
68 Conversely, atypical social cognition is a feature of several mental health disorders (Gkika et
69 al., 2018; Green & Leitman, 2008; Hoernagl & Hofer, 2014; Patin & Hurlemann, 2015;
70 Weightman et al., 2014) where individuals report biased perceptions of oneself which persist
71 in response to social feedback (Gilboa-Schechtman et al., 2017; Hoffmann et al., 2024;
72 Kirchner et al., 2025; Koban et al., 2017; Korn et al., 2012; Kube, 2023). Although previous
73 studies have explored the computational basis of social evaluation learning and its relation to
74 various mental health disorders, to date, none have specifically taken a transdiagnostic
75 approach when doing so.

76

77 Social learning can be formally investigated using the framework of reinforcement learning
78 (RL) (Sutton & Barto, 2018), a form of error-based learning driven by differences between
79 actual and expected outcomes (i.e., prediction error). Specifically, individuals form a social
80 prediction error (Joiner et al., 2017; Zhang et al., 2020; Zhang & Gläscher, 2020) - the
81 difference between expected social feedback and the social feedback actually received.
82 Application of the RL framework has revealed distinct biases in social feedback learning in
83 several psychiatric disorders. For example, individuals with social anxiety and borderline
84 personality disorder are biased towards learning from negative social feedback (Korn et al.,
85 2016; Müller-Pinzler et al., 2019; Zabag et al., 2022, 2024) and demonstrate reduced positivity
86 regarding one's own attributes (Hoffmann et al., 2024; Hopkins et al., 2021; Korn et al., 2012,
87 2016). In social anxiety, this is specifically linked to the maintenance of a negative self-view
88 and low self-esteem (Koban et al., 2017, 2023). Similarly, in contrast to healthy participants,
89 depressed individuals do not update their beliefs more strongly following positive compared
90 to negative social feedback (Hobbs et al., 2022; Hoffmann et al., 2024). This blunted learning

91 appears to reflect the domain-general effect of anhedonia which extends to the context of
92 social interactions (Barkus & Badcock, 2019).

93

94 Psychiatric disorders may present differently within categorical diagnoses (Feczko et al.,
95 2019), reflecting a complex behavioral profile which cannot be appropriately captured by a
96 single dimension or questionnaire. Indeed, diagnostic manuals, despite being standard tools,
97 have been critiqued for their categorical approach (Adriaens & De Block, 2013; Kendler, 2009;
98 Kendler et al., 2011), leading to the development and advocacy of alternative approaches with
99 broader, multi-faceted definitions of mental health (Borsboom, 2017; Insel et al., 2010; Kotov
100 et al., 2017). Compared to categorical assessments, a transdiagnostic approach to
101 psychopathology (Gillan et al., 2016; Insel et al., 2010; Robbins et al., 2012; Tanaka, 2024)
102 favours a broader definition of mental illness, spanning normal variations in psychopathology
103 in a general population sample. Focusing on dimensional traits that cut across multiple
104 disorders, latent psychological and cognitive constructs or factors of mental health can be
105 extracted that more accurately map onto candidate cognitive processes (Gillan & Seow,
106 2020). This approach can further be combined with the theory-driven computational modeling
107 of behavior such as RL models (Sohail & Zhang, 2024), to infer the construct-specific
108 computational processes that characterize a given factor (Wise et al., 2023). Combining
109 transdiagnostic assessments of psychopathology with computational models have identified
110 the specific cognitive processes underlying altered behavior in mental health disorders,
111 including model-based planning (Gillan et al., 2016), metacognition (Rouault et al., 2018),
112 reward processing (Suzuki et al., 2021) and uncertainty (Norbury et al., 2018).

113

114 In the context of social feedback learning, existing research has primarily taken a categorical
115 approach often focusing on a single disorder or trait. However, a transdiagnostic approach
116 can uncover latent factors across trait dimensions, allowing for shared cognitive components
117 to be extracted. In the current study we adopted such a perspective, using a questionnaire
118 battery in a nonclinical sample assessing diverse domains of psychopathology known to affect
119 social behavior including autism (Chevallier et al., 2012), social anxiety (Hofmann, 2007),
120 depression (Rygula et al., 2015; Tse & Bond, 2004), apathy (Le Heron et al., 2019), narcissism
121 (Morf & Rhodewalt, 2001), and borderline personality disorder (Herpertz & Bertsch, 2014).
122 Specifically, we used exploratory factor analysis (EFA), a data-driven approach yielding
123 composite trait factors that can be related to behavioral measures and computational
124 parameters. For our behavioral task, we used the Social Evaluation Learning Task (SELT)
125 (Button et al., 2012, 2015), a well-established paradigm requiring participants to learn whether

126 a fictional computer persona liked or disliked them across three feedback conditions differing
127 by valence.

128

129 In line with previous studies extracting latent factors from similar trait measures (Gillan et al.,
130 2016; Oka et al., 2025), we anticipated transdiagnostic factors reflecting social engagement
131 and mood to be extracted from our questionnaire battery. We also hypothesised these factors
132 to differentially moderate social learning effects. Specifically, we expected for social
133 engagement to negatively correlate with a negative learning bias, and for a depressive-mood
134 factor to positively correlate with reduced learning from positive feedback. These hypotheses
135 stem from previous studies reporting these biases for categorical assessments of social
136 anxiety and depression (Button et al., 2012; Hoffmann et al., 2024; Hopkins et al., 2021). Using
137 a within-subject design across feedback conditions and continuous trait measures, this study
138 ultimately provides a nuanced look at how overlapping psychiatric traits relate to the ability to
139 learn from social feedback.

140

141 **Methods**

142 **Participants**

143 A priori sample size estimation using G*Power 3 software (Faul et al., 2007) suggested a
144 sample of 191 would be needed to detect a small effect size ($f = 0.20$) with 80% power and α
145 = .05 for a repeated-measures ANOVA. Participants were eligible for the study if they were
146 18-to-25 years old, proficient in English, and had no history of brain trauma, cognitive
147 impairment, neurological conditions or psychological disorders. Participants were recruited
148 using the University of Birmingham Research Participation Scheme (Sona Systems) and flyers
149 on campus and were tested online using the online experiment platform Gorilla
150 (www.gorilla.sc) (Anwyl-Irvine et al., 2020). Participants received one academic credit point
151 credit for completing the study. Data collection comprised 198 participants, which after
152 excluding five participants who failed attention checks, resulted in a final sample (Table 1) of
153 193 participants (female = 173; age = 19.4 ± 2.0). All participants provided informed consent
154 prior to participation, and the study was approved by the University of Birmingham Research
155 Ethics Committee (ERN_20-1897AP15).

156

157 **Table 1**

158 Participant demographics for the study

159 *Demographic characteristics of the study participants, including group sizes (n) and*
160 *percentages (%) for sex assigned at birth, gender identity, and ethnic group.*

	<i>n</i>	%
<i>Sex Assigned at birth</i>		
<i>Male</i>	20	10.4
<i>Female</i>	173	89.6
<i>Gender</i>		
<i>Male</i>	18	9.3
<i>Female</i>	173	89.6
<i>Transgender</i>	1	0.5
<i>Gender Variant</i>	1	0.5
<i>Ethnic Group</i>		
<i>White</i>	88	45.6
<i>Asian or Asian British</i>	59	30.6
<i>Black, Black British, Caribbean or African</i>	23	11.9
<i>Mixed or Multiple</i>	13	6.7
<i>Other</i>	10	5.2
Total	193	

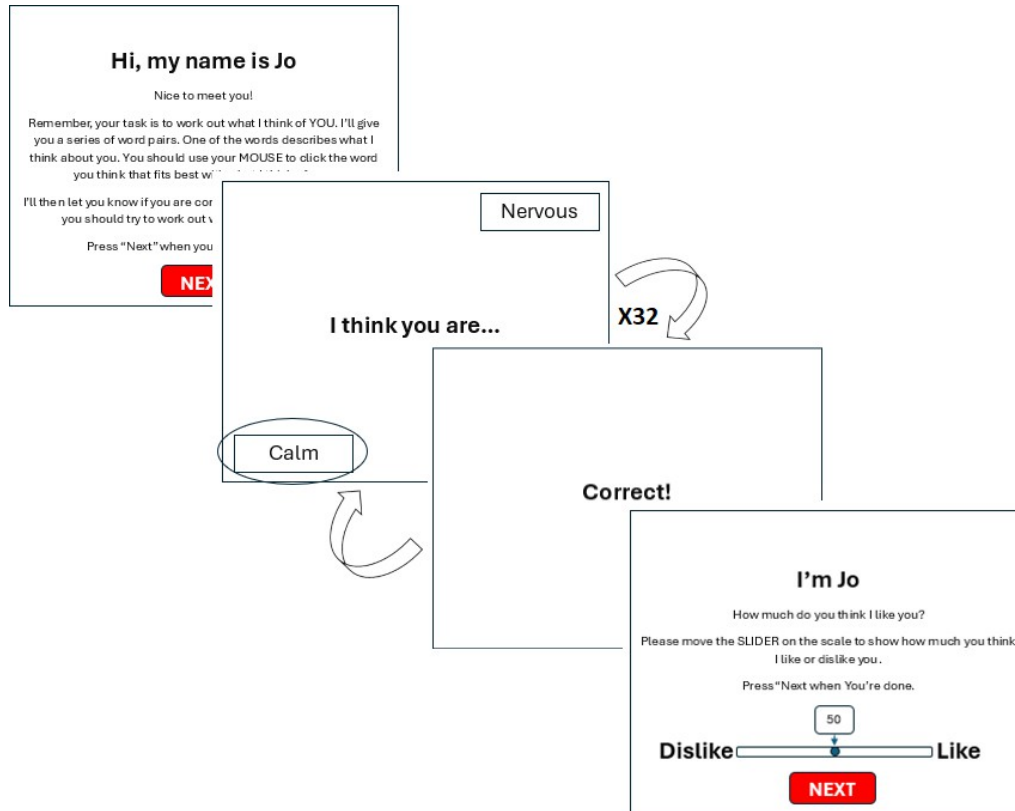
161

162 **Measures**

163 **Social Evaluation Learning Task**

164 Participants completed the Social Evaluation Learning Task (Button et al., 2012, 2015), a
165 probabilistic learning paradigm involving social feedback. In our experiment, we included three
166 self-referential test blocks, each with a computer profile (i.e., persona) that participants had to
167 learn either liked, disliked, or was neutral towards the participants (**Figure 1**). The order of the
168 three blocks was randomized and counterbalanced across participants, with each block
169 consisting of 32 trials. In each trial, participants were presented with the sentence “I think you
170 are...” and a different positive/negative word pair randomly selected from a list of 64 pairs
171 (e.g., nervous vs. calm). The placement (i.e., lower left vs. upper right) of the two words on
172 the screen was counterbalanced across trials, to avoid motor habituation. Each persona
173 followed a unique feedback contingency rule: the positive word in the word pair was the correct
174 choice 80% of the time for the ‘Like’ persona, 20% for the ‘Dislike’ persona, and 50% for the
175 ‘Neutral’ persona, whereas the negative word in the word pair was the correct choice 20% of
176 the time for the ‘Like’ persona, 80% for the ‘Dislike’ persona, and 50% for the ‘Neutral’ persona.
177 On each trial, participants were instructed to choose which word in each pair described the
178 persona’s opinion toward them, to which they received feedback according to the

179 contingencies for each persona. The feedback contingencies were not communicated to the
 180 participants; instead, participants needed to work these out using trial and error. At the end of
 181 each block, participants were asked to rate how much the persona liked them on a scale from
 182 0 ('Dislike') to 100 ('Like').
 183



184
 185 **Fig 1. The Social Evaluation Learning Task.** Participants were instructed to determine the
 186 computer persona's perspective of themselves. Participants played three blocks of 32 trials
 187 across three conditions, differing in the probability of positive/negative words being the correct
 188 choice. Feedback was provided after each trial. At the end of each block, participants were
 189 asked to provide a single rating on how much the persona disliked them using a Likert scale
 190 from 0 (Dislike) to 100 (Like).

191
 192 **Questionnaires**

193 Immediately following the behavioral task, participants completed six self-report
 194 questionnaires to assess a wide range of psychopathology previously associated with deficits
 195 in social learning, namely the Autism Spectrum Quotient (AQ-10) (Allison et al., 2012), Social
 196 Phobia Inventory (SPIN) (Connor et al., 2000), Beck Depression Inventory (BDI-II) (Beck et
 197 al., 1996), Apathy Motivation Index (AMI) (Ang et al., 2017), Narcissistic Admiration and

198 Rivalry Questionnaire (NARQ) (Back et al., 2013), and Borderline Personality Disorder (BPD)
199 checklist (First et al., 1997).

200

201 Data analysis approach

202 **Data quality measures**

203 An attention check was randomly assigned to two of the three blocks, where participants were
204 asked to select a specific word: “Please select <word> to show you are paying attention...”.

205 Participants who failed either of the two attention checks were removed from the analysis.

206 Trials with a reaction time greater than 5000ms or less than 300ms were also excluded from
207 all analyses.

208

209 **Behavioral measures**

210 Three primary dependent variables were selected as indicators of task behavior: choice
211 accuracy, reaction time, and end-of-block persona rating. Choice accuracy was defined as

212 selecting the word whose valence was associated with the feedback contingency; positive

213 words for the ‘Like’ condition, negative words for the ‘Dislike’ condition. For the neutral

214 condition, there is no actually correct response. We coded the ‘Like’ option as correct response

215 in this condition (‘Dislike’ would reveal same results) to enable joint analysis of all conditions.

216 Raw untransformed reaction times were used for repeated-measures ANOVA, as subject-

217 level mean RTs were normally distributed across all conditions (Shapiro-Wilk tests: $W > 0.985$;

218 skewness = 0.09-0.24). For trial-level mixed-effects models, we used log-transformed RT

219 values. All mixed-effects models were fitted using the {lme4} package in R (Bates et al., 2015).

220

221 **Computational modeling**

222 A reinforcement learning framework was applied to quantify the computational processes

223 underlying task behavior. Computational models included a simple Rescorla-Wagner model

224 (Rescorla & Wagner, 1972), where for each trial (t), the value (V) of the chosen option was

225 updated by the prediction error (PE) – the difference between the received outcome (O , coded

226 as -1 and 1) and the value:

227

$$228 \quad PE_t = O_t - V_t \quad (1)$$

229

230 Subsequently, the value of the next trial was adjusted by the prediction error, weighted by the

231 learning rate (α , $0 < \alpha < 1$):

232

233
$$V_{t+1} = V_t + \alpha PE_t \tag{2}$$

234

235 Considering that participants may learn differently from positive versus negative social
 236 feedback (Zhang et al., 2020), a second RL model was tested with two distinct learning rates
 237 (α^+ and α^-) for positive and negative outcomes:

238

239
$$\begin{aligned} V_{t+1} &= V_t + \alpha^+ PE_t, & O_t &> 0 \\ V_{t+1} &= V_t + \alpha^- PE_t, & O_t &< 0 \end{aligned} \tag{3}$$

240

241 For both models, a softmax choice function was employed to compute the action probability
 242 based on the action values, where during each trial, the action probability of choosing the
 243 positive word over the negative word was defined as:

244

245
$$p(Positive)_t = \frac{1}{1 + e^{-\beta(V(Positive)_t - V(Negative)_t)}} \tag{4}$$

246

247 where the inverse temperature (β , $\beta > 0$) regulates the stochasticity of how value computation
 248 influences choice. Choices were coded according to word valence (positive word = 1, negative
 249 word = 2), with action values initialised at zero at the start of each condition block.

250

251 Model estimation was conducted under the hierarchical Bayesian framework using the {rstan}
 252 (Carpenter et al., 2017), and {hBayesDM} (Ahn et al., 2017) R packages. Posterior inference
 253 was conducted through Markov Chain Monte Carlo (MCMC) sampling, using four MCMC
 254 chains with 1000 post-warmup iterations per chain. R-hat values for all parameters were
 255 inspected to examine model convergence (Gelman & Rubin, 1992) and were below 1.02.
 256 Model comparison was conducted by computing the leave-one-out information criterion
 257 (LOOIC) using the {loo} package (Vehtari et al., 2017), a commonly applied criterion for
 258 Bayesian analysis balancing model fit and complexity. Lower LOOIC scores signified superior
 259 model performance, and the model with the lowest LOOIC was chosen for further analyses.

260

261 **Exploratory factor analysis**

262 An exploratory factor analysis was conducted to identify transdiagnostic factors across the six
 263 psychiatric questionnaires used. To account for mixed item types (binary and ordinal),
 264 polychoric correlations were computed using the {polycor} package. Factor analysis was
 265 performed using the 'fa' function from the {psych} package in R with maximum likelihood
 266 estimation and oblique (oblimin) rotation (Costello & Osborne, 2005). Eighty of the eighty-two

267 items (two removed due to zero-variance) across the six questionnaires were included in the
268 analysis. A Cattell-Nelson-Gorsuch test was used through the {nFactors} package to
269 determine the optimal number of factors. Factor loadings ≥ 0.3 were deemed interpretable.
270 Factor scores were calculated using regression scoring weights and standardized for
271 subsequent analyses. To uncover transdiagnostic influences on behavioral and computational
272 outcomes, multiple regression analyses were conducted where transdiagnostic factor scores
273 and task condition, and their interaction served as independent variables in models, to predict
274 choice accuracy, reaction times, end-of-block ratings, and learning parameters as dependent
275 variables.

276

277 **Results**

278 Behavioral

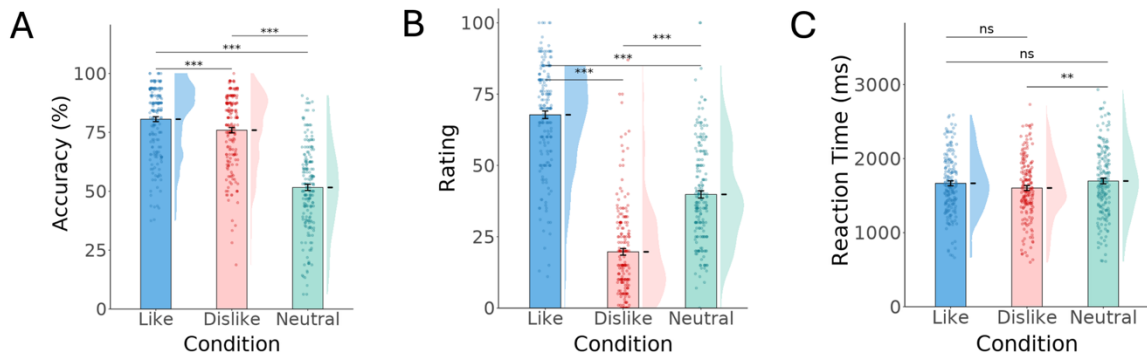
279 To test whether feedback valence affected learning accuracy, we examined choice accuracy
280 across the three persona conditions. We initially tested whether accuracy across valanced
281 conditions exceeded chance level (50%), where one-sample t-tests confirmed that accuracy
282 exceeded chance (50%) in both valanced conditions (Like: $t_{(192)} = 28.94$, $p < 0.001$; Dislike:
283 $t_{(192)} = 23.02$, $p < 0.001$) but not Neutral ($t_{(192)} = 1.24$, $p = 0.215$). (**Figure 2**). One-way repeated
284 measures ANOVA further revealed a significant main effect of condition on choice accuracy
285 ($F_{(2, 384)} = 208.41$, $p < .001$, $\eta^2p = 0.52$). Post-hoc comparisons with Tukey corrections showed
286 that the accuracy was the highest in the 'Like' condition ($M = 80.6 \pm 14.7\%$), followed by those
287 from the 'Dislike' ($M = 75.9 \pm 15.6\%$) and 'Neutral' ($M = 51.6 \pm 18.1\%$) conditions (**Figure 2A**;
288 Like vs. Neutral: $t_{(192)} = 17.259$, $p < 0.001$, Dislike vs. Neutral: $t_{(192)} = 14.597$, $p < 0.001$, Like
289 vs Dislike: $t_{(192)} = 3.982$, $p < 0.001$).

290

291 As a manipulation check, we examined end-of-block ratings of how much the feedback agent
292 liked them across conditions. One-sample t-tests against the scale midpoint confirmed that
293 Like ratings ($M = 67.75 \pm 18.37$) were significantly above the midpoint ($t_{(192)} = 13.42$, $p < .001$),
294 whilst Dislike ($M = 39.83 \pm 17.56$; $t_{(192)} = -25.98$, $p < .001$), and Neutral ratings ($M = 19.71 \pm$
295 16.20) were significantly below the midpoint ($t_{(192)} = -8.04$, $p < .001$). After applying a
296 Greenhouse-Geisser correction ($\epsilon = 0.924$) due to a sphericity violation ($W = 0.918$, $p < .001$),
297 a one-way repeated measures ANOVA revealed a significant effect of condition on ratings
298 ($F_{(1.85, 354.82)} = 331.93$, $p < .001$, $\eta^2p = 0.634$). Post-hoc Tukey-corrected pairwise comparisons
299 demonstrated significant differences between all conditions (all $p < .001$) (**Figure 2B**).

300

301 Finally, to examine whether feedback valence influenced decision time, we tested for an effect
 302 of persona condition on reaction time (RT). A one-way repeated measures ANOVA showed a
 303 significant effect of condition on RT ($F_{(2,384)} = 4.18, p = 0.016$), with post-hoc Tukey-corrected
 304 comparisons showing participants responded significantly faster in the 'Dislike' ($M = 1604 \pm$
 305 476ms) condition compared to 'Neutral' ($M = 1698 \pm 529\text{ms}$) ($p = 0.009$), while other pairwise
 306 comparisons were not significant ('Like' condition: $M = 1667 \pm 466\text{ms}$, **Figure 2C**).
 307



308 **Fig 2. Behavioral results for the Social Evaluation Learning Task.** Feedback valence
 309 demonstrated a significant effect upon all conditions with choice accuracy (A) and final ratings
 310 (B), and between Neutral and Dislike conditions for reaction time (C).
 311

312 *** p < 0.01, *** p < 0.001*

313
 314 **Computational modeling**

315 To determine whether participants exhibited asymmetric learning from positive versus
 316 negative feedback, we compared two computational models: a standard Rescorla-Wagner
 317 (RW) model with a single learning rate, and a dual-learning rate model with separate
 318 parameters for positive (α^+) and negative (α^-) feedback. Model comparison revealed the dual-
 319 learning model performed better than the simple RW model ($\Delta\text{LOOIC} = 271.09$), confirming
 320 that participants processed positive and negative social feedback differently. Parameters from
 321 the winning model showed good-to-excellent recovery, with posterior predictive checks
 322 accurately replicating the key patterns in our behavioural data (see Supplementary Materials:
 323 Supplementary Figure 1 and Supplementary Figure 2).

324
 325 After identifying the winning model, we tested whether learning rates varied as a function of
 326 feedback valence and persona condition. A learning rate valence (2: positive vs negative) by
 327 condition (3: Like, Dislike, Neutral) repeated measures ANOVA revealed a significant main
 328 effect for learning rate valence ($F_{(1,192)} = 545.63, p < 0.001$) and condition ($F_{(1.34, 256.33)} = 71.05,$
 329 $p < 0.001$) (**Figure 3A**). There was a significant interaction ($F_{(1.35, 260.07)} = 265.38, p < 0.001$),

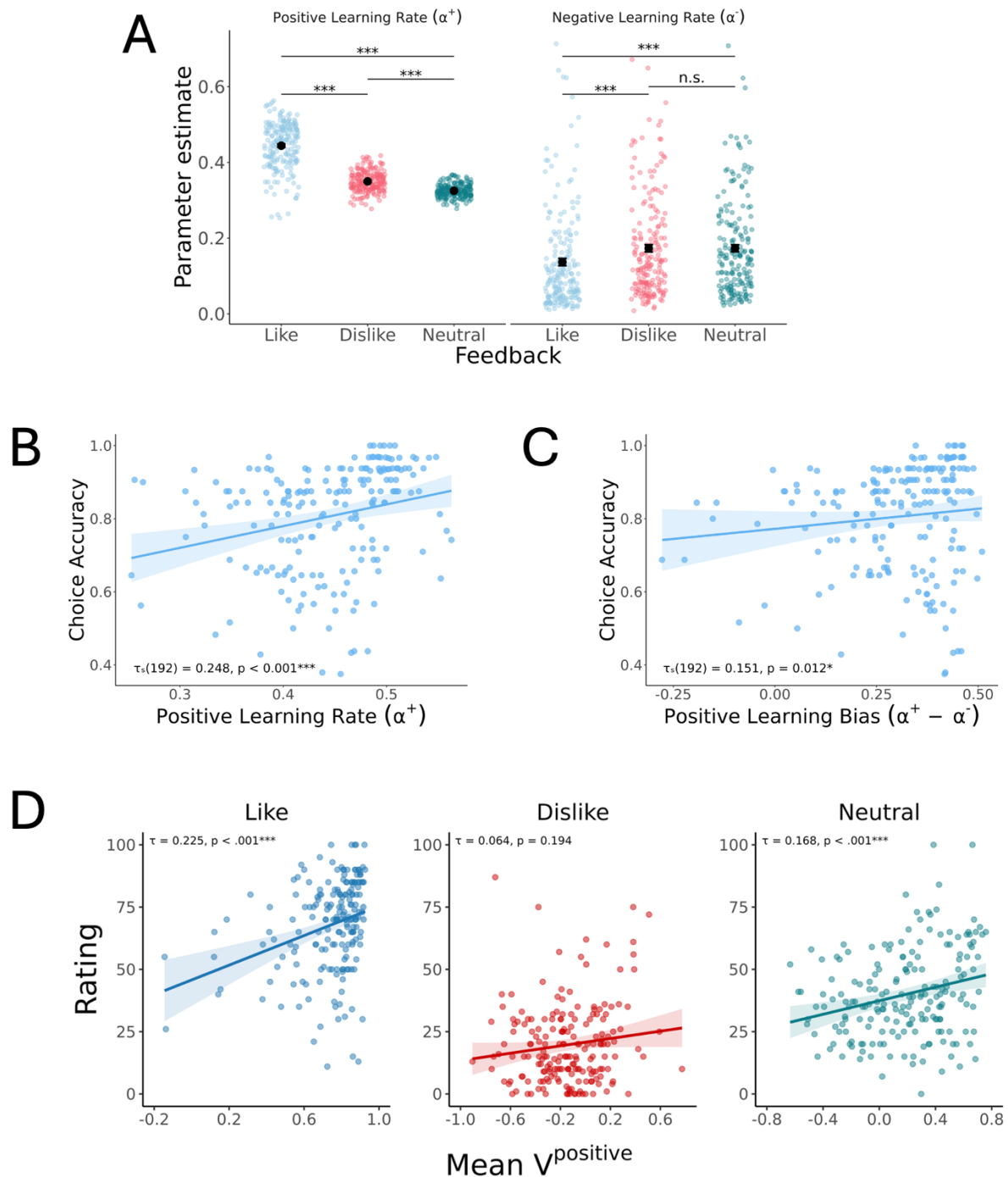
330 indicating that the relationship between feedback valence and learning varied across
331 conditions. Post-hoc pairwise comparisons using Tukey's adjustment showed the learning rate
332 for positive feedback was significantly higher in the 'Like' than 'Neutral' ($t_{(192)} = 34.38$, $p <$
333 0.001) and 'Dislike' ($t_{(192)} = 24.35$, $p < 0.001$) conditions, and in 'Dislike' than 'Neutral' ($t_{(192)} =$
334 25.29 , $p < 0.001$). For negative learning rates, 'Neutral' and 'Dislike' conditions did not differ
335 ($t_{(192)} = 0.01$, $p = 0.999$), but both differed significantly from 'Like' (Neutral vs Like: $t_{(192)} = 5.74$,
336 $p < 0.001$; Like vs Dislike: $t_{(192)} = 4.51$, $p < 0.001$), with the 'Like' condition showing higher
337 negative learning rates.

338

339 To examine whether individual differences in learning rates predicted behavioral performance,
340 we tested correlations between computational parameters and choice accuracy across
341 conditions. Both the positive learning rate (α^+ ; Kendall's $\tau = 0.249$, $p < 0.001$) and the
342 difference between the positive and negative learning rates - the positive learning bias - ($\alpha^+ -$
343 α^- ; Kendall's $\tau = 0.150$, $p = 0.014$) were positively associated with accuracy in the 'Like'
344 condition after Bonferroni correction (**Figure 3B-C**). However, no such relationship was
345 observed for the negative learning rate (all $\tau = 0.249$, $p > 0.36$). These results highlight a direct
346 link between the capacity to learn from positive social feedback and its recognition in social
347 contexts.

348

349 As a follow-up analysis, we computed the mean trial-averaged learned value of the positive-
350 word option from the winning model (V^{positive}) and participants' end-of-block ratings of how
351 much the feedback agent liked them (0–100), for each feedback condition (See
352 Supplementary Materials). This way, we can link the learning signal in the decisions task to
353 the end-of-block valence rating. Results showed that mean V^{positive} differed substantially across
354 conditions in the expected direction: Like ($M = 0.741 \pm 0.184$), Neutral ($M = 0.184 \pm 0.337$),
355 and Dislike ($M = -0.145 \pm 0.288$). Kendall's τ correlations between mean V^{positive} and explicit
356 ratings were significant in the Neutral ($\tau = 0.168$, $z = 3.27$, $p = .001$, $r = 0.26$) and Like ($\tau =$
357 0.225 , $z = 4.38$, $p < .001$, $r = 0.34$) conditions, but not in the Dislike condition ($\tau = 0.064$, $z =$
358 1.30 , $p = .194$, $r = 0.10$) (**Figure 3D**). These results provide an important validation check for
359 the winning model, highlighting a strong relationship between the model-derived learning
360 signal and observed persona ratings.



361

362 **Fig 3. Parameter estimates from the winning model and associations with behavior. (A)**

363 Group-level parameter estimates for negative and positive learning rates from the winning

364 model. Individual subject estimates shown as coloured points; black points indicate group

365 means with error bars representing standard error of the mean. **(B)** Positive learning rate

366 significantly correlated with choice accuracy in the Like condition (Kendall's $\tau = 0.248, p <$

367 0.001 , Bonferroni-corrected across six planned comparisons). **(C)** Positive learning bias ($\alpha^+ -$

368 α^-) significantly correlated with choice accuracy in the Like condition (Kendall's $\tau = 0.151, p =$

369 0.012 , Bonferroni-corrected). **(D)** Correlations between mean V^{positive} and subject ratings for

370 each feedback condition. Each point represents a single participant. Shaded regions represent
371 95% confidence intervals.

372

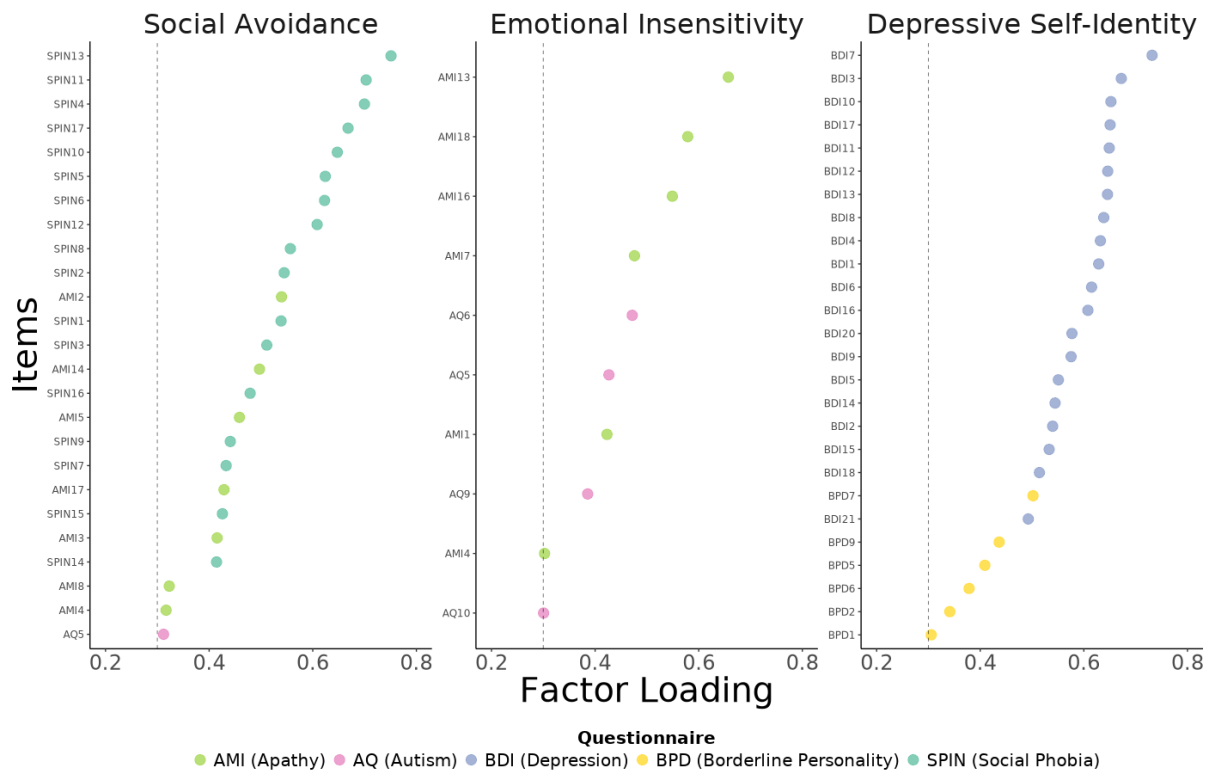
373 Exploratory factor analysis and associations with social learning

374 To identify transdiagnostic dimensions of psychopathology that might relate to social feedback
375 learning, we conducted an exploratory factor analysis (EFA) on the six questionnaires
376 administered. Exploratory factor analysis using maximum likelihood estimation with oblique
377 rotation, polychoric correlations to account for binary items, and the Cattell-Nelson-Gorsuch
378 test revealed that a 3-factor solution best accounted for the questionnaire data, explaining
379 25.8% of the variance.

380

381 Based on the strongest individual item loadings (**Figure 4**), factors were labelled accordingly
382 as 'Social Avoidance', with higher scores reflecting reduced social motivation and increased
383 avoidance, 'Emotional Insensitivity', where higher scores indicate high apathy and difficulties
384 evaluating and empathizing with others, and 'Depressive Self-Identity', where higher scores
385 correspond to low mood, low self-esteem and negative self-identity. Specifically, for the Social
386 Avoidance factor, the highest average loadings came from the SPIN questionnaire ($M = 0.57,$
387 ± 0.11), with significant contributions from AMI ($M = 0.19 \pm 0.23$) and AQ-10 ($M = 0.16 \pm 0.14$).
388 The 'Emotional Insensitivity' factor showed only moderate loadings from AMI ($M = 0.20 \pm 0.26$)
389 and AQ ($M = 0.17 \pm 0.22$), with no other questionnaires contributing meaningfully. Finally,
390 Depressive Self-Identity was dominated by items from the BDI questionnaire ($M = 0.59 \pm 0.10$),
391 followed by BPD ($M = 0.36 \pm 0.09$), with no other questionnaires reaching the 0.15 average
392 loading threshold.

393



394

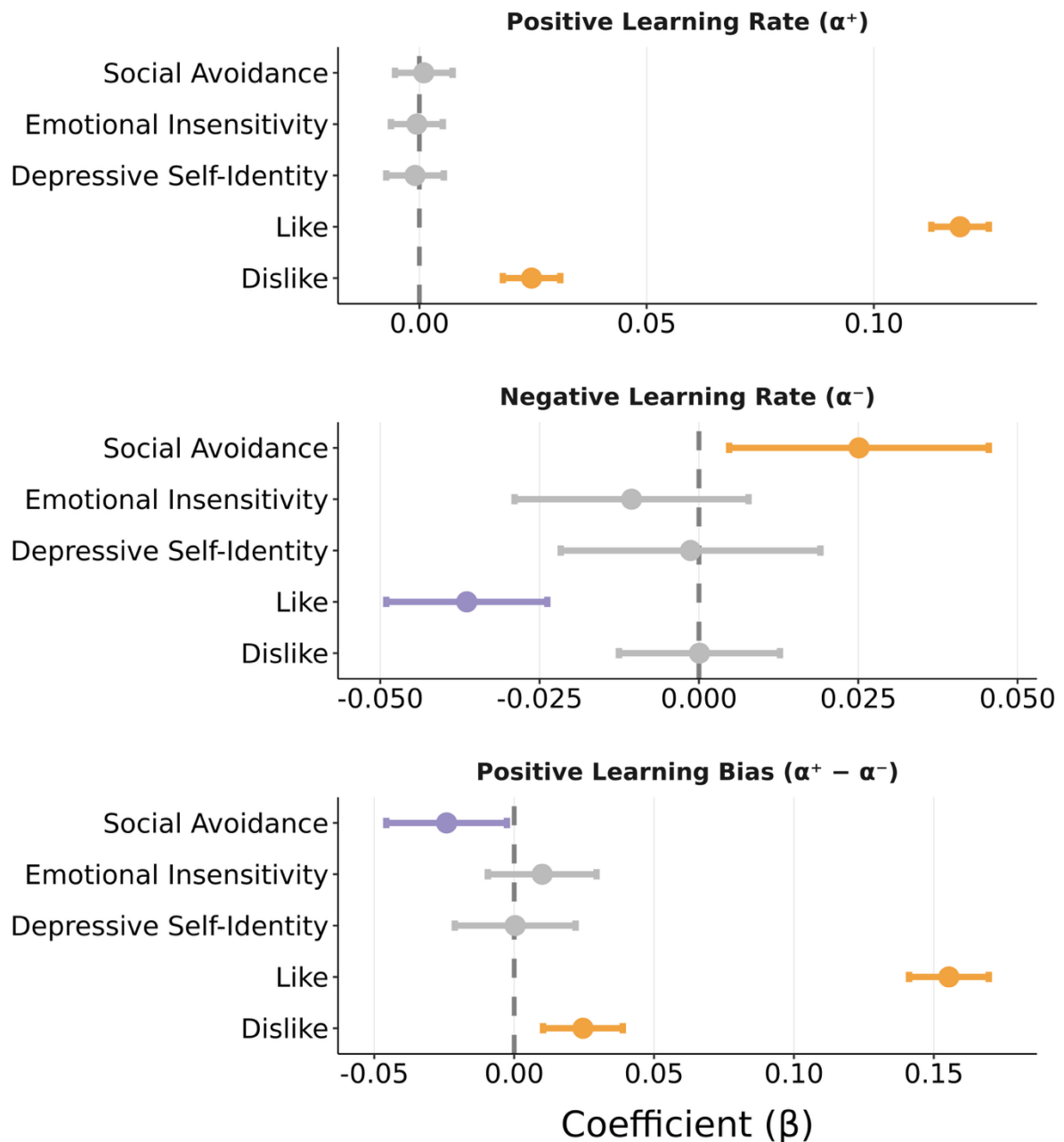
395 **Fig 4. Factor scores (loadings > 0.3) for the three-factor solution.** Factors were labelled
 396 'Social Avoidance' featuring items measuring social phobia and motivation, 'Emotional
 397 Insensitivity' consisting of items reflecting apathy and autistic traits, and 'Depressive Self-
 398 Identity' primarily loading items assessing depression and personality disorder. No significant
 399 loadings (> 0.3) were observed for the Narcissistic Admiration and Rivalry Questionnaire
 400 (NARQ).

401

402 **Transdiagnostic dimensions and condition-specific learning**

403 To examine whether transdiagnostic factors predicted computational mechanisms underlying
 404 social feedback learning, we fit linear mixed-effects models with computational parameters
 405 from the winning model as dependent variables. Models included condition (Like/Dislike
 406 relative to Neutral), transdiagnostic factor scores (Social Avoidance, Emotional Insensitivity,
 407 Depressive Self-Identity), and their interactions as fixed factors, with random intercepts for
 408 subjects to account for the repeated measures structure (formula: $DV \sim \text{condition} \times (\text{FA1} +$
 409 $\text{FA2} + \text{FA3}) + (1|\text{subID})$) (**Figure 5**). Condition effects on learning parameters were robust,
 410 with the Like condition significantly increased positive learning rates ($\beta = 0.119, p < 0.001$)
 411 and decreased negative learning rates ($\beta = -0.036, p < 0.001$), while the Dislike condition
 412 significantly increased both positive learning rates ($\beta = 0.025, p < 0.001$) and positive learning
 413 bias ($\beta = 0.025, p < 0.001$). Additionally, the 'Social Avoidance' factor significantly predicted
 414 negative learning rates ($\beta = 0.025, p = 0.016$) and learning rate difference ($\beta = -0.024, p =$

415 0.029) in opposite directions. Specifically, individuals with higher scores on this 'Social
 416 Avoidance' factor learned more from negative social feedback, reflecting a negative learning
 417 bias. No other transdiagnostic factors significantly predicted computational learning
 418 parameters, suggesting that the negative learning bias observed is specific to traits related to
 419 social avoidance and low social motivation.



420
 421 **Fig 5. Main effect coefficients for linear regressions predicting computational**
 422 **parameters from the winning model.** Top: Predictors of the positive learning rate (α^+)
 423 Middle: Predictors of the negative learning rate (α^-). Bottom: Predictors of the positive learning
 424 bias ($\alpha^+ - \alpha^-$). Yellow = significant positive effects ($p < 0.05$), Purple = significant negative
 425 effects ($p < 0.05$), grey = non-significant effects ($p \geq 0.05$). Error bars = 95% CI.

426

427 We next sought to understand whether the three transdiagnostic factors moderated the effect
428 of feedback valence on behavioral performance. To test this, we fit a series of three-way
429 mixed-effects regression models with transdiagnostic factors, learning rate parameters (α^+
430 and α^-), and valence conditions as predictors for task performance (**Figure 6**).

431

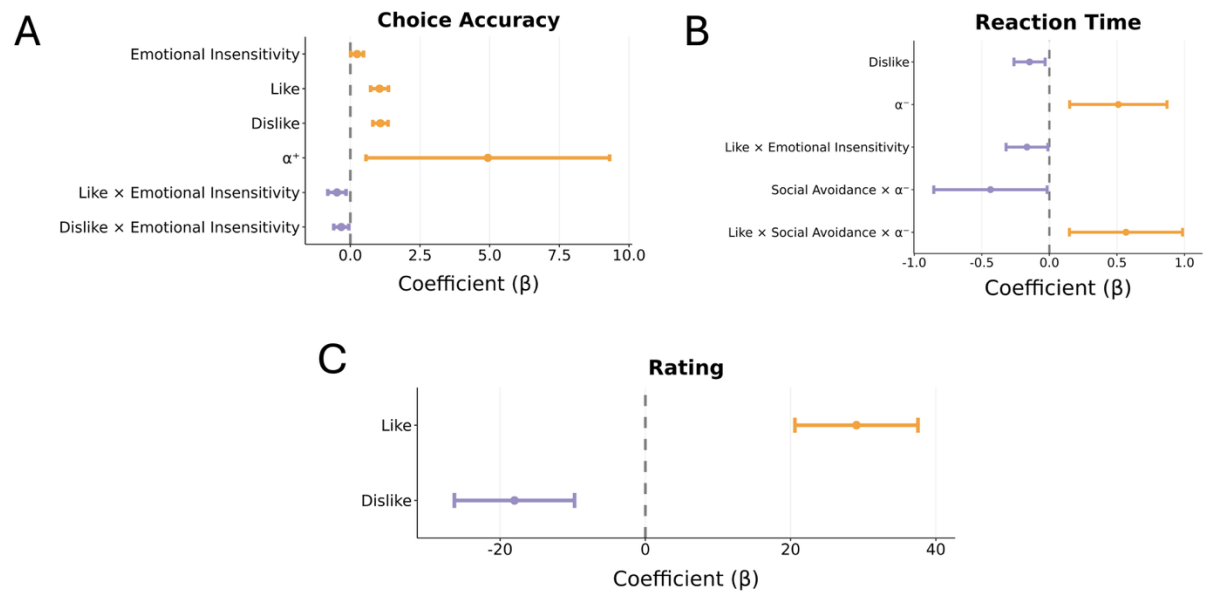
432 First, to examine effects on choice accuracy (**Figure 6A**), we fit a mixed-effects logistic
433 regression model predicting choice accuracy, which revealed significant effects for both 'Like'
434 ($\beta = 1.049$, $p < 0.001$) and 'Dislike' conditions ($\beta = 1.080$, $p < 0.001$), consistent with our
435 repeated measures ANOVA analysis (cf. Figure 2A). Additionally, we observed a significant
436 positive main effect of the positive learning rate parameter ($\beta = 4.935$, $p = 0.027$), indicating
437 that individuals who learn more rapidly from positive feedback showed higher overall choice
438 accuracy. A significant main effect of 'Emotional Insensitivity' was also observed ($\beta = 0.243$,
439 $p = 0.041$), suggesting higher scores on this factor were associated with increased choice
440 accuracy. However, significant two-way interactions between the 'Like' ($\beta = -0.480$, $p = 0.004$)
441 and 'Dislike' ($\beta = -0.323$, $p = 0.020$) conditions with 'Emotional Insensitivity' revealed that
442 higher scores on this factor attenuated the benefits in accuracy typically observed in value-
443 congruent learning contexts. No other two-way interactions with condition or learning rate
444 parameters were significant, and no three-way interactions reached significance (all $p > 0.05$).

445

446 Next, we fit linear mixed-effects models to determine whether transdiagnostic factors and
447 learning rate parameters influenced decision time (**Figure 6B**). This revealed a significant
448 main effect of the 'Dislike' condition ($\beta = -0.147$, $p = 0.014$), indicating faster responses. A
449 significant main effect of the negative learning rate parameter was also observed ($\beta = 0.510$,
450 $p = 0.006$), indicating that individuals who learned more rapidly from negative feedback
451 showed slower overall response times. A significant two-way interaction between Social
452 Avoidance and the negative learning rate parameter ($\beta = -0.436$, $p = 0.043$) linked higher
453 Social Avoidance scores with a reduced relationship between negative learning rate and
454 response time, whilst a significant two-way interaction between the 'Like' condition and
455 'Emotional Insensitivity' ($\beta = -0.166$, $p = 0.037$) indicated that higher scores on this factor were
456 associated with faster responses specifically in the 'Like' condition. Furthermore, a significant
457 three-way interaction between the 'Like' condition, Social Avoidance, and the negative
458 learning rate parameter ($\beta = 0.567$, $p = 0.009$) revealed that individuals with higher Social
459 Avoidance scores and higher negative learning rates exhibited slower response times in the
460 'Like' condition.

461
462
463
464
465
466
467
468

Similarly, linear mixed-effects models predicting participants' end-of-block ratings (**Figure 6C**) revealed significant main effects of task condition for both 'Like' ($\beta = 14.54$, $p < 0.001$) and 'Dislike' ($\beta = -9.01$, $p = 0.001$) personas, consistent with our manipulation of social feedback valence. No significant main effects of transdiagnostic factors or learning rate parameters were observed, and no two-way or three-way interactions between factors, learning rate parameters, and conditions reached significance (all $p > 0.05$).



469
470
471
472
473
474
475
476

Fig 6. Significant three-way mixed-effects model coefficients. (A) Beta coefficients for choice accuracy logistic regression showing condition main effects and interactions. (B) Beta coefficients for reaction time showing condition main effects and interactions. (C) Beta coefficients for ratings showing only condition main effects. Yellow = significant positive effects ($p < 0.05$), Purple = significant negative effects ($p < 0.05$), grey = non-significant effects ($p \geq 0.05$). Error bars = 95% CI.

477 Discussion

478 The current study investigated how behavioral traits affect how individuals learn from social
479 feedback utilizing a transdiagnostic approach and computational modeling under the
480 reinforcement learning framework. Model comparison revealed that a modified Rescorla-
481 Wagner model with separate learning rates for positive and negative feedback fit the data
482 better than a standard model, indicating asymmetric learning processes. This aligns with
483 previous findings that reinforcement learning involves differential updating from positive
484 versus negative (social) outcomes (Behrens et al., 2008; Zhang et al., 2020). These learning

485 rates were found to vary by task condition, with the learning rate from positive feedback
486 highest in the 'Like' condition. Individual differences in positive learning rates and positive
487 learning bias were also both associated with higher choice accuracy in the 'Like' condition.
488 This suggests that individuals who learn more strongly from positive feedback make more
489 accurate evaluative judgements in socially positive contexts.

490

491 Our factor analysis extracted three robust transdiagnostic dimensions: 'Social Avoidance',
492 'Emotional Insensitivity' and 'Depressive Self-Identity'. These dimensions partly overlap with
493 previous studies extracting transdiagnostic factors of psychopathology from general
494 population samples (Gillan et al., 2016; Hoven et al., 2023; Oka et al., 2025). On-the-other-
495 hand, the presence of an 'Emotional Insensitivity' factor not reported previously reflects the
496 current study's specific focus towards assessing social evaluative learning. Indeed, this factor
497 loaded items from the Autism Spectrum Quotient (AQ-10) and Apathy Motivation Index (AMI),
498 questionnaires not included together in previous studies. The extraction of this factor highlights
499 the ability to understand and interpret the emotions of others as an important transdiagnostic
500 dimension of social learning.

501

502 Our regression analyses found significant associations between extracted transdiagnostic
503 dimensions and task-related assessments of social evaluation learning. Specifically, we show
504 that a latent Social Avoidance dimension, cutting across social anxiety and related traits, is
505 selectively associated with a negative learning bias and slower response times towards
506 positive social feedback. This finding extends previous research using categorical
507 assessments of social anxiety and borderline personality disorder (BPD) demonstrating similar
508 learning biases using reinforcement learning models (Hoffmann et al., 2024; Koban et al.,
509 2017, 2023; Korn et al., 2016) and increased deliberation in response to positive social
510 information (Hunter et al., 2022; Van Der Molen et al., 2014). The Social Avoidance factor
511 strongly loaded items from the Social Phobia Inventory (SPIN) and those from the 'Social
512 Motivation' subscale of the Apathy Motivation Index (AMI). Lowered social motivation - the
513 interest and desire to engage in social interactions (Chevallier et al., 2012; Geen, 1991;
514 Kanterman & Shamay-Tsoory, 2025) - is commonly reported among individuals with social
515 anxiety (Abplanalp et al., 2022; Alden & Taylor, 2010; Blay et al., 2021; Geen, 1991; Yang et
516 al., 2024), and co-morbidly presents in autism spectrum disorder (ASD) (Bagg et al., 2024;
517 Briot et al., 2020; Spain et al., 2018), schizophrenia (Bershad & de Wit, 2023; Catalano &
518 Green, 2023; Hajdúk et al., 2024) and attention deficit hyperactivity disorder (ADHD) (Martin
519 et al., 2024). It is proposed that lowered social motivation distinctly arises from changes to

520 cognitive processes associated with effort-based decision-making, which lead to individuals
521 underestimating the benefits and overestimating the effort costs associated with social
522 interactions (Catalano & Green, 2023; Kanterman & Shamay-Tsoory, 2025). Recent accounts
523 have identified the motivation to socially engage as a key transdiagnostic measure given its
524 comorbid presentation across multiple disorders (Barkus & Badcock, 2019) that may have
525 distinct clinical and neurobiological correlates and respond differently to treatments (Chetcuti
526 et al., 2026). By extracting this specific measure and revealing an association with learning
527 parameters, we highlight a potential target candidate for social motivation.

528

529 A second transdiagnostic factor, Emotional Insensitivity, reflecting lowered empathy and
530 problems inferring the emotions of others, attenuated performance in both Like and Dislike
531 conditions. As empathy requires accurate emotional identification (Bird & Viding, 2014;
532 Shamay-Tsoory & Hertz, 2022), individuals with reduced empathy are hypothesized to be less
533 capable of discerning changes in other's emotions (Coll et al., 2017). Indeed, reduced
534 empathy and emotion recognition have been linked to poorer performance on complex social
535 cognitive tasks through impaired interpretation of others' affective states (Gonzalez-Gadea et
536 al., 2014). This is commonly observed among individuals with autism spectrum disorder (ASD)
537 who show difficulties in social interactions (Bolis & Schilbach, 2017), due to challenges with
538 understanding the intentions of others (Rosenthal et al., 2019). In the current study, measures
539 from the Autism Quotient (AQ-10) and 'Emotional Sensitivity' subscale of the Apathy
540 Motivation Index (AMI) loaded strongly onto the Emotional Insensitivity factor. Apathy is a
541 multi-dimensional construct, consisting of executive, emotional, and initiation domains, each
542 supported by distinct cognitive and neural mechanisms (Dickson & Husain, 2022). Most
543 relevant to social behaviour is emotional apathy, characterised by emotional blunting, reduced
544 empathy, and altered social interactions. Emotional apathy is robustly associated with
545 prosocial behaviour across different contexts (Contreras-Huerta et al., 2022), with higher
546 levels predicting a lowered tendency to help out others (Lockwood et al., 2017a). Yet, few
547 studies have specifically investigated its role with social reward learning. In a study of
548 dementia patients, those with high emotional apathy we found to demonstrate impaired
549 learning towards social rewards (Wong et al., 2023). Our results further show that difficulties
550 with interpreting the emotions of others directly impairs the ability to accurately evaluate
551 positive and negative social feedback about the self. Therefore, emotional sensitivity is an
552 important trait for correctly identifying and responding to the social preferences of others.
553 Recent accounts also highlight individual differences in empathy and apathy, where affective
554 empathy and emotional motivation are underpinned by the same latent factor (Lockwood et

555 al., 2017b). Our transdiagnostic approach was able to extract these socially-relevant
556 measures from the AMI and AQ-10 that may be obscured in categorical assessments or
557 questionnaires. We subsequently advocate for similar dimensional measures to be used when
558 assessing social behaviour and cognition.

559

560 The present study has some limitations. Importantly, our sample is mainly comprised of
561 females and consisted exclusively of university students with a mean age of approximately
562 19. This may explain the lack of any significant results with the Depressive Self-Identity factor,
563 which strongly loaded items assessing depressive symptoms centred on guilt, lack of
564 satisfaction and suicidal ideation. The reported Beck Depression Inventory (BDI) scores in our
565 sample ($M = 12.52 \pm 9.34$) primarily reflect minimal (0-13) or mild (14-19) depression (Beck et
566 al., 1996). However, depressive symptomatology may only emerge with more severe cases
567 than those present in our non-clinical student population. This aligns with dimensional
568 perspectives suggesting that certain cognitive-affective impairments may only manifest
569 beyond a clinical threshold (Insel et al., 2010). Future studies should therefore aim to capture
570 a wider range of depressive symptoms across a more diverse sample.

571

572 Taken together, combining a transdiagnostic measure of behavior with computational
573 modelling, our study identifies key behavioral signatures of biased social evaluation learning.
574 Computational modelling revealed that participants exhibited asymmetric learning from
575 positive versus negative social feedback, with learning rates varying by feedback valence. A
576 transdiagnostic factor reflecting social anxiety and motivational deficits predicted enhanced
577 learning from negative feedback and a reduced positive learning bias, demonstrating biased
578 social learning towards negatively salient information. An 'Emotional Insensitivity' factor
579 capturing difficulties in empathy and mentalizing was associated with reduced choice
580 accuracy in valenced conditions and faster reaction times in response to positive social
581 feedback, suggesting reduced deliberative processing. Together, these findings demonstrate
582 how transdiagnostic traits influence social evaluative learning, identifying candidate
583 computational processes that can be targeted in future studies.

584 **CRedit authorship contribution statement**

585 A.S. – Methodology, Software, Validation, Formal analysis, Data Curation, Writing - Original
586 Draft, Writing - Review & Editing, Visualization

587 A.M. – Methodology, Validation, Formal analysis, Investigation, Writing - Review & Editing

588 M.S. – Methodology, Validation, Formal analysis, Investigation, Writing - Review & Editing

589 L.P. – Methodology, Validation, Formal analysis, Investigation, Writing - Review & Editing

590 Y.L. – Validation, Formal analysis, Investigation, Writing - Review & Editing

591 G.R. – Validation, Formal analysis, Investigation, Writing - Review & Editing

592 A.C.-F. – Validation, Formal analysis, Investigation, Writing - Review & Editing

593 J.R. – Validation, Formal analysis, Investigation, Writing - Review & Editing

594 H.L. – Validation, Formal analysis, Investigation, Writing - Review & Editing

595 T.S. – Validation, Formal analysis, Investigation, Writing - Review & Editing

596 A.C.S. – Validation, Formal analysis, Writing - Review & Editing, Visualization, Supervision

597 L.Z. – Conceptualization, Methodology, Software, Validation, Formal analysis, Resources,
598 Data Curation, Writing - Review & Editing, Visualization, Supervision, Project administration,
599 Funding acquisition

600

601 **Ethical standards**

602 The authors assert that all procedures contributing to this work comply with the ethical
603 standards of the relevant national and institutional committees on human experimentation and
604 with the Helsinki Declaration of 1975, as revised in 2008.

605

606 **Data availability statement**

607 All code and data for data analysis and figure generation are publicly available at
608 [10.5281/zenodo.19161237](https://doi.org/10.5281/zenodo.19161237).

609

610 **Acknowledgements**

611 A.S was funded by an MRC AIM iCASE Studentship (MR/W007002/1). A.C.S. was funded by
612 a Walter Benjamin Fellowship from the German Research Foundation (DFG/534697423). L.Z.
613 was partially supported by the Wellcome Trust (228268/Z/23/Z) and the Royal Society
614 (IES\R3\243253).

615

616 **Competing interests**

617 We declare no conflict of interest.

618

619 **References**

- 620 Abplanalp, S. J., Mote, J., Uhlman, A. C., Weizenbaum, E., Alvi, T., Tabak, B. A., & Fulford,
621 D. (2022). Parsing social motivation: Development and validation of a self-report
622 measure of social effort. *Journal of Mental Health, 31*(3), 366–373.
623 <https://doi.org/10.1080/09638237.2021.1952948>
- 624 Adriaens, P. R., & De Block, A. (2013). Why We Essentialize Mental Disorders. *The Journal*
625 *of Medicine and Philosophy: A Forum for Bioethics and Philosophy of Medicine, 38*(2),
626 107–127. <https://doi.org/10.1093/jmp/jht008>
- 627 Ahn, W.-Y., Haines, N., & Zhang, L. (2017). Revealing Neurocomputational Mechanisms of
628 Reinforcement Learning and Decision-Making With the hBayesDM Package.
629 *Computational Psychiatry (Cambridge, Mass.), 1*, 24–57.
630 https://doi.org/10.1162/CPSY_a_00002
- 631 Alden, L. E., & Taylor, C. T. (2010). Interpersonal processes in social anxiety disorder. In
632 *Interpersonal processes in the anxiety disorders: Implications for understanding*
633 *psychopathology and treatment* (pp. 125–152). American Psychological Association.
634 <https://doi.org/10.1037/12084-005>
- 635 Allison, C., Auyeung, B., & Baron-Cohen, S. (2012). Toward Brief “Red Flags” for Autism
636 Screening: The Short Autism Spectrum Quotient and the Short Quantitative Checklist
637 in 1,000 Cases and 3,000 Controls. *Journal of the American Academy of Child &*
638 *Adolescent Psychiatry, 51*(2), 202-212.e7. <https://doi.org/10.1016/j.jaac.2011.11.003>
- 639 Ang, Y.-S., Lockwood, P., Apps, M. A. J., Muhammed, K., & Husain, M. (2017). Distinct
640 Subtypes of Apathy Revealed by the Apathy Motivation Index. *PLOS ONE, 12*(1),
641 e0169938. <https://doi.org/10.1371/journal.pone.0169938>
- 642 Anwyl-Irvine, A. L., Massonnié, J., Flitton, A., Kirkham, N., & Evershed, J. K. (2020). Gorilla in
643 our midst: An online behavioral experiment builder. *Behavior Research Methods,*
644 *52*(1), 388–407. <https://doi.org/10.3758/s13428-019-01237-x>
- 645 Back, M. D., Küfner, A. C. P., Dufner, M., Gerlach, T. M., Rauthmann, J. F., & Denissen, J. J.
646 A. (2013). Narcissistic admiration and rivalry: Disentangling the bright and dark sides
647 of narcissism. *Journal of Personality and Social Psychology, 105*(6), 1013–1037.
648 <https://doi.org/10.1037/a0034431>
- 649 Bagg, E., Pickard, H., Tan, M., Smith, T. J., Simonoff, E., Pickles, A., Carter Leno, V., &
650 Bedford, R. (2024). Testing the social motivation theory of autism: The role of co-
651 occurring anxiety. *Journal of Child Psychology and Psychiatry, 65*(7), 899–909.
652 <https://doi.org/10.1111/jcpp.13925>

- 653 Barkus, E., & Badcock, J. C. (2019). A Transdiagnostic Perspective on Social Anhedonia.
654 *Frontiers in Psychiatry, 10*. <https://doi.org/10.3389/fpsy.2019.00216>
- 655 Bates, D., Mächler, M., Bolker, B., & Walker, S. (2015). Fitting Linear Mixed-Effects Models
656 Using lme4. *Journal of Statistical Software, 67*, 1–48.
657 <https://doi.org/10.18637/jss.v067.i01>
- 658 Bauer, R. A. (1967). Societal Feedback. *The ANNALS of the American Academy of Political
659 and Social Science, 373*(1), 180–192. <https://doi.org/10.1177/000271626737300109>
- 660 Beck, A. T., Steer, R. A., Ball, R., & Ranieri, W. F. (1996). Comparison of Beck Depression
661 Inventories-IA and-II in Psychiatric Outpatients. *Journal of Personality Assessment,
662 67*(3), 588–597. https://doi.org/10.1207/s15327752jpa6703_13
- 663 Behrens, T. E. J., Hunt, L. T., Woolrich, M. W., & Rushworth, M. F. S. (2008). Associative
664 learning of social value. *Nature, 456*(7219), 245–249.
665 <https://doi.org/10.1038/nature07538>
- 666 Bershada, A. K., & de Wit, H. (2023). Social Psychopharmacology: Novel Approaches to Treat
667 Deficits in Social Motivation in Schizophrenia. *Schizophrenia Bulletin, 49*(5), 1161–
668 1173. <https://doi.org/10.1093/schbul/sbad094>
- 669 Bird, G., & Viding, E. (2014). The self to other model of empathy: Providing a new framework
670 for understanding empathy impairments in psychopathy, autism, and alexithymia.
671 *Neuroscience & Biobehavioral Reviews, 47*, 520–532.
672 <https://doi.org/10.1016/j.neubiorev.2014.09.021>
- 673 Blay, Y., Keshet, H., Friedman, L., & Gilboa-Schechtman, E. (2021). Interpersonal motivations
674 in social anxiety: Weakened approach and intensified avoidance motivations for
675 affiliation and social-rank. *Personality and Individual Differences, 170*, 110449.
676 <https://doi.org/10.1016/j.paid.2020.110449>
- 677 Bolis, D., & Schilbach, L. (2017). Beyond one Bayesian brain: Modeling intra- and inter-
678 personal processes during social interaction: Commentary on “Mentalizing
679 homeostasis: The social origins of interoceptive inference” by Fotopoulou & Tsakiris.
680 *Neuropsychoanalysis, 19*(1), 35–38. <https://doi.org/10.1080/15294145.2017.1295215>
- 681 Borsboom, D. (2017). A network theory of mental disorders. *World Psychiatry, 16*(1), 5–13.
682 <https://doi.org/10.1002/wps.20375>
- 683 Boyd, R., Richerson, P. J., & Henrich, J. (2011). The cultural niche: Why social learning is
684 essential for human adaptation. *Proceedings of the National Academy of Sciences,
685 108*(supplement_2), 10918–10925. <https://doi.org/10.1073/pnas.1100290108>
- 686 Briot, K., Jean, F., Jouni, A., Geoffray, M.-M., Ly-Le Moal, M., Umbricht, D., Chatham, C.,
687 Murtagh, L., Delorme, R., Bouvard, M., Leboyer, M., & Amestoy, A. (2020). Social

688 Anxiety in Children and Adolescents With Autism Spectrum Disorders Contribute to
689 Impairments in Social Communication and Social Motivation. *Frontiers in Psychiatry*,
690 11. <https://doi.org/10.3389/fpsyt.2020.00710>

691 Button, K. S., Browning, M., Munafò, M. R., & Lewis, G. (2012). Social inference and social
692 anxiety: Evidence of a fear-congruent self-referential learning bias. *Journal of Behavior
693 Therapy and Experimental Psychiatry*, 43(4), 1082–1087.
694 <https://doi.org/10.1016/j.jbtep.2012.05.004>

695 Button, K. S., Kounali, D., Stapinski, L., Rapee, R. M., Lewis, G., & Munafò, M. R. (2015). Fear
696 of Negative Evaluation Biases Social Evaluation Inference: Evidence from a
697 Probabilistic Learning Task. *PLOS ONE*, 10(4), e0119456.
698 <https://doi.org/10.1371/journal.pone.0119456>

699 Carpenter, B., Gelman, A., Hoffman, M. D., Lee, D., Goodrich, B., Betancourt, M., Brubaker,
700 M. A., Guo, J., Li, P., & Riddell, A. (2017). Stan: A Probabilistic Programming
701 Language. *Journal of Statistical Software*, 76, 1. <https://doi.org/10.18637/jss.v076.i01>

702 Catalano, L. T., & Green, M. F. (2023). Social Motivation in Schizophrenia: What's Effort Got
703 to Do With It? *Schizophrenia Bulletin*, 49(5), 1127–1137.
704 <https://doi.org/10.1093/schbul/sbad090>

705 Chetcuti, L., Hardan, A. Y., Frazier, T. W., Loth, E., McPartland, J. C., Youngstrom, E. A., &
706 Uljarevic, M. (2026). Advancing scientific understanding of the drive to socially engage:
707 From broad constructs to transdiagnostic 'building blocks'. *Molecular Psychiatry*, 31(1),
708 599–609. <https://doi.org/10.1038/s41380-025-03185-9>

709 Chevallier, C., Kohls, G., Troiani, V., Brodtkin, E. S., & Schultz, R. T. (2012). The social
710 motivation theory of autism. *Trends in Cognitive Sciences*, 16(4), 231–239.
711 <https://doi.org/10.1016/j.tics.2012.02.007>

712 Coll, M.-P., Viding, E., Rütgen, M., Silani, G., Lamm, C., Catmur, C., & Bird, G. (2017). Are
713 we really measuring empathy? Proposal for a new measurement framework.
714 *Neuroscience & Biobehavioral Reviews*, 83, 132–139.
715 <https://doi.org/10.1016/j.neubiorev.2017.10.009>

716 Connor, K. M., Davidson, J. R. T., Churchill, L. E., Sherwood, A., Weisler, R. H., & Foa, E.
717 (2000). Psychometric properties of the Social Phobia Inventory (SPIN): New self-rating
718 scale. *British Journal of Psychiatry*, 176(4), 379–386.
719 <https://doi.org/10.1192/bjp.176.4.379>

720 Contreras-Huerta, L. S., Lockwood, P. L., Bird, G., Apps, M. A. J., & Crockett, M. J. (2022).
721 Prosocial behavior is associated with transdiagnostic markers of affective sensitivity in
722 multiple domains. *Emotion*, 22(5), 820–835. <https://doi.org/10.1037/emo0000813>

723 Costello, A. B., & Osborne, J. (2005). Best practices in exploratory factor analysis: Four
724 recommendations for getting the most from your analysis. *Practical Assessment,*
725 *Research, and Evaluation, 10*(1). <https://doi.org/10.7275/jyj1-4868>

726 Diaconescu, A. O., Mathys, C., Weber, L. A. E., Daunizeau, J., Kasper, L., Lomakina, E. I.,
727 Fehr, E., & Stephan, K. E. (2014). Inferring on the Intentions of Others by Hierarchical
728 Bayesian Learning. *PLOS Computational Biology, 10*(9), e1003810.
729 <https://doi.org/10.1371/journal.pcbi.1003810>

730 Dickson, S. S., & Husain, M. (2022). Are there distinct dimensions of apathy? The argument
731 for reappraisal. *Cortex, 149*, 246–256. <https://doi.org/10.1016/j.cortex.2022.01.001>

732 Elder, J., Davis, T., & Hughes, B. L. (2022). Learning About the Self: Motives for Coherence
733 and Positivity Constrain Learning From Self-Relevant Social Feedback. *Psychological*
734 *Science, 33*(4), 629–647. <https://doi.org/10.1177/09567976211045934>

735 Faul, F., Erdfelder, E., Lang, A.-G., & Buchner, A. (2007). G*Power 3: A flexible statistical
736 power analysis program for the social, behavioral, and biomedical sciences. *Behavior*
737 *Research Methods, 39*(2), 175–191. <https://doi.org/10.3758/BF03193146>

738 Feczko, E., Miranda-Dominguez, O., Marr, M., Graham, A. M., Nigg, J. T., & Fair, D. A. (2019).
739 The Heterogeneity Problem: Approaches to Identify Psychiatric Subtypes. *Trends in*
740 *Cognitive Sciences, 23*(7), 584–601. <https://doi.org/10.1016/j.tics.2019.03.009>

741 FeldmanHall, O., & Chang, L. J. (2018). Social Learning. In *Goal-Directed Decision Making*
742 (pp. 309–330). Elsevier. <https://doi.org/10.1016/B978-0-12-812098-9.00014-0>

743 First, M. B., Gibbon, M., Spitzer, R. L., Williams, J. B. W., & Benjamin, L. S. (1997). *Structured*
744 *Clinical Interview for DSM-IV Axis II Personality Disorders, (SCID-II)*.
745 <https://cir.nii.ac.jp/crid/1370004237531794560>

746 Geen, R. G. (1991). Social motivation. *Annual Review of Psychology, 42*, 377–399.
747 <https://doi.org/10.1146/annurev.ps.42.020191.002113>

748 Gelman, A., & Rubin, D. B. (1992). Inference from Iterative Simulation Using Multiple
749 Sequences. *Statistical Science, 7*(4). <https://doi.org/10.1214/ss/1177011136>

750 Gilboa-Schechtman, E., Keshet, H., Livne, T., Berger, U., Zabag, R., Hermesh, H., & Marom,
751 S. (2017). Explicit and implicit self-evaluations in social anxiety disorder. *Journal of*
752 *Abnormal Psychology, 126*(3), 285–290. <https://doi.org/10.1037/abn0000261>

753 Gillan, C. M., Kosinski, M., Whelan, R., Phelps, E. A., & Daw, N. D. (2016). Characterizing a
754 psychiatric symptom dimension related to deficits in goal-directed control. *eLife, 5*,
755 e11305. <https://doi.org/10.7554/eLife.11305>

- 756 Gillan, C. M., & Seow, T. X. F. (2020). Carving Out New Transdiagnostic Dimensions for
757 Research in Mental Health. *Biological Psychiatry: Cognitive Neuroscience and*
758 *Neuroimaging*, 5(10), 932–934. <https://doi.org/10.1016/j.bpsc.2020.04.013>
- 759 Gkika, S., Wittkowski, A., & Wells, A. (2018). Social cognition and metacognition in social
760 anxiety: A systematic review. *Clinical Psychology & Psychotherapy*, 25(1), 10–30.
761 <https://doi.org/10.1002/cpp.2127>
- 762 Gonzalez-Gadea, M. L., Herrera, E., Parra, M., Gomez Mendez, P., Baez, S., Manes, F., &
763 Ibanez, A. (2014). Emotion recognition and cognitive empathy deficits in adolescent
764 offenders revealed by context-sensitive tasks. *Frontiers in Human Neuroscience*, 8.
765 <https://doi.org/10.3389/fnhum.2014.00850>
- 766 Green, M. F., & Leitman, D. I. (2008). Social Cognition in Schizophrenia. *Schizophrenia*
767 *Bulletin*, 34(4), 670–672. <https://doi.org/10.1093/schbul/sbn045>
- 768 Hajdúk, M., Abplanalp, S. J., Jimenez, A. M., Fisher, M., Haut, K. M., Hooker, C. I., Lee, H.,
769 Ventura, J., Nahum, M., & Green, M. F. (2024). Linking social motivation, general
770 motivation, and social cognition to interpersonal functioning in schizophrenia: Insights
771 from exploratory graph analysis. *European Archives of Psychiatry and Clinical*
772 *Neuroscience*, 274(6), 1385–1393. <https://doi.org/10.1007/s00406-023-01733-4>
- 773 Herpertz, S. C., & Bertsch, K. (2014). The social-cognitive basis of personality disorders.
774 *Current Opinion in Psychiatry*, 27(1), 73.
775 <https://doi.org/10.1097/YCO.000000000000026>
- 776 Hobbs, C., Vozarova, P., Sabharwal, A., Shah, P., & Button, K. (2022). Is depression
777 associated with reduced optimistic belief updating? *Royal Society Open Science*, 9(2),
778 190814. <https://doi.org/10.1098/rsos.190814>
- 779 Hoertnagl, C. M., & Hofer, A. (2014). Social cognition in serious mental illness. *Current Opinion*
780 *in Psychiatry*, 27(3), 197. <https://doi.org/10.1097/YCO.0000000000000055>
- 781 Hoffmann, J. A., Hobbs, C., Moutoussis, M., & Button, K. S. (2024). Lack of optimistic bias
782 during social evaluation learning reflects reduced positive self-beliefs in depression
783 and social anxiety, but via distinct mechanisms. *Scientific Reports*, 14(1), 22471.
784 <https://doi.org/10.1038/s41598-024-72749-6>
- 785 Hofmann, S. G. (2007). Cognitive Factors that Maintain Social Anxiety Disorder: A
786 Comprehensive Model and its Treatment Implications. *Cognitive Behaviour Therapy*,
787 36(4), 193–209. <https://doi.org/10.1080/16506070701421313>
- 788 Hopkins, A. K., Dolan, R., Button, K. S., & Moutoussis, M. (2021). A Reduced Self-Positive
789 Belief Underpins Greater Sensitivity to Negative Evaluation in Socially Anxious

790 Individuals. *Computational Psychiatry (Cambridge, Mass.)*, 5(1), 21–37.
791 <https://doi.org/10.5334/cpsy.57>

792 Hoven, M., Luijckes, J., Denys, D., Rouault, M., & van Holst, R. J. (2023). How do confidence
793 and self-beliefs relate in psychopathology: A transdiagnostic approach. *Nature Mental*
794 *Health*, 1(5), 337–345. <https://doi.org/10.1038/s44220-023-00062-8>

795 Hunter, L. E., Meer, E. A., Gillan, C. M., Hsu, M., & Daw, N. D. (2022). Increased and biased
796 deliberation in social anxiety. *Nature Human Behaviour*, 6(1), 146–154.
797 <https://doi.org/10.1038/s41562-021-01180-y>

798 Insel, T., Cuthbert, B., Garvey, M., Heinssen, R., Pine, D. S., Quinn, K., Sanislow, C., & Wang,
799 P. (2010). Research Domain Criteria (RDoC): Toward a New Classification Framework
800 for Research on Mental Disorders. *American Journal of Psychiatry*, 167(7), 748–751.
801 <https://doi.org/10.1176/appi.ajp.2010.09091379>

802 Joiner, J., Piva, M., Turrin, C., & Chang, S. W. C. (2017). Social learning through prediction
803 error in the brain. *Npj Science of Learning*, 2(1), 8. [https://doi.org/10.1038/s41539-017-](https://doi.org/10.1038/s41539-017-0009-2)
804 [0009-2](https://doi.org/10.1038/s41539-017-0009-2)

805 Kanterman, A., & Shamay-Tsoory, S. (2025). From social effort to social behavior: An
806 integrated neural model for social motivation. *Neuroscience & Biobehavioral Reviews*,
807 173, 106170. <https://doi.org/10.1016/j.neubiorev.2025.106170>

808 Kendal, R. L., & Watson, R. (2023). Adaptive Social Learning: Social Learning Strategies and
809 their Applications. In J. J. Tehrani, J. Kendal, & R. L. Kendal (Eds), *Oxford Handbook*
810 *of Cultural Evolution* (1st edn, pp. 193–207). Oxford University Press.
811 <https://doi.org/10.1093/oxfordhb/9780198869252.013.14>

812 Kendler, K. S. (2009). An historical framework for psychiatric nosology. *Psychological*
813 *Medicine*, 39(12), 1935–1941. <https://doi.org/10.1017/S0033291709005753>

814 Kendler, K. S., Zachar, P., & Craver, C. (2011). What kinds of things are psychiatric disorders?
815 *Psychological Medicine*, 41(6), 1143–1150.
816 <https://doi.org/10.1017/S0033291710001844>

817 Kirchner, L., Kube, T., Berg, M., Eckert, A.-L., Straube, B., Endres, D., & Rief, W. (2025).
818 Social expectations in depression. *Nature Reviews Psychology*, 4(1), 20–34.
819 <https://doi.org/10.1038/s44159-024-00386-x>

820 Koban, L., Andrews-Hanna, J. R., Ives, L., Wager, T. D., & Arch, J. J. (2023). Brain mediators
821 of biased social learning of self-perception in social anxiety disorder. *Translational*
822 *Psychiatry*, 13(1), Article 1. <https://doi.org/10.1038/s41398-023-02587-z>

823 Koban, L., Schneider, R., Ashar, Y. K., Andrews-Hanna, J. R., Landy, L., Moscovitch, D. A.,
824 Wager, T. D., & Arch, J. J. (2017). Social anxiety is characterized by biased learning

825 about performance and the self. *Emotion (Washington, D.C.)*, 17(8), 1144–1155.
826 <https://doi.org/10.1037/emo0000296>

827 Korn, C. W., La Rosée, L., Heekeren, H. R., & Roepke, S. (2016). Social feedback processing
828 in borderline personality disorder. *Psychological Medicine*, 46(3), 575–587.
829 <https://doi.org/10.1017/S003329171500207X>

830 Korn, C. W., Prehn, K., Park, S. Q., Walter, H., & Heekeren, H. R. (2012). Positively Biased
831 Processing of Self-Relevant Social Feedback. *Journal of Neuroscience*, 32(47),
832 16832–16844. <https://doi.org/10.1523/JNEUROSCI.3016-12.2012>

833 Kotov, R., Krueger, R. F., Watson, D., Achenbach, T. M., Althoff, R. R., Bagby, R. M., Brown,
834 T. A., Carpenter, W. T., Caspi, A., Clark, L. A., Eaton, N. R., Forbes, M. K., Forbush,
835 K. T., Goldberg, D., Hasin, D., Hyman, S. E., Ivanova, M. Y., Lynam, D. R., Markon,
836 K., ... Zimmerman, M. (2017). The Hierarchical Taxonomy of Psychopathology
837 (HiTOP): A dimensional alternative to traditional nosologies. *Journal of Abnormal*
838 *Psychology*, 126(4), 454–477. <https://doi.org/10.1037/abn0000258>

839 Kube, T. (2023). Biased belief updating in depression. *Clinical Psychology Review*, 103,
840 102298. <https://doi.org/10.1016/j.cpr.2023.102298>

841 Le Heron, C., Holroyd, C. B., Salamone, J., & Husain, M. (2019). Brain mechanisms underlying
842 apathy. *Journal of Neurology, Neurosurgery & Psychiatry*, 90(3), 302–312.
843 <https://doi.org/10.1136/jnnp-2018-318265>

844 Lockwood, P. L., Hamonet, M., Zhang, S. H., Ratnavel, A., Salmony, F. U., Husain, M., &
845 Apps, M. A. J. (2017a). Prosocial apathy for helping others when effort is required.
846 *Nature Human Behaviour*, 1(7), 0131. <https://doi.org/10.1038/s41562-017-0131>

847 Lockwood, P. L., Ang, Y.-S., Husain, M., & Crockett, M. J. (2017b). Individual differences in
848 empathy are associated with apathy-motivation. *Scientific Reports*, 7(1), 17293.
849 <https://doi.org/10.1038/s41598-017-17415-w>

850 Lundgren, D. C. (2004). Social Feedback and Self-Appraisals: Current Status of the Mead-
851 Cooley Hypothesis. *Symbolic Interaction*, 27(2), 267–286.
852 <https://doi.org/10.1525/si.2004.27.2.267>

853 Martin, R., McKay, E., & Kirk, H. (2024). Lowered social motivation is associated with
854 adolescent attention deficit hyperactivity disorder and social anxiety symptoms.
855 *Clinical Child Psychology and Psychiatry*, 29(1), 338–352.
856 <https://doi.org/10.1177/13591045231218475>

857 Morf, C. C., & Rhodewalt, F. (2001). Unraveling the Paradoxes of Narcissism: A Dynamic Self-
858 Regulatory Processing Model. *Psychological Inquiry*, 12(4), 177–196.
859 https://doi.org/10.1207/S15327965PLI1204_1

860 Müller-Pinzler, L., Czekalla, N., Mayer, A. V., Stolz, D. S., Gazzola, V., Keyzers, C., Paulus,
861 F. M., & Krach, S. (2019). Negativity-bias in forming beliefs about own abilities.
862 *Scientific Reports*, 9(1), 14416. <https://doi.org/10.1038/s41598-019-50821-w>

863 Norbury, A., Robbins, T. W., & Seymour, B. (2018). Value generalization in human avoidance
864 learning. *eLife*, 7, e34779. <https://doi.org/10.7554/eLife.34779>

865 Oka, T., Sasaki, A., & Kobayashi, N. (2025). A transdiagnostic dimensional approach to
866 behavioral dysregulation: Examining self-reported reward and punishment sensitivity
867 across psychopathology. *Journal of Affective Disorders*, 387, 119493.
868 <https://doi.org/10.1016/j.jad.2025.119493>

869 Olsson, A., Knapska, E., & Lindström, B. (2020). The neural and computational systems of
870 social learning. *Nature Reviews Neuroscience*, 21(4), Article 4.
871 <https://doi.org/10.1038/s41583-020-0276-4>

872 Patin, A., & Hurlemann, R. (2015). Social Cognition. In K. M. Kantak & J. G. Wettstein (Eds),
873 *Cognitive Enhancement* (pp. 271–303). Springer International Publishing.
874 https://doi.org/10.1007/978-3-319-16522-6_10

875 Peters, A., Helming, H., Bruchmann, M., Wiegandt, A., Straube, T., & Schindler, S. (2024).
876 How and when social evaluative feedback is processed in the brain: A systematic
877 review on ERP studies. *Cortex*, 173, 187–207.
878 <https://doi.org/10.1016/j.cortex.2024.02.003>

879 Rescorla, R. A., & Wagner, A. R. (1972). A theory of Pavlovian conditioning: Variations in the
880 effectiveness of reinforcement and non-reinforcement. *Classical Conditioning, Current*
881 *Research and Theory*, 2, 64–69.

882 Robbins, T. W., Gillan, C. M., Smith, D. G., Wit, S. de, & Ersche, K. D. (2012). Neurocognitive
883 endophenotypes of impulsivity and compulsivity: Towards dimensional psychiatry.
884 *Trends in Cognitive Sciences*, 16(1), 81–91. <https://doi.org/10.1016/j.tics.2011.11.009>

885 Rosenthal, I. A., Hutcherson, C. A., Adolphs, R., & Stanley, D. A. (2019). Deconstructing
886 Theory-of-Mind Impairment in High-Functioning Adults with Autism. *Current Biology*,
887 29(3), 513-519.e6. <https://doi.org/10.1016/j.cub.2018.12.039>

888 Rouault, M., Seow, T., Gillan, C. M., & Fleming, S. M. (2018). Psychiatric Symptom
889 Dimensions Are Associated With Dissociable Shifts in Metacognition but Not Task
890 Performance. *Biological Psychiatry, Translating Biology to Treatment in*
891 *Schizophrenia*, 84(6), 443–451. <https://doi.org/10.1016/j.biopsych.2017.12.017>

892 Rygula, R., Clarke, H. F., Cardinal, R. N., Cockcroft, G. J., Xia, J., Dalley, J. W., Robbins, T.
893 W., & Roberts, A. C. (2015). Role of Central Serotonin in Anticipation of Rewarding

894 and Punishing Outcomes: Effects of Selective Amygdala or Orbitofrontal 5-HT
895 Depletion. *Cerebral Cortex*, 25(9), 3064–3076. <https://doi.org/10.1093/cercor/bhu102>

896 Schröder, A., Czekalla, N., Mayer, A. V., Zhang, L., Stolz, D. S., Korn, C. W., Diekelmann, S.,
897 Luebber, F., Paulus, F. M., Müller-Pinzler, L., & Krach, S. (2025). Initial Expectations
898 and Confidence Affect the Formation of Novel Self-Beliefs and Their Revision. *Open*
899 *Mind*, 9, 1576–1596. <https://doi.org/10.1162/OPMI.a.36>

900 Shamay-Tsoory, S. G., & Hertz, U. (2022). Adaptive Empathy: A Model for Learning Empathic
901 Responses in Response to Feedback. *Perspectives on Psychological Science*, 17(4),
902 1008–1023. <https://doi.org/10.1177/17456916211031926>

903 Sohail, A., & Zhang, L. (2024). Informing the treatment of social anxiety disorder with
904 computational and neuroimaging data. *Psychoradiology*, 4, kkae010.
905 <https://doi.org/10.1093/psyrad/kkae010>

906 Spain, D., Sin, J., Linder, K. B., McMahon, J., & Happé, F. (2018). Social anxiety in autism
907 spectrum disorder: A systematic review. *Research in Autism Spectrum Disorders*, 52,
908 51–68. <https://doi.org/10.1016/j.rasd.2018.04.007>

909 Sutton, R. S., & Barto, A. G. (2018). *Reinforcement learning: An introduction, 2nd edition*. The
910 MIT Press.

911 Suzuki, S., Yamashita, Y., & Katahira, K. (2021). Psychiatric symptoms influence reward-
912 seeking and loss-avoidance decision-making through common and distinct
913 computational processes. *Psychiatry and Clinical Neurosciences*, 75(9), 277–285.
914 <https://doi.org/10.1111/pcn.13279>

915 Tanaka, M. (2024). Beyond the boundaries: Transitioning from categorical to dimensional
916 paradigms in mental health diagnostics. *Advances in Clinical and Experimental*
917 *Medicine*, 33(12), 1295–1301. <https://doi.org/10.17219/acem/197425>

918 Tse, W. S., & Bond, A. J. (2004). The Impact of Depression on Social Skills: A Review. *The*
919 *Journal of Nervous and Mental Disease*, 192(4), 260.
920 <https://doi.org/10.1097/01.nmd.0000120884.60002.2b>

921 Van Der Molen, M. J. W., Poppelaars, E. S., Van Hartingsveldt, C. T. A., Harrewijn, A., Gunther
922 Moor, B., & Westenberg, P. M. (2014). Fear of negative evaluation modulates
923 electrocortical and behavioral responses when anticipating social evaluative feedback.
924 *Frontiers in Human Neuroscience*, 7. <https://doi.org/10.3389/fnhum.2013.00936>

925 Vehtari, A., Gelman, A., & Gabry, J. (2017). Practical Bayesian model evaluation using leave-
926 one-out cross-validation and WAIC. *Statistics and Computing*, 27(5), 1413–1432.
927 <https://doi.org/10.1007/s11222-016-9696-4>

- 928 Weightman, M. J., Air, T. M., & Baune, B. T. (2014). A Review of the Role of Social Cognition
929 in Major Depressive Disorder. *Frontiers in Psychiatry*, 5.
930 <https://doi.org/10.3389/fpsy.2014.00179>
- 931 Wise, T., Robinson, O. J., & Gillan, C. M. (2023). Identifying Transdiagnostic Mechanisms in
932 Mental Health Using Computational Factor Modeling. *Biological Psychiatry, Emerging*
933 *Topics in Computational Psychiatric Research*, 93(8), 690–703.
934 <https://doi.org/10.1016/j.biopsych.2022.09.034>
- 935 Wong, S., Wei, G., Husain, M., Hodges, J. R., Piguet, O., Irish, M., & Kumfor, F. (2023). Altered
936 reward processing underpins emotional apathy in dementia. *Cognitive, Affective, &*
937 *Behavioral Neuroscience*, 23(2), 354–370. [https://doi.org/10.3758/s13415-022-01048-](https://doi.org/10.3758/s13415-022-01048-2)
938 2
- 939 Yang, Y., Zhou, Y., Zhang, H., Kou, H., Zhao, J., Tian, J., & Guo, C. (2024). Social anxiety
940 undermines prosocial behaviors when required effort. *International Journal of Clinical*
941 *and Health Psychology*, 24(4), 100533. <https://doi.org/10.1016/j.ijchp.2024.100533>
- 942 Zabag, R., Gilboa-Schechtman, E., & Levy-Gigi, E. (2022). Reacting to changing environment:
943 Updating patterns in social anxiety. *Behaviour Research and Therapy*, 157, 104159.
944 <https://doi.org/10.1016/j.brat.2022.104159>
- 945 Zabag, R., Rinck, M., Becker, E., Gilboa-Schechtman, E., & Levy-Gigi, E. (2024). Although I
946 know it: Social anxiety is associated with a deficit in positive updating even when the
947 cost of avoidance is Obvious. *Journal of Psychiatric Research*, 169, 279–283.
948 <https://doi.org/10.1016/j.jpsychires.2023.11.041>
- 949 Zhang, L., & Gläscher, J. (2020). A brain network supporting social influences in human
950 decision-making. *Science Advances*, 6(34), eabb4159.
951 <https://doi.org/10.1126/sciadv.abb4159>
- 952 Zhang, L., Lengersdorff, L., Mikus, N., Gläscher, J., & Lamm, C. (2020). Using reinforcement
953 learning models in social neuroscience: Frameworks, pitfalls and suggestions of best
954 practices. *Social Cognitive and Affective Neuroscience*, 15(6), 695–707.
955 <https://doi.org/10.1093/scan/nsaa089>
- 956

957 **Supplementary material**

958 **Questionnaire battery**

959 Autism Spectrum Quotient (AQ-10) (Allison et al., 2012), a 10-item screening instrument for
960 autistic traits, where higher scores indicate greater social communication difficulties and
961 rigidity. Items are rated as “Agree” or “Disagree,” and scored dichotomously (0 or 1), yielding
962 a total score range of 0–10.

963

964 Social Phobia Inventory (SPIN) (Connor et al., 2000), a 17-item measure of social anxiety
965 symptoms. Participants rate how much each statement (e.g., fear of embarrassment) applies
966 to them, with higher scores indicating greater social anxiety. Each item is rated on a 5-point
967 Likert scale (0 = “Not at all” to 4 = “Extremely”), producing a total score range of 0–68.

968

969 Beck Depression Inventory (BDI-II) (Beck et al., 1996), a 21-item inventory assessing
970 depressive symptoms such as sadness, pessimism, and anhedonia. Higher scores reflect
971 more severe depressive tendencies, and each item is rated from 0 to 3, leading to a total score
972 range of 0–63.

973

974 Apathy Motivation Index (AMI) (Ang et al., 2017), an 18-item scale measuring apathy across
975 behavioral and social domains. We used the total score as an index of general apathy (higher
976 scores = greater apathy/lower motivation). Items are rated on a 5-point Likert scale (0 =
977 “Strongly disagree” to 4 = “Strongly agree”), giving a total score range of 0–72.

978

979 Narcissistic Admiration and Rivalry Questionnaire (NARQ) (Back et al., 2013), a short 6-item
980 version of a narcissism scale, with higher scores indicating stronger narcissistic tendencies.
981 Items are rated on a 6-point scale (1 = “Disagree strongly” to 6 = “Agree strongly”), so total
982 scores range from 6 to 36.

983

984 Borderline Personality Disorder (BPD) checklist (First et al., 1997), a 10-item yes/no screening
985 checklist for BPD features, where higher scores indicate greater severity. One item with no
986 variance in our sample was removed, so scores were based on the remaining 9 items, with a
987 total score range of 0–9.

988

989 **Winning model specification**

990 The winning model implemented a hierarchical Rescorla-Wagner framework with asymmetric
 991 learning rates for reward (α^+) and punishment (α^-), alongside an inverse temperature
 992 parameter (τ) governing choice stochasticity. Parameters were estimated using Bayesian
 993 hierarchical modeling with non-centered parameterization and effect coding for within-subject
 994 condition contrasts.

995

996 The model defines the choice process by:

997

$$998 \quad \text{Choice}_{(s,t+1)} \sim \text{categorical_logit}(\tau_{(s,c)} \cdot V_{(s,t)}) \quad (1)$$

999

1000 where $V_{(s,t)}$ are action values updated according to the outcome valence:

1001

$$1002 \quad \begin{aligned} V_{(t+1)} &= V_{(t)} + \alpha_{(s,c)}^+ \cdot (R - V_{(t)}) && \text{if } R > 0 \\ V_{(t+1)} &= V_{(t)} + \alpha_{(s,c)}^- \cdot (R - V_{(t)}) && \text{if } R < 0 \end{aligned} \quad (2)$$

1003

1004 with hierarchical priors:

1005

$$1006 \quad \mu_{\alpha^+}, \mu_{\alpha^-}, \mu_{\tau} \sim \text{Normal}(0, 1)$$

$$1007 \quad \sigma_{\alpha^+}, \sigma_{\alpha^-}, \sigma_{\tau} \sim \text{Exponential}(2)$$

$$1008 \quad \varepsilon_{\alpha^+}, \varepsilon_{\alpha^-}, \varepsilon_{\tau} \sim \text{Normal}(0, 1)$$

1009

1010 **Parameter recovery**

1011 After model fitting, we confirmed the identifiability of parameters through parameter recovery.

1012 We repeated the following procedure for 20 iterations. Denoting ϕ as a generic parameter, we

1013 iterated the following steps:

1014

1015 (a) We randomly drew one joint sample of group-level parameters from the joint posterior

1016 group-level distribution of the winning model, a Rescorla-Wagner model with separate learning

1017 rates for positive and negative outcomes. That is, a group-level mean (μ_{ϕ}) and a group-level

1018 standard deviation (σ_{ϕ}) of the parameter ϕ . We repeated this procedure for all group-level

1019 parameters of the winning model:

1020

$$1021 \quad \mu_{\phi}, \sigma_{\phi} \sim p(\mu_{\phi}, \sigma_{\phi} \mid D) \quad (3)$$

1022

1023 (b) Next, we simulated 193 synthetic participants whose parameters were sampled using the
1024 hierarchical structure implemented in the Stan model. Individual-level parameters ϕ_i were
1025 generated using non-centered parameterization in probit-transformed space, then
1026 transformed to natural parameter bounds via effect coding and the cumulative normal
1027 distribution function, preserving the model's hierarchical structure and parameter constraints.
1028 We repeated this procedure for all individual-level parameters of the winning model.

1029

1030 (c) Then, we used the winning model as a generative tool to simulate behavioral data for our
1031 social evaluation learning task, namely, to simulate choices for 96 trials per participant (32
1032 trials per condition). Individuals' choices (denoted as D_i) were sampled from the sampling
1033 distribution conditional on individual-level parameters (ϕ_i) from the previous step (i.e.,
1034 likelihood function):

1035

$$1036 \quad D_i \sim p(D_i, \sigma_i) \quad (4)$$

1037

1038 (d) We fit the model to the simulated data (D_i) in the same way as we did for the real data
1039 obtaining parameter estimates (i.e., posterior distributions) at both the group-level (e.g., $\hat{\mu}_\phi$,
1040 $\hat{\sigma}_\phi$) and the individual level (e.g., $\hat{\phi}_i$).

1041

1042 (e) Finally, we compared whether the posterior distributions recovered the true data-
1043 generating parameters. At the group level, we assessed whether the true parameter values
1044 fell within the 95% highest density interval (HDI) of the recovered posterior distributions for
1045 each iteration, and report coverage rates (percentage of iterations where true value falls within
1046 95% HDI) across the 20 iterations (Supplementary Figure 1A). At the individual level, we
1047 computed Spearman's rank correlations (ρ) between the true and recovered parameters for
1048 each parameter and condition combination within each iteration, then aggregated these
1049 correlations across iterations using Fisher's Z transformation to obtain mean correlation
1050 coefficients (Supplementary Figure 1B).

1051

1052 The parameter recovery analyses initially revealed differential recoverability across individual-
1053 level parameters, with positive learning rates showing lower correlations compared to negative
1054 learning rates. Examining the fitted parameters revealed extremely low between-subject
1055 variability for the positive learning rate (CV = 0.05-0.14), whereas in contrast, the negative
1056 learning rate showed high variability (CV = 0.74-0.98) (Supplementary Table 1).

1057

Condition	α^+ mean	α^+ CV	α^- mean	α^- CV
Neutral	0.32 (0.17)	0.05	0.11 (0.89)	0.74
Like	0.44 (0.43)	0.14	0.04 (0.99)	0.98
Dislike	0.34 (0.20)	0.08	0.08 (0.70)	0.77

1058

1059 **Supplementary Table 1. Parameter estimates for the winning model.** Mean values are
1060 accompanied by the standard deviation (SD) in brackets. Mean values are group-level
1061 posterior means in natural space [0,1]. SD values are group-level posterior standard
1062 deviations in probit space. Coefficients of variation (CV) are calculated from individual-level
1063 posterior means in natural space to reflect between-subject variability.

1064

1065 To test whether low variability of this parameter caused the observed differences with
1066 recovery, we re-ran parameter recovery after artificially increased the group-level standard
1067 deviation of α^+ while keeping all other parameters at empirical values. Specifically, we
1068 calibrated α^+ probit-space SDs to achieve CVs matching α^- levels (Neutral: SD = 1.00, CV =
1069 0.74; Like: SD = 2.90, CV = 0.88; Dislike: SD = 1.20, CV = 0.78). Increasing α^+ variability
1070 subsequently improved its recovery across all three conditions (Neutral: $r = 0.24$ to $r = 0.57$;
1071 Like: $r = 0.35$ to $r = 0.68$; Dislike: $r = 0.26$ to $r = 0.57$). In contrast, α^- recovery remained stable
1072 ($r = 0.60$ - 0.66).

1073

1074 **Posterior predictive checks**

1075 To validate that the winning model could generate behavioral patterns consistent with the
1076 observed data, we conducted posterior predictive checks using a generative simulation
1077 approach. We selected 4,000 posterior samples from the winning model's MCMC chains,
1078 where for each sample, we extracted the individual-level parameters (α^+ , α^- , τ) for all 193
1079 participants across all three conditions and used the model to simulate a complete synthetic
1080 dataset of 96 trials per participant (32 trials per condition), maintaining the same task structure
1081 and feedback contingencies as in the original experiment. We then analysed each simulated
1082 dataset using the same statistical procedures applied to the observed data. Specifically, we
1083 calculated mean choice accuracy for each condition and conducted repeated-measures
1084 ANOVA with post-hoc pairwise comparisons (Tukey-corrected). To evaluate model fit, we
1085 calculated Bayesian p-values as the proportion of predicted statistics more extreme than the
1086 corresponding observed statistic (two-tailed), with p-values < 0.05 indicating circumstances
1087 where the model fails to capture the observed data pattern.

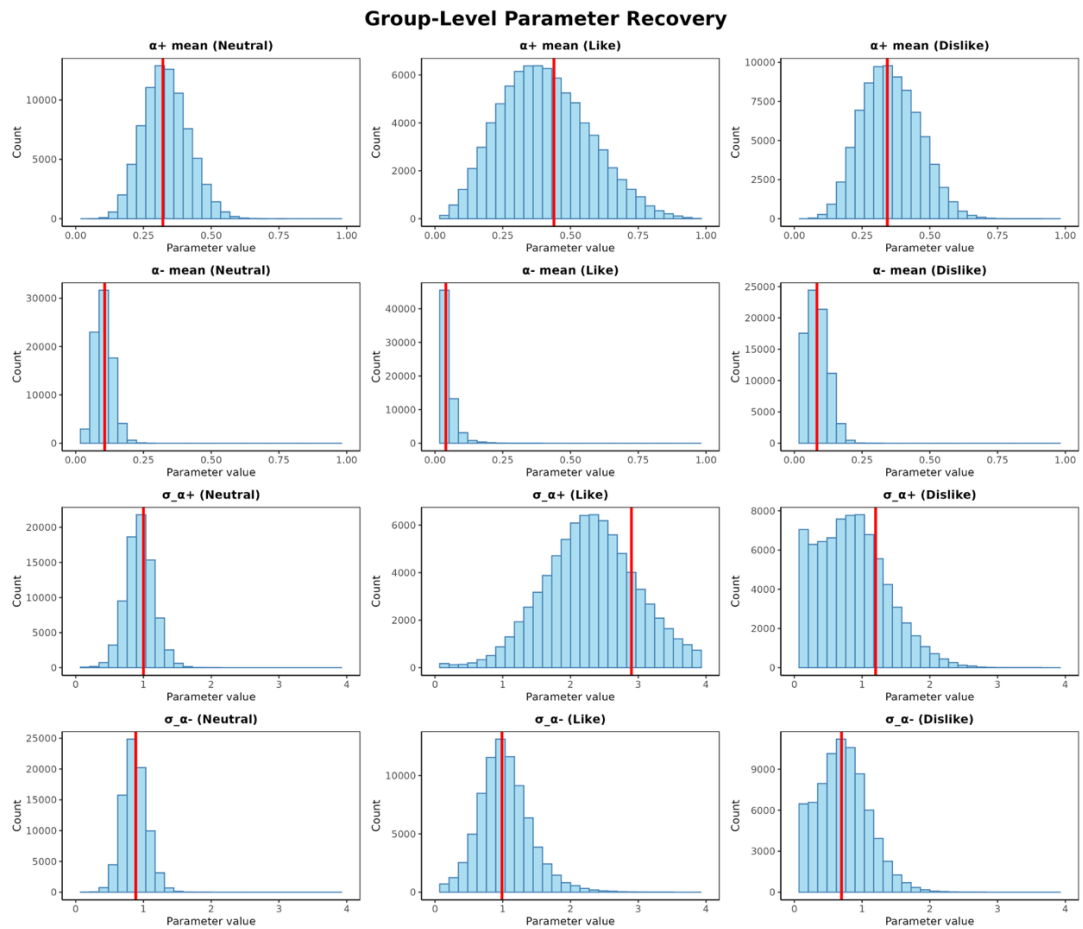
1088

1089 Posterior predictive checks (PPCs) on the empirical data showed deviations in the Like
1090 condition, where the model underpredicted accuracy. Model predictions accurately captured
1091 the Neutral (observed = 49.8%, predicted = 50.0%, 95% HDI: [48.8%, 51.3%]) and Dislike
1092 (observed = 76.2%, predicted = 76.2%, 95% HDI: [74.2%, 78.2%]), with observed values
1093 falling within the 95% prediction intervals. However, the model systematically underestimated
1094 accuracy in the Like condition (observed = 80.9%, predicted = 76.2%, 95% HDI: [73.6%,
1095 78.6%]), with the observed value falling outside the prediction interval. Subsequently, a
1096 weaker predicted condition effect (Bayesian $p = 0.005$) and smaller predicted pairwise
1097 differences were observed for Neutral vs. Like ($p = 0.001$) and Like vs. Dislike ($p = 0.003$).

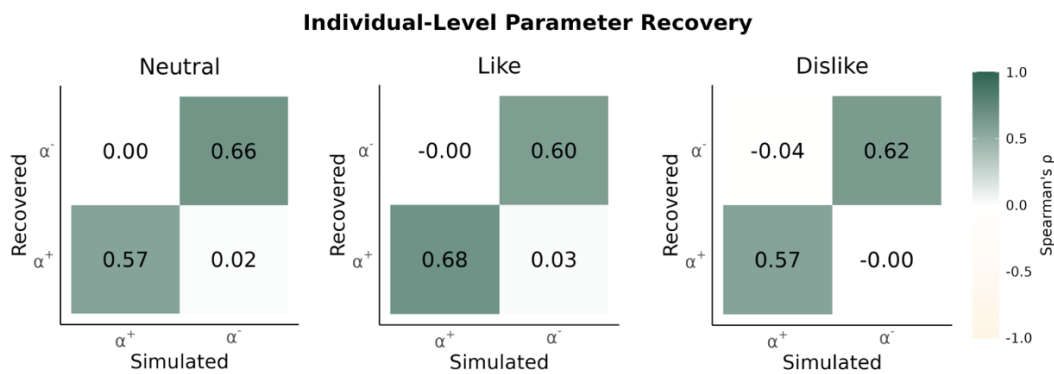
1098

1099 To test whether this also reflected low α^+ variance, we fit the model to simulated data with
1100 inflated α^+ variance as described above, and re-generated 4,000 posterior predictive datasets.
1101 Doing so significantly improved the similarity between the observed and predicted data for the
1102 Like condition (observed = 67.6%, predicted = 67.9%, 95% HDI: [65.7%, 70.0%], Bayesian p -
1103 value = 0.756). The remaining conditions were also well-recovered: Neutral (observed =
1104 50.0%, predicted = 50.0%, 95% HDI: [47.9%, 52.2%], $p = 0.982$) and Dislike (observed =
1105 70.7%, predicted = 70.8%, 95% HDI: [68.9%, 72.7%], $p = 0.863$), with all observed values
1106 falling within their respective 95% prediction intervals. The model also accurately captured the
1107 overall condition effect (observed $F = 94.17$, predicted mean $F = 101.61$, 95% HDI: [71.73,
1108 136.80], $p = 0.664$) and all pairwise differences between conditions (all Bayesian p -values $>$
1109 0.75; Supplementary Figure 2).

A



B



1110

1111 **Supplementary Figure 1. Parameter recovery for the winning computational model.**

1112 Group-level mean parameters were drawn from the empirical posterior of the winning model,

1113 with between-subject variance for α^+ inflated to match empirically plausible levels ($CV \approx 0.74$ –

1114 0.98). The data-generating model was then fitted to synthetic behavioral data from 193

1115 simulated participants, and parameter estimates were compared to the true data-generating

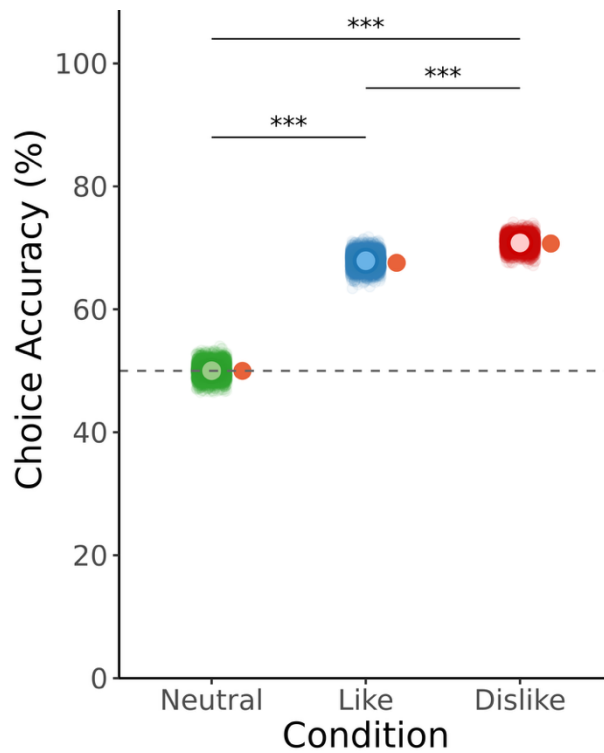
1116 values. **(A)** At the group level, true mean and group-level standard deviation parameters fell

1117 within the 95% HDI of recovered posterior distributions across all conditions. Note that SD

1118 parameters are reported in probit space. **(B)** At the individual level, parameter recovery

1119 showed good identifiability for both the negative learning rate ($p = 0.60-0.66$) and positive
1120 learning rate ($p = 0.57-0.68$)."

1121



1122

1123 **Supplementary Figure 2. Posterior predictive checks of the winning model.** Each
1124 condition has the distribution of predicted mean choice accuracy across 4,000 posterior
1125 predictive datasets plotted per condition (Neutral = green, Like = blue, Dislike = red), within
1126 the mean predicted accuracy. The orange dot denotes the observed mean accuracy from the
1127 simulated dataset, representing the empirical estimate. Significance brackets reflect pairwise
1128 differences in the predicted data. *** $p < .001$

1129

1130 **Subjective ratings and learned value**

1131 To verify that participants' explicit social appraisals tracked the feedback contingencies across
1132 conditions, we examined whether model-derived learned values predicted self-reported
1133 ratings of how much the feedback agent liked them. This served as a manipulation check,
1134 confirming that the computational learning signals captured by the model corresponded to
1135 participants' subjective experience of the social feedback.

1136

1137 For each participant and condition, we reconstructed the trial-by-trial trajectory of V^{positive} - the
1138 learned value of the positive-word option - by replaying the Rescorla-Wagner update rule

1139 forward using each participant's posterior mean α^+ and α^- estimates. Specifically, V^{positive} was
 1140 initialised to zero at the start of each block and updated on each valid trial according to:

1141

$$1143 \quad V_t^{\text{positive}} = V_{t-1}^{\text{positive}} + \alpha^\pm \cdot (r_t - V_{t-1}^{\text{positive}}) \quad (5)$$

1142

1144 where α^+ was applied following positive outcomes and α^- following negative outcomes. The
 1145 mean V^{positive} across all valid trials within each block was taken as the learned social value for
 1146 that condition. This quantity reflects the degree to which participants learned that the agent
 1147 responded positively to their choices: values approaching +1 indicate strong learned
 1148 association between positive choices and positive feedback, while negative values indicate
 1149 the inverse.

1150

1151 Explicit ratings were collected via a continuous slider (0–100) at the end of each block, with
 1152 participants indicating how much they felt the feedback agent liked them. We computed
 1153 Kendall's τ correlations between mean V^{positive} and ratings separately for each condition.

1154

1155 **Supplementary tables**

1156

Predictor	Choice Accuracy	Reaction Time	Rating
Intercept	0.290* (0.122)	7.423*** (0.066)	36.607*** (3.848)
Like	1.049*** (0.164)	-0.105 (0.074)	29.071*** (4.324)
Dislike	1.080*** (0.140)	-0.147* (0.059)	-18.017*** (4.221)
α^+	4.935* (2.230)	2.247 (1.260)	-64.323 (75.272)
α^-	-0.117 (0.472)	0.510** (0.184)	-9.918 (10.725)
Depressive Self-Identity	0.008 (0.122)	-0.047 (0.078)	2.074 (4.604)
Social Avoidance	-0.022 (0.128)	0.032 (0.077)	-5.060 (4.606)
Emotional Insensitivity	0.243* (0.119)	0.042 (0.069)	-0.711 (4.056)
Like x α^+	-0.697 (2.316)	-2.226 (1.203)	94.351 (78.096)
Dislike x α^+	-1.421 (2.826)	-1.658 (1.144)	5.432 (89.455)
Like x Depressive Self-Identity	0.033 (0.177)	0.026 (0.087)	-1.935 (5.132)
Dislike x Depressive Self-Identity	0.045 (0.136)	0.058 (0.070)	0.042 (4.987)
Like x Social Avoidance	0.013 (0.183)	-0.011 (0.086)	2.736 (5.218)
Dislike x Social Avoidance	0.076 (0.144)	-0.027 (0.070)	3.409 (4.968)
Like x Emotional Insensitivity	-0.480** (0.167)	-0.166* (0.079)	-2.841 (4.615)
Dislike x Emotional Insensitivity	-0.323* (0.139)	-0.109 (0.062)	3.621 (4.440)
α^+ x Depressive Self-Identity	-0.698 (2.082)	-0.159 (1.482)	24.555 (89.024)
α^+ x Social Avoidance	1.576 (2.432)	0.451 (1.541)	-111.455 (94.878)
α^+ x Emotional Insensitivity	3.414 (2.097)	2.336 (1.320)	-28.313 (78.797)

Like x α^-	-0.046 (0.678)	-0.251 (0.189)	16.568 (14.527)
Dislike x α^-	0.963 (0.653)	0.041 (0.183)	8.668 (14.948)
Depressive Self-Identity x α^-	-0.492 (0.521)	0.378 (0.210)	-10.270 (12.059)
Social Avoidance x α^-	0.897 (0.528)	-0.436* (0.214)	19.899 (12.408)
Emotional Insensitivity x α^-	-0.328 (0.490)	0.179 (0.197)	6.092 (11.392)
Like x α^+ x Depressive Self-Identity	1.259 (2.188)	0.302 (1.427)	-44.807 (92.670)
Dislike x α^+ x Depressive Self-Identity	0.366 (2.445)	0.768 (1.338)	61.081 (104.559)
Like x α^+ x Social Avoidance	-2.297 (2.622)	-0.939 (1.496)	153.429 (98.570)
Dislike x α^+ x Social Avoidance	2.118 (3.014)	0.342 (1.438)	66.422 (111.565)
Like x α^+ x Emotional Insensitivity	-3.126 (2.324)	-1.688 (1.264)	42.010 (82.483)
Dislike x α^+ x Emotional Insensitivity	-2.417 (2.698)	-1.597 (1.189)	20.662 (92.636)
Like x Depressive Self-Identity x α^-	0.709 (0.768)	-0.079 (0.217)	13.631 (16.758)
Dislike x Depressive Self-Identity x α^-	-0.434 (0.723)	-0.102 (0.205)	17.190 (16.849)
Like x Social Avoidance x α^-	-0.737 (0.738)	0.567** (0.214)	-20.751 (16.307)
Dislike x Social Avoidance x α^-	-0.517 (0.767)	0.297 (0.215)	-23.714 (17.618)
Like x Emotional Insensitivity x α^-	0.089 (0.712)	0.134 (0.202)	-16.903 (15.501)
Dislike x Emotional Insensitivity x α^-	0.369 (0.710)	0.083 (0.200)	-4.951 (16.364)

1157

1158 **Supplementary Table 2. Regression coefficients predicting behavioral task**

1159 **performance.** Models consisted of a fixed-effects logistic regression for choice accuracy,

1160 linear mixed-effects model for reaction time, and standard linear regression for likeability

1161 ratings. All models included condition (Like, Dislike relative to Neutral), transdiagnostic factor

1162 scores and their interactions as predictors. Mixed-effects models included random intercepts

1163 for subjects. Model formula: DV ~ condition * α_{pos} * (fa1 + fa2 + fa3) + condition *

1164 α_{neg} * (fa1 + fa2 + fa3) + (1 + condition | Subject). Coefficients for factor scores represent

1165 semi-standardized estimates (standardized predictors, unstandardized outcomes) with

1166 standard errors in parentheses. * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

1167

Predictor	α^+	α^-	α^{diff}
Intercept	0.325*** (0.003)	0.173*** (0.009)	0.152*** (0.010)
Like	0.119*** (0.003)	-0.036*** (0.006)	0.155*** (0.007)
Dislike	0.025*** (0.003)	0.000 (0.006)	0.025*** (0.007)
Depressive Self-Identity	-0.001 (0.003)	-0.001 (0.010)	0.000 (0.011)
Social Avoidance	0.001 (0.003)	0.025* (0.010)	-0.024* (0.011)
Emotional Insensitivity	-0.001 (0.003)	-0.011 (0.009)	0.010 (0.010)
Like x Depressive Self-Identity	0.003 (0.004)	0.002 (0.007)	0.001 (0.008)
Dislike x Depressive Self-Identity	-0.001 (0.004)	-0.005 (0.007)	0.003 (0.008)
Like x Social Avoidance	0.002 (0.004)	-0.006 (0.007)	0.008 (0.008)
Dislike x Social Avoidance	-0.001 (0.004)	0.001 (0.007)	-0.002 (0.008)
Like x Emotional Insensitivity	-0.002 (0.003)	0.006 (0.006)	-0.008 (0.007)

Dislike x Emotional Insensitivity 0.000 (0.003) -0.000 (0.006) 0.000 (0.007)

1168

1169 **Supplementary Table 3. Regression coefficients predicting computational parameters**

1170 **from the winning model.** The table shows results from three separate linear mixed-effects

1171 regression models with condition (Like/Dislike relative to Neutral), transdiagnostic factor

1172 scores (Depressive Self-Identity, Social Avoidance, Emotional Insensitivity), and their

1173 interactions predicting each computational parameter from the winning reinforcement learning

1174 model. Outcome variables across the four models were positive learning rate (α^+), negative

1175 learning rate (α^-) and positive learning bias ($\alpha^+ - \alpha^-$). Model formula: $DV \sim \text{condition} \times (\text{fa1} +$

1176 $\text{fa2} + \text{fa3}) + (1|\text{subID})$. Values reflect semi-standardized regression coefficients (standardized

1177 predictors, unstandardized outcomes) with standard errors in parentheses. * $p < 0.05$, ** $p <$

1178 0.01 , *** $p < 0.001$

1179

Item	Depressive Self-Identity	Social Avoidance	Emotional Insensitivity
NARQ1			
NARQ2			
NARQ3			
NARQ4			
NARQ5			
NARQ6			
SPIN1		0.539	
SPIN2		0.545	
SPIN3		0.511	
SPIN4		0.7	
SPIN5		0.624	
SPIN6		0.623	
SPIN7		0.433	
SPIN8		0.557	
SPIN9		0.441	
SPIN10		0.647	
SPIN11		0.703	
SPIN12		0.608	
SPIN13		0.751	
SPIN14		0.414	
SPIN15		0.426	
SPIN16		0.479	
SPIN17		0.668	
AMI1			0.423
AMI2		0.54	
AMI3		0.415	

AMI4		0.317	0.302
AMI5		0.458	
AMI6			
AMI7			0.476
AMI8		0.323	
AMI9			
AMI10			
AMI11			
AMI12			
AMI13			0.657
AMI14		0.497	
AMI15			
AMI16			0.549
AMI17		0.429	
AMI18			0.579
AQ1			
AQ2			
AQ3			
AQ4			
AQ5		0.312	0.427
AQ6			0.472
AQ7			
AQ8			
AQ9			0.386
AQ10			0.301
BDI1	0.629		
BDI2	0.54		
BDI3	0.672		
BDI4	0.632		
BDI5	0.551		
BDI6	0.615		
BDI7	0.732		
BDI8	0.638		
BDI9	0.575		
BDI10	0.652		
BDI11	0.649		
BDI12	0.646		
BDI13	0.646		
BDI14	0.544		
BDI15	0.533		
BDI16	0.608		
BDI17	0.651		
BDI18	0.514		
BDI19			

BDI20	0.577
BDI21	0.493
BPD1	0.305
BPD2	0.341
BPD3	
BPD5	0.409
BPD6	0.378
BPD7	0.502
BPD8	
BPD9	0.437

1180 **Supplementary Table 4. Factor loadings from the exploratory factor analysis three-**
1181 **factor solution.** Only loadings ≥ 0.30 are displayed, consistent with standard reporting
1182 thresholds for exploratory factor analysis.