

1 **Title:**

2 Computational and neural mechanisms of affected beliefs

3 **Authors:**

4 Laura Müller-Pinzler¹, Nora Czekalla¹, Annalina V Mayer¹, Alexander Schröder¹, David S

5 Stolz¹, Frieder M Paulus¹, Sören Krach¹

6 **Email:**

7 Laura Müller-Pinzler* laura.muellerpinzler@uni-luebeck.de

8 Nora Czekalla n.czekalla@uni-luebeck.de

9 Annalina V Mayer ann.mayer@uni-luebeck.de

10 Alexander Schröder a.schroeder@uni-luebeck.de

11 David S Stolz david.stolz@uni-luebeck.de

12 Frieder M Paulus frieder.paulus@uni-luebeck.de

13 Sören Krach soeren.krach@uni-luebeck.de

14 **Affiliations:**

15 1: Department of Psychiatry and Psychotherapy, Social Neuroscience Lab, University of

16 Lübeck, Ratzeburger Allee 160, D-23538 Lübeck, Germany

17 **Short Title:**

18 ***Corresponding Author:**

19 Dr. Laura Müller-Pinzler

20 Department of Psychiatry and Psychotherapy, Social Neuroscience Lab

21 University of Lübeck, Ratzeburger Allee 160, D-23538 Lübeck, Germany

22 Phone: +49 45131017529

23

24

Abstract

25 The feedback people receive on their behavior shapes the process of belief formation
26 and self-efficacy in mastering a given task. The neural and computational mechanisms of how
27 the subjective value of these beliefs and corresponding affect bias the learning process are yet
28 unclear. Here we investigate this question during learning of self-efficacy beliefs using fMRI,
29 pupillometry, computational modeling and individual differences in affective experience.
30 Biases in the formation of self-efficacy beliefs were associated with affect, pupil dilation and
31 neural activity within the anterior insula, amygdala, VTA/SN, and mPFC. Specifically, neural
32 and pupil responses map the valence of the prediction errors in correspondence to the
33 experienced affect and learning bias people show during belief formation. Together with the
34 functional connectivity dynamics of the anterior insula within this network our results hint
35 towards neural and computational mechanisms that integrate affect in the process of belief
36 formation.

37

38

39

Introduction

40 Self-efficacy expectation is a person's subjective conviction that he or she can
41 overcome challenging situations by own actions (Bandura, 1977). To successfully perform
42 goal directed actions, humans must learn from incoming information and thereby form beliefs
43 about the world and about themselves enmeshed in this world. According to economic theory,
44 learning should result in accurate beliefs representing an internal model of the world that is
45 suitable to inform decision making. Novel theoretical frameworks however emphasize that
46 besides its instrumentality (i.e. being accurate) the belief itself may carry intrinsic value
47 (Bromberg-Martin and Sharot, 2020) thereby shaping the learning process and how people
48 ultimately arrive at their beliefs (Sharot and Garrett, 2016). Here, affective states, e.g.
49 happiness about one's own good health prognosis, represent intrinsic values that individuals
50 are inclined to optimize during belief formation (Bromberg-Martin and Sharot, 2020; Hughes
51 and Zaki, 2015). To prove this entanglement of affect and belief formation, we applied a
52 learning task that induces affective reactions during the process of forming conceptually novel
53 beliefs about own abilities in mastering a task (Czekalla et al., 2021; Müller-Pinzler et al.,
54 2019). We focused on the primary affective states elicited by self-related beliefs – the self-
55 conscious emotions of embarrassment and pride – and their impact on the belief. By having
56 experimental control over failures and successes during the process of belief formation, we
57 were able to assess how experienced affect relates to computational mechanisms of belief
58 formation and the underlying neural systems activity, explaining the shifts of preferences for
59 information of positive or negative valence during learning.

60 Affective states are considered to guide cognitive processing, representing embodied
61 and experiential information about the positive or negative value of what people encounter
62 (Frijda, 1987; Storbeck and Clore, 2008). These internal affective information are proposed to
63 be integrated with external information shaping beliefs that, rather than being objective, are
64 motivated and biased by subjective feelings towards the belief itself, which constitutes a

65 recursive influence of beliefs and affective states on each other (Bromberg-Martin and Sharot,
66 2020; Loewenstein, 2006). Previous studies support aspects of the Bromberg-Martin & Sharot
67 framework by showing that internal beliefs and external feedback can elicit emotions like
68 happiness, pride, or embarrassment (Müller-Pinzler et al., 2015; Rutledge et al., 2014, 2016;
69 Stolz et al., 2020). Affective states also alter decision making (Charpentier et al., 2016a,
70 2016b; Stolz et al., 2020) and cognitive processes like situational judgments or learning styles
71 (Storbeck and Clore, 2008). Social anxiety, low self-esteem, or depression, which are likely
72 associated with more negative affective reactions towards self-related beliefs, also bias belief
73 formation (Koban et al., 2017; Korn et al., 2014; Müller-Pinzler et al., 2019; Will et al., 2020)
74 supporting the overall rationale of the formation of “affected beliefs”, that is, the notion that
75 beliefs are fundamentally shaped by affective experiences. The question however remains
76 which neurophysiological mechanisms can explain how emotions brought up during learning
77 are associated with biases in belief formation.

78 Neuroscientific studies provided initial evidence that common brain areas map the
79 value of stimuli, actions, and their motivational relevance during social and non-social learning
80 and decision making (Chib et al., 2009; Ruff and Fehr, 2014). Prediction errors, that is, the
81 mismatch of prior expectation and a situation’s outcome, which are being minimized by
82 updating beliefs during learning, are generally processed in the dopaminergically innervated
83 ventral striatum, but also in the orbitofrontal cortex or the amygdala during learning (King-
84 Casas et al., 2005; O’Doherty, 2004; Ruff and Fehr, 2014; Schultz et al., 1997). However,
85 more recent findings suggest that there are distinct and unique neural computations potentially
86 reflecting the impact of motivational and emotional processes that are prominent during belief
87 formation. For example, studies could show that distinct value-related neural processes in
88 subregions of the anterior cingulate cortex (ACC) are recruited depending on whether
89 information about oneself or another agent is processed (Lockwood and Wittmann, 2018;
90 Lockwood et al., 2016). In other studies, activity in the ventral striatum was modulated if the

91 social context changed from a private to a public situation, suggesting that the presence or
92 absence of others influenced sensitivity to the reward value of certain decisions (Izuma et al.,
93 2010). Biases specific for self-related belief updating that are absent when people learn about
94 another person (Kuzmanovic et al., 2016), have been related to differences in tracking of
95 negative prediction errors (Sharot et al., 2011). Here, the ventromedial prefrontal cortex
96 (vmPFC) shows valence specific encoding of self-related information, which has been shown
97 to predict optimistic biases in updating (Kuzmanovic et al., 2016, 2018).

98 Affective states triggered after personal failures or successes are particularly important
99 when people acquire novel self-concepts (Hopkins et al., 2021) and develop an initial
100 understanding of themselves as being self-efficacious individuals in a novel task environment.
101 Central to the entanglement of affect and such self-efficacy beliefs is the assumption that
102 people are highly motivated to perform well and maintain or even construct a positively shaped
103 self-image (Markus and Wurf, 1987; Sedikides and Gregg, 2008). Within this process,
104 performance feedback elicits self-conscious emotions such as pride in case of success (Stolz et
105 al., 2020; Tangney et al., 2007; Williams and DeSteno, 2008), but also embarrassment if one
106 fails to achieve the expected outcome (Miller, 1996; Müller-Pinzler et al., 2015; Tangney et
107 al., 2007). It has already been shown that these emotions are not only a consequence of the
108 situation but also directly affect behavior. Pride experiences, associated with increased
109 functional coupling between the ventral striatum and cortical midline structures (Stolz et al.,
110 2020), thus functions as a motivator to continue one's effort (Williams and DeSteno, 2008). In
111 contrast, embarrassment experiences, as signified by increased functional connectivity between
112 brain areas involved in "Theory of Mind" (Kanske et al., 2015) and arousal processing systems
113 within the ventral anterior insula (vAI) and amygdala (Müller-Pinzler et al., 2015), rather lead
114 individuals to stop the current behavior, withdraw, and appease others (Apsler, 1975; Feinberg
115 et al., 2012). For the process of belief formation, it has been argued that specifically the dorsal
116 mediofrontal cortex (dmMFC), the ventral and dorsal anterior insula (vAI/ dAI), and amygdala,

117 brain areas involved during action monitoring as well as emotional processing, integrate
118 affective states with outcome information (Koban and Pourtois, 2014). Therefore, among
119 others, the AI has been regarded as an integrative hub for motivated cognition and emotional
120 behavior (Koban and Pourtois, 2014; Wager and Feldman Barrett, 2017). Similarly,
121 dopaminergic midbrain nuclei in the ventral tegmental area and substantia nigra (VTA/ SN)
122 are associated with attention processes and at the same time events (i.e. reward cues) that are
123 of motivational significance specifically during learning (Adcock et al., 2006; Schultz, 1998).

124 While current frameworks support the idea that intrinsic outcomes such as affective
125 states impact the process of belief formation (Bromberg-Martin and Sharot, 2020; Hughes and
126 Zaki, 2015), studies on this issue currently do not probe this framework as a whole. We aim
127 to bridge this gap by showing how emotional states relate to biases in self-related beliefs, that
128 is, the formation of self-efficacy and how they shift preferences for information of positive or
129 negative valence. To do so, we tested the effects of individual differences in affective reactions
130 to the task and learning behavior. Affective states were evoked during of learning self-efficacy
131 beliefs in a conceptually novel task environment. Using trial-by-trial updates of performance
132 expectations, we computed prediction error learning rates by fitting computational learning
133 models revealing valence specific learning biases. As predicted by current frameworks,
134 individual differences in the experience of the emotions embarrassment and pride were
135 distinctively related to biases in learning of self-related beliefs. Biased belief updating and
136 affect were jointly related to neural processing of valence specific prediction errors in the AI,
137 amygdala, VTA/ SN and mPFC as well as pupillary reactivity in favor of information that is
138 preferably used to update the belief. Increases in valence specific functional connectivity of the
139 dAI with the amygdala, VTA/ SN and mPFC support the idea of an integrative mechanism of
140 affective and attentional processes within the dAI. These findings provide insights into brain
141 networks involved in computational biases shaped by emotional experiences and coherently
142 support current theoretical frameworks integrating affective experiences in the process of belief

143 formation.

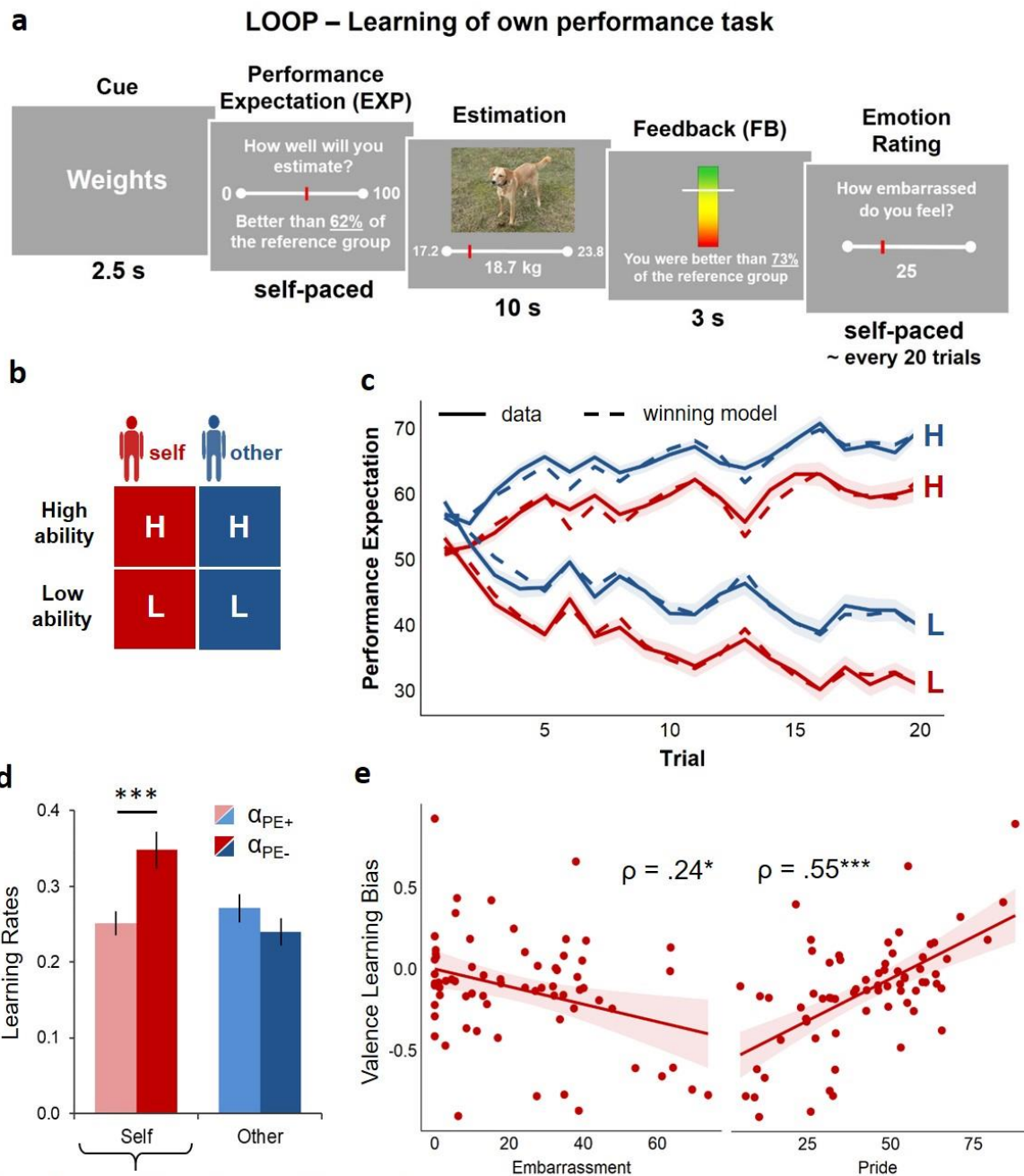
144

145 **Results**

146 **Experimental design**

147 Of the participants in our overall sample more than half completed the task in the MRI
148 while additionally eye-tracking data were assessed during scanning, the other half completed
149 the task as a behavioral study (see methods section for details). In our self-efficacy learning
150 experiment, the “learning of own performance” (LOOP) task (Müller-Pinzler et al., 2019),
151 participants were repeatedly confronted with manipulated feedback on their own and another
152 person’s cognitive estimation ability. In different domains (e.g., estimating the weight of
153 animals) participants were led to form novel beliefs about their performance-related self-
154 efficacy. We invited each participant together with a confederate to take part in the “cognitive
155 estimation” experiment in our neuroimaging lab. The participant performed the task in the
156 MRI scanner, while the confederate (introduced as another participant) was located in an
157 adjacent room to simultaneously perform the task there. During the task, participants were
158 asked to estimate specific characteristics of different properties (e.g. heights of buildings or
159 weights of animals). After each trial, they received a manipulated performance feedback for
160 their last estimation (see **Figure 1a**). During the entire experiment participants took turns in
161 performing the estimation task themselves (“Self” condition) or observing the other
162 participant performing (“Other” condition). At the beginning of each trial participants were
163 requested to rate either their own or the other person’s expected performance for the upcoming
164 trial, enabling us to examine the process of self- and other-related belief formation. The design
165 of the LOOP task provided a High Ability and a Low Ability condition which resulted in
166 overall four feedback conditions: Agent condition (Self vs Other) x Ability condition (High
167 Ability vs Low Ability; see **Figure 1b** and methods for a detailed description of the task). In
168 previous studies we showed that over time participants adapted their expected performance

169 ratings according to the feedback allowing for an assessment of valence specific self- and
 170 other-related learning processes (Czekalla et al., 2021; Müller-Pinzler et al., 2019).



171

$$\text{Valence Learning Bias} = (\alpha_{PE+} - \alpha_{PE-}) / (\alpha_{PE+} + \alpha_{PE-})$$

172 **Figure 1.** Trial sequence and timing, experimental conditions, modeling of learning behavior,
 173 learning rates and their association with self-conscious emotions. **a)** A cue in the beginning of
 174 each trial indicated the following estimation category as well as the agent who's turn it was.
 175 After providing their performance expectation ratings (EXP) participants received an
 176 estimation question, followed by the corresponding performance feedback. After
 177 approximately every 20 trials participants were asked to rate their current emotional state
 178 (pride, embarrassment, happiness, stress/ arousal). **b)** The LOOP task contained two
 179 experimental factors, Ability level (High ability vs Low ability) and Agent (Self vs Other),
 180 resulting in four feedback learning conditions that are distinguishable via different estimation

181 question types (e.g. estimation of weights of animals or heights of buildings). **c)** Predicted and
182 actual performance expectation ratings across time. The behavioral data indicate that
183 participants adapted their performance expectation ratings (solid lines) to the provided
184 feedback, thus learning about their allegedly distinct performance levels in the two ability
185 conditions for themselves and the other person. The winning valence specific learning model
186 captured the participants' behavior as indicated by a close match of actual performance
187 expectations with the predicted data (dashed lines). Shaded areas represent the standard errors
188 for the actual performance expectations for each trial. **d)** Learning rates derived from the
189 winning Valence Model indicate that there was a bias towards increased updating in response
190 to negative prediction errors (α_{PE-}) in contrast to positive prediction errors (α_{PE+}) for self-
191 related learning. Bars represent mean learning rates, error bars depict +/- one standard error;
192 *** = $p < .001$, indicates a significant negativity bias during self-related learning. **e)**
193 Correlation plots and spearman correlations of self-related Valence Learning Bias and
194 embarrassment as well as pride experience during the experiment. * = $p < .05$, *** = $p < .001$.
195

196 **Model free behavioral analyses reveal more negative self-evaluation**

197 We first performed a model free analysis to capture the basic effects we observed in
198 our behavioral data. Analyses of behavioral data and learning rates are based on the combined
199 fMRI ($n=39$) and behavioral sample ($n=30$; overall $N=69$). The Trial x Ability condition x
200 Agent condition x Group ANOVA revealed a significant main effect of Ability condition
201 ($F_{(1,67)}=175.51, p<.001$) and interaction of Trial x Ability condition ($F_{(19,1273)}=108.87, p<.001$)
202 indicating that participants adapted their performance expectation ratings over time according
203 to the feedback provided in each Ability condition (see **Figure 1c**). There was a significant
204 main effect of Agent condition ($F_{(1,67)}=44.70, p=.001$) which indicated that participants
205 evaluated their own performance more negatively than the other's performance. There was no
206 significant interaction of Agent condition x Ability condition ($F_{(1,67)}=0.67, p=.415$). The three-
207 way interaction of Trial x Agent condition x Ability condition ($F_{(19,1273)}=1.60, p=.047$) showed
208 a significant effect hinting towards differential learning patterns between the Ability
209 conditions for Self vs Other. There was a significant main effect of Group ($F_{(1,67)}=4.32,$
210 $p=.041$) indicating slightly higher ratings in the fMRI sample, but there was no interaction of
211 Group with any of the effects reported above ($p>.174$).

212 **Selection of computational models for self-related learning**

213 In a next step, we modeled the participants' behavior by means of learning rates.

214 Therefore, we used the trial-specific expectation ratings for Self and Other to model prediction
215 error (PE) update learning and assess differences between updating behavior for information
216 of positive vs negative valence. Our model space contained different models allowing us to
217 assess the importance of valence specific learning rates in contrast to unbiased learning
218 between conditions and agents (**Figure S1**). In line with our previous studies, an extended
219 version of the Valence Model, including separate learning rates for positive and negative PEs
220 for Self vs Other was the winning model (Model 8; for a more detailed description of this
221 model and the whole model space see methods section). It received the highest sum PSIS-
222 LOO score (approximate leave-one-out cross-validation (LOO) using Pareto-smoothed
223 importance sampling (PSIS)) (Vehtari et al., 2016) out of all models (for all PSIS-LOO scores
224 see **Supplementary Table S1**). In addition, Bayesian model selection (BMS) resulted in a
225 protected exceedance probability of $pxp > .999$ for this model and a Bayesian Omnibus Risk
226 of $BOR < .001$. Thus, the extended Valence Model was selected for all further analyses of
227 learning parameters allowing for a comparison of valence specific learning rates. The time
228 courses of performance expectation ratings predicted by our winning model successfully
229 captured trial-by-trial changes in the actual expectations due to PE updates within each of the
230 ability conditions at the individual subject level ($R^2 = 0.46 \pm 0.28$ [$M \pm SD$]) supporting the
231 validity of the model in describing the subjects' learning behavior. In addition to revealing PE
232 valence specific learning, which could not directly be assessed via model free behavioral
233 analyses, posterior predictive checks also confirmed that the winning model captured the core
234 effects in our model free analysis (see **Supplementary Results** and **Figure 1c**).

235 **Replication of the negativity bias for self-related learning**

236 Participants showed higher learning rates when learning about themselves compared
237 to learning about another person (main effect of Agent: $F_{(1,67)} = 5.77$, $p = .019$). There was also
238 a main effect of PE valence [pos| neg] ($F_{(1,67)} = 5.22$, $p = .025$; indicating the categorical
239 comparison of learning rates for positive vs negative PEs) and a significant interaction of

240 Agent x PE valence [pos| neg] ($F_{(1,67)}=21.47, p<.001$) which replicates earlier findings of a
241 bias towards more negative updating when learning about one's own performance ($t_{(68)}=-4.85,$
242 $p<.001, M\alpha_{PE+Self}=0.25, SD=0.13; M\alpha_{PE-Self}=0.35, SD=0.20$) (Müller-Pinzler et al., 2019).
243 Learning about the other person's performance did not show a significant bias towards more
244 negative updating ($t_{(68)}=1.53, p=.128; M\alpha_{PE+Other}=0.27, SD=0.16; M\alpha_{PE-Other}=0.24, SD=0.15;$
245 see **Figure 1d**). There was no significant main effect or interaction for Group ($p>.097$).

246 **Associations of self-related learning with self-conscious emotions.**

247 Individual differences in the overall experience of embarrassment and pride during the
248 task were used as between-subject measures to quantify associations between learning
249 behavior and affect. Embarrassment and pride ratings were only weakly correlated ($\rho_{(68)}=-.10,$
250 $p=.436$), indicating that the experience of embarrassment and pride during the task represent
251 two rather independent affective components with respect to the self-related feedback. The
252 bias in self-related learning (Valence Learning Bias= $(\alpha_{Self/PE+} - \alpha_{Self/PE-})/(\alpha_{Self/PE+} + \alpha_{Self/PE-})$)
253 (Müller-Pinzler et al., 2019; Niv et al., 2012; Palminteri et al., 2017) was negatively linked to
254 the reported experience of embarrassment during the task ($\rho_{(68)}=-.24, p=.043$; fMRI
255 subsample: $\rho_{(38)}=-.44, p=.005$), that is, more negative updating behavior was associated with
256 increased embarrassment ratings. In contrast, the Valence Learning Bias was positively linked
257 to the emotion of pride ($\rho_{(68)}=.55, p<.001$; fMRI subsample: $\rho_{(38)}=.47, p=.002$). A regression
258 predicting the Valence Learning Bias with both affect ratings simultaneously showed
259 independent effects of pride ($\beta=0.56, t_{(66)}=5.81, p<.001$) and embarrassment ($\beta=-0.22, t_{(66)}=-$
260 $2.30, p=.025; R^2=.41, F_{(1,66)}=22.90, p<.001$). This indicates that the experience of self-
261 conscious emotions during successful and unsuccessful performances are tied to the way
262 people updated their self-efficacy beliefs (see **Figure 1e**). Further, the way participants
263 processed the performance feedback in order to update their self-related ability beliefs was
264 associated with their self-esteem. People with higher self-esteem showed more positive
265 updating, $\rho_{(68)}=.33, p=.006$ (fMRI subsample: $\rho_{(38)}=.35, p=.030$), which strengthens the

266 assumption that prior beliefs about the self have a direct impact on how individuals learn novel
267 information about new abilities (Müller-Pinzler et al., 2019; Rouault et al., 2019).

268 **Pupil dilation slopes are associated with surprise and valence of prediction errors in line**
269 **with a negative learning bias**

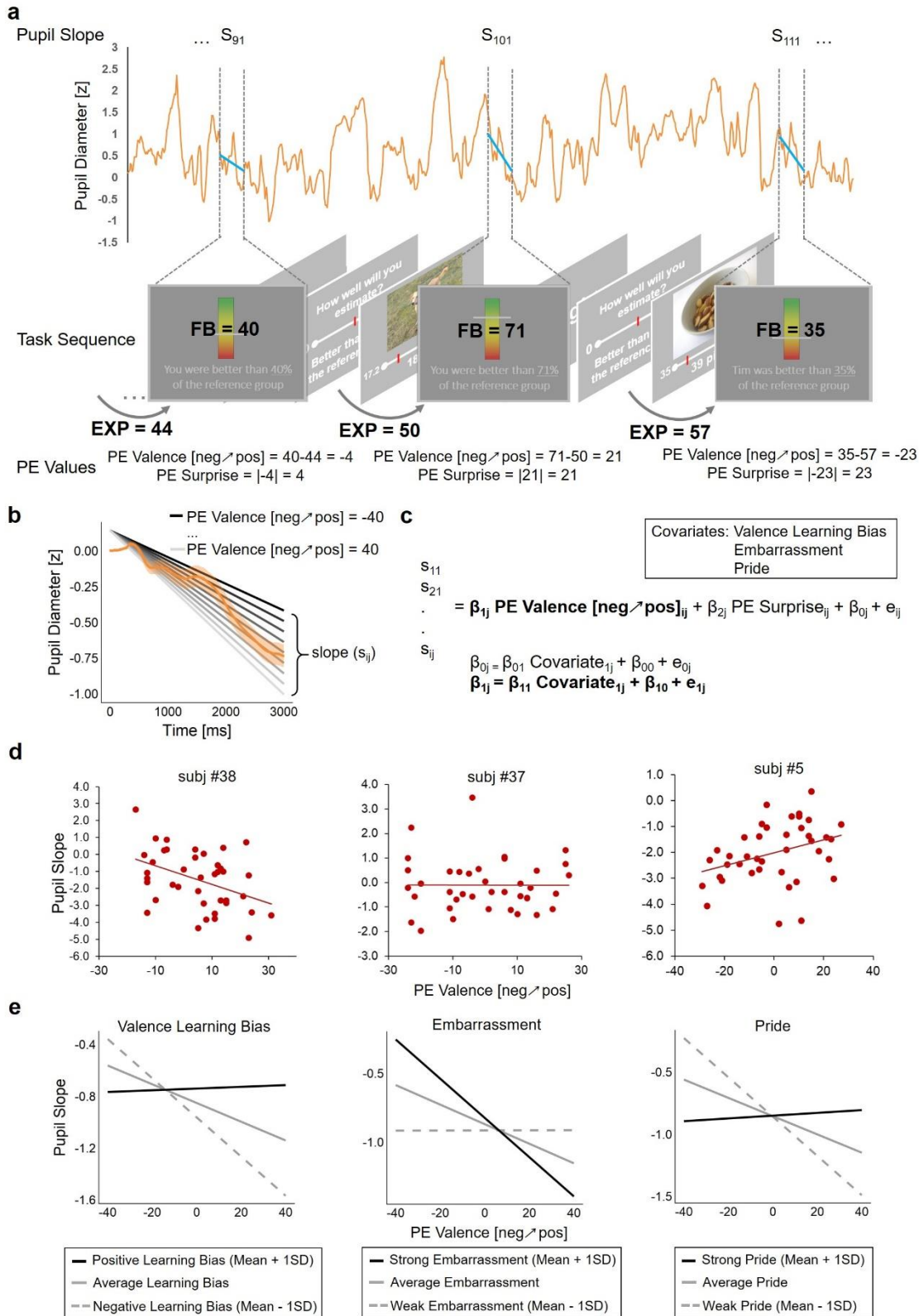
270 Prior research has successfully linked changes in pupil diameter to surprise, PEs and
271 learning (Koenig et al., 2018; Preuschoff, 2011) as well as emotional experiences and arousal
272 (Bradley et al., 2008; Müller-Pinzler et al., 2015). To corroborate our assumption that changes
273 in pupil diameter, as indicated by the slope of the change in pupil size during self-related
274 feedback presentation, reflect increased arousal or attention associated with greater PEs we
275 regressed trial-by-trial variability in the pupil slope on PE surprise (continuous effect of
276 unsigned PEs) and PE valence [neg↗pos] (continuous effect of signed PEs; see **Figure 2a**)
277 (Rouhani and Niv, 2021). The linear mixed model showed a significant positive effect for PE
278 surprise ($t_{(1406)}=2.20$, $p=.028$) and a negative effect for PE valence [neg↗pos] ($t_{(31.2)}=-2.50$,
279 $p=.018$; see **Figure 2b**). First, we observed an effect of PE surprise in the sense that the more
280 surprising the feedback was with respect to trial-by-trial prior expectations the more pupil
281 dilation increased, in line with previous findings on pupil dilation in response to surprising
282 events (Preuschoff, 2011). Second, the results indicate that pupil dilation was greater with
283 decreasing PE values, linking negative PEs rather than positive PEs to greater dilation,
284 potentially indicating increased arousal and attention towards negative PEs in line with the
285 negativity bias we found in learning rates.

286 **Pupil dilation response to prediction error valence is associated with affect and learning**
287 **bias**

288 It has been suggested that pupil dilation not only reflects differences between stimuli
289 but similarly individual biases during decision making (see **Figure 2d** for examples of
290 individual differences; de Gee, Knapen, & Donner, 2014). To corroborate our assumption that
291 pupil slopes should reflect increased arousal or attention associated with PEs that are

292 preferably used for updating by each individual (i.e. individuals preferably learning from
293 negative vs positive PEs, experiencing more embarrassment/ less pride) we thus introduced
294 individual differences in learning and self-conscious emotions as between-subject covariates
295 into the linear mixed models assessing trial-by-trial pupil slopes (see **Figure 2c**). These
296 analyses demonstrated that individuals who experienced more embarrassment showed
297 stronger pupil dilations for negative compared to positive PEs while in individuals with lower
298 embarrassment pupil slopes did not differ between positive and negative PEs (significant
299 interaction of embarrassment and PE valence [neg↗pos]: $t_{(31.8)}=-2.57, p=.015$; no main effect
300 for embarrassment: $t_{(34)}=-0.42, p=.680$; see **Figure 2e**). Effects were reversed when pride
301 ratings were included in the model instead of embarrassment ratings (interaction pride and PE
302 valence [neg↗pos] $t_{(32.8)}=3.14, p=.004$; main effect of pride: $t_{(34.1)}=0.04, p=.971$). The Valence
303 Learning Bias modulated the relationship between PE valence [neg↗pos] and pupil slopes in
304 the same way (interaction Valence Learning Bias and PE valence [neg↗pos] $t_{(31.3)}=2.96,$
305 $p=.006$; main effect of Valence Learning Bias: $t_{(34.3)}=1.05, p=.300$) indicating that people with
306 a more negative Valence Learning Bias had greater pupil dilation for negative PEs whereas
307 people with no or positive bias, showed less differentiation in the pupil dilation in response
308 to the valence of the PE.

309



310

311 **Figure 2.** Association of pupil slopes with prediction error (PE) valence and individual pupil
 312 response differences explained by differences in Valence Learning Bias, embarrassment and
 313 pride experience. **a**) Exemplary pupil diameter trace over three trials for one subject (orange
 314 line) and trial specific fitted linear slopes (blue lines) for the feedback phase of each trial. PE
 315 values are calculated with the participant's current performance expectation (EXP) and the
 316 following feedback value (FB) and PE valence [neg ↗ pos] represents the signed PE while PE

317 surprise represents the unsigned PE. **b)** More negative PEs are associated with greater pupil
318 slopes compared to positive PEs. The average pupil diameter trace during feedback is depicted
319 in orange, shaded area represents +/- one standard error. Pupil slopes for the different levels
320 of PEs (from black = negative to grey = positive) were predicted by the multi-level model
321 containing PE valence [neg↗ pos] and PE surprise as predictors as described in the methods
322 section and depicted in **c**. **c)** Description of the multi-level model assessing the association of
323 PE valence [neg↗ pos] and PE surprise with within-subject trial-by-trial pupil slopes and the
324 impact of Valence Learning Bias, embarrassment and pride experience as between-subject
325 second level covariates explaining differences in the associations on the within-subject level
326 (cross-level-interaction; indicated in bold). **d)** Three exemplary scatter plots for three different
327 subjects show the association of pupil slopes with PE valence [neg↗ pos] and illustrate the
328 variance between subjects. Subj#38 shows increased pupil slopes for negative PEs in line with
329 the group level effect. Subj#37 shows no association and subj#5 shows increased pupil slopes
330 for positive PEs. **e)** Illustration of the impact of the three between-subject covariates, Valence
331 Learning Bias (left), embarrassment (middle) and pride experience (right) explaining
332 differences in the associations of PE valence [neg↗ pos] and pupil slope. The plots show the
333 data as predicted by the multi-level models for the mean covariate as well as the mean
334 covariate +/- 1 standard deviation (SD). The results show that, for example, individuals with
335 a more negative Valence Learning Bias (gray dashed line, left), increased experience of
336 embarrassment (black line, middle) and decreased experience of pride (gray dashed line, right)
337 showed a stronger pupil response to negative vs positive PEs.
338

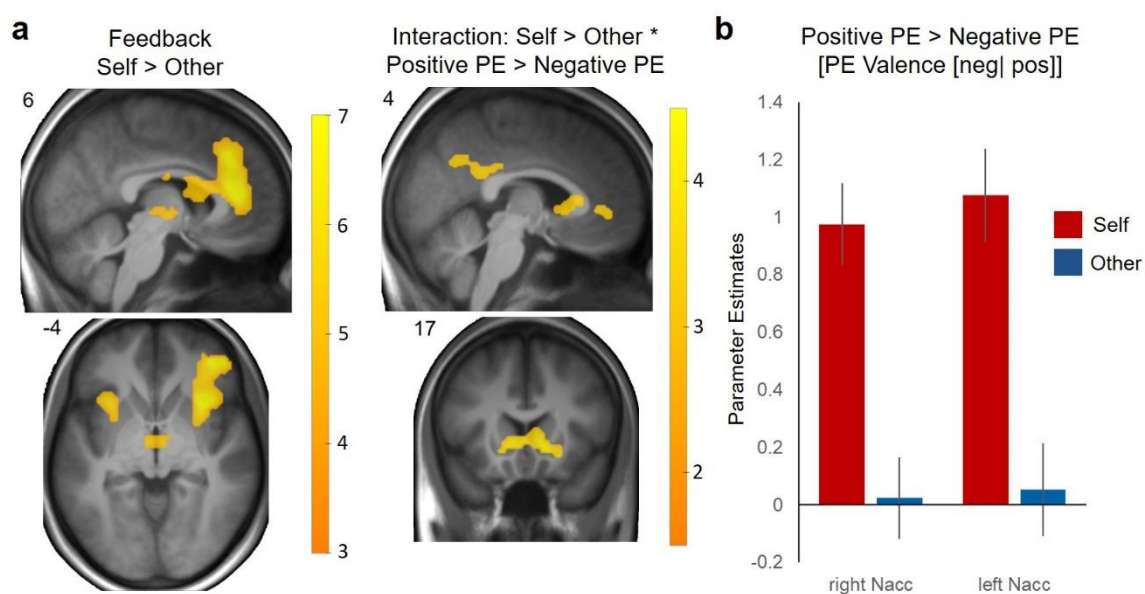
339 **Neural activations associated with feedback processing indicate a specific role of** 340 **feedback valence during self-related learning**

341 In a next step, we examined the brain processes that underlie how people form self-
342 and other-related ability beliefs. Therefore, we first compared neural activation during
343 feedback processing as measured with fMRI. We found that the bilateral insula, anterior
344 cingulate cortex, and thalamus (amongst others, see **Figure 3a** and **Supplementary Table**
345 **S2**) were activated significantly stronger if feedback was related to one's own performance as
346 compared to the other person (i.e. Agent effect). This finding of heightened activity in brain
347 regions that have been linked to arousal but also self-agency potentially reflects a difference
348 in the subjective salience of self- vs other-related information (Craig, 2009; Späti et al., 2014;
349 Sperduti et al., 2011). Feedback for the Other as compared to the Self resulted in stronger
350 activation of the left and right middle temporal gyrus and precuneus/ middle cingulate gyrus
351 (**Supplementary Table S2**).

352 Second, we compared self-related positive vs negative feedback to examine how the

353 valence of information affected neural processing (categorical PE valence [pos| neg] effect).
354 We found significantly stronger activations of the left and right nucleus accumbens/ ventral
355 striatum (NAcc/VS), bilateral angular gyrus, medial prefrontal cortex (mPFC), and precuneus/
356 posterior cingulate cortex (PCC) for positive vs. negative PE valence [pos| neg] (see
357 **Supplementary Table S2**). This valence effect was unique for processing of self-related
358 information and absent when feedback was related to the other person's performance (no
359 significant clusters for the PE valence [pos| neg] effect for Other; $p < .001$). The opposite
360 contrast, negative vs positive PE valence [pos| neg], yielded no significant activations, neither
361 for self-related nor for other-related information. When testing the interaction of Agent x PE
362 valence [pos| neg] we found increased activation for self-related positive vs negative feedback
363 ([Self positive PE > Self negative PE] > [Other positive PE > Other negative PE]) in the
364 angular gyrus (see **Supplementary Table S2**), and on a more lenient threshold also the
365 bilateral NAcc/VS, the precuneus/ PCC, and precentral gyrus (cluster-wise FWE-corrected
366 with $p < .05$ at a cluster forming threshold of $p < .001$; see **Figure 3a/ b** and **Supplementary**
367 **Table S3**).

368



369

370 **Figure 3.** Neural activations associated with feedback processing. **a)** Self-related feedback vs
371 other-related feedback (left) was associated with an increased activation of the mPFC/ ACC,

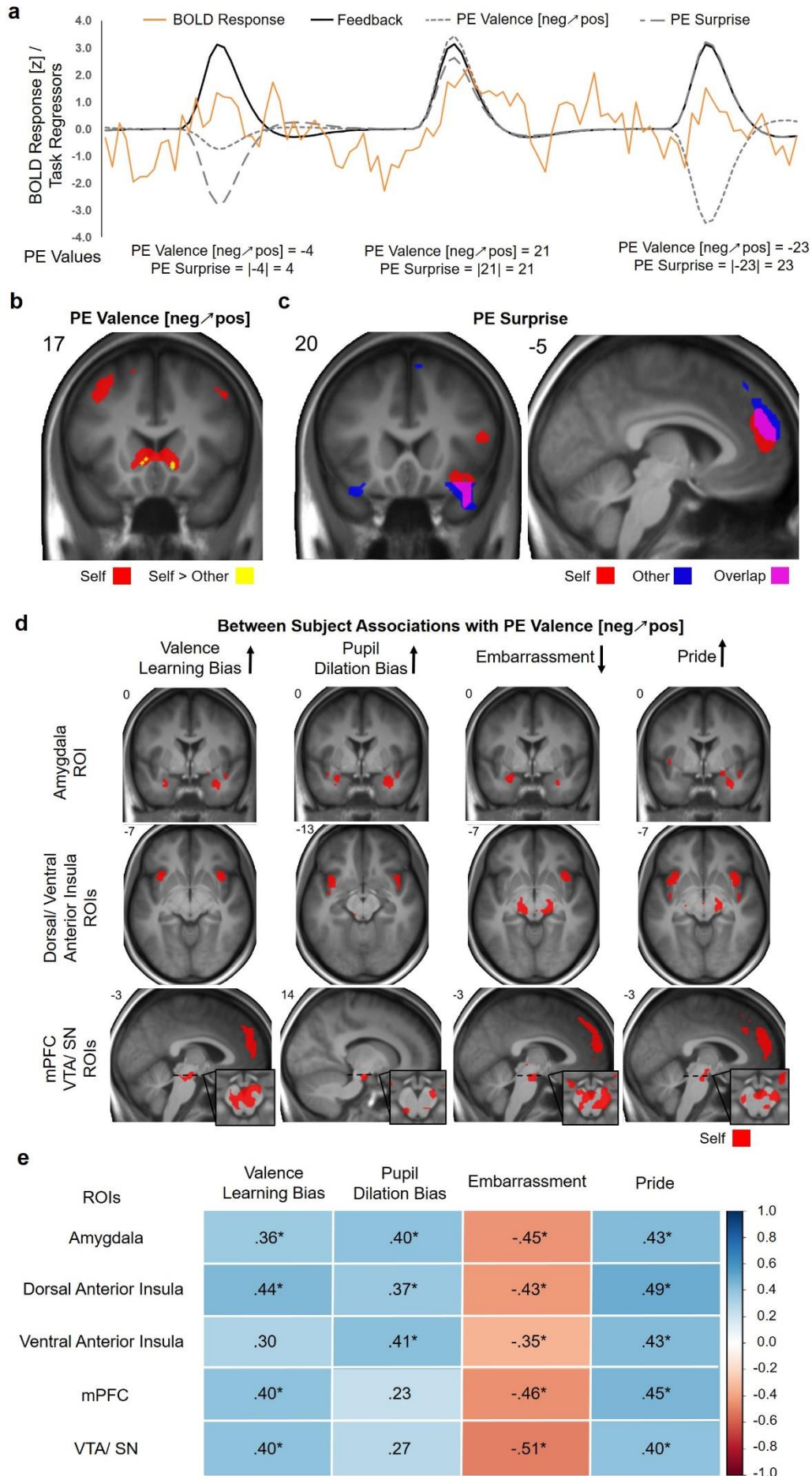
372 bilateral anterior insula and thalamus among other regions ($p < .05$, FWE-corrected). The
373 interaction of Agent and PE valence [pos|neg] (([Self positive PE > Self negative PE] > [Other
374 positive PE > Other negative PE]; right) resulted in activation of the angular gyrus, the
375 bilateral NAcc/Vs, the precuneus/ posterior cingulate cortex, and precentral gyrus (cluster-
376 wise FWE-corrected with $p < .05$ at a cluster forming threshold of $p < .001$). **b)** Parameter
377 estimates for the differences between feedback for positive vs negative PEs derived from the
378 peak voxels of the interaction effect depict the interaction in the left NAcc/ VS [x, y, z: -9 20
379 -1] and right NAcc/ VS [x, y, z: 12 20 -1]. Self-related feedback resulted in a valence specific
380 activation while other-related feedback did not.
381

382 **Common neural activations associated with prediction error surprise and distinct** 383 **activations for self-related prediction error valence**

384 In a second step we assessed the effects of continuous trial-by-trial PE surprise and PE
385 valence [neg↗ pos] as parametric weights to assess neural aspects of learning more
386 specifically (see **Figure 4a**). Increased PE surprise was associated with greater activation of
387 the mPFC for Self and Other as well as clusters in the bilateral insula/ temporal pole/ frontal
388 orbital gyrus on a more lenient threshold (cluster-wise FWE-corrected with $p < .05$ at a cluster
389 forming threshold of $p < .001$; see **Figure 4c** and **Supplementary Table S4**). There was no
390 significant difference between Self and Other ($p < .001$), potentially indicating that there is a
391 common process of error tracking mapped within the same brain regions independent of the
392 agent.

393 Assessing PE valence [neg↗ pos] revealed a distinct pattern for self- and other-related
394 learning. Self-related PE valence [neg↗ pos] was positively associated with increased
395 activation of the NAcc/Vs, mPFC, bilateral angular gyrus/ superior parietal lobule/ lateral
396 occipital gyrus and precentral gyrus, showing stronger activation for positive vs negative PEs
397 (**Figure 4b** and **Supplementary Table S5**). There was no effect for other-related PE valence
398 [neg↗ pos] and a direct comparison of self vs other-related PE valence [neg↗ pos] effects
399 showed stronger associations in the NAcc/Vs for Self (right: x, y, z: 12, 17, -4, $t_{(38)} = 5.23$; k
400 = 2; left: x, y, z: -9, 26, -1, $t_{(38)} = 5.77$, k = 19; all coordinates in MNI space), supporting that
401 the valence of the feedback has a specific value when feedback refers to the self as compared

402 to another person. Although behavioral data and learning rates clearly stress the importance
403 of negative over positive PEs, there were no significant negative associations with PE valence
404 [neg↗pos] in the neural data ($p < .001$). To test if the activations associated with self-related
405 PE valence [neg↗pos] were actually related to PEs and not only to the feedback value alone
406 we ran an additional control analysis including parametric weights for trial-by-trial feedback
407 and performance expectation ratings instead of PE values (Zhang et al., 2020). This model
408 replicated the findings showing positive associations within the reported brain regions,
409 including the NAcc/ VS, with trial-by-trial self-related feedback values and negative
410 associations with prior expectations as it would be expected for PE-related effects
411 **(Supplementary Table S6).**



413 **Figure 4.** Common neural activations associated with prediction error (PE) surprise, distinct
414 activations for self-related PE valence [neg↗ pos] and individual response differences to PE
415 valence [neg↗ pos] explained by differences in Valence Learning Bias, embarrassment and
416 pride experience, and pupil dilation. **a)** Exemplary BOLD response over three trials for one
417 subject (orange line) and regressors for the feedback phase of each trial (back line; originally
418 separate regressors for self- and other-related feedback are here combined for displaying
419 purposes). PE valence [neg↗ pos] (small dashed) and PE surprise (large dashed) are added as
420 parametric modulators in addition to the feedback regressors. PE values are calculated as
421 shown in **Figure 2.** **b)** PE valence [neg↗ pos] was associated with increased activation of the
422 NAcc/VS, mPFC, bilateral angular gyrus/ superior parietal lobule/ lateral occipital gyrus and
423 precentral gyrus for Self. Activation of the NAcc/VS was stronger for Self vs Other. **c)** PE
424 surprise was associated with activation of the mPFC and the bilateral insula/ temporal pole/
425 frontal orbital gyrus for Self and Other. **d)** Neural tracking of self-related PE valence [neg↗
426 pos] in the predefined regions of interest (ROIs: amygdala, ventral and dorsal anterior insula,
427 mPFC, VTA/ SN) was modulated by between-subject variables. Specifically, individuals with
428 a more negative Valence Learning bias showed relatively stronger neural responses to
429 negative as compared to positive PEs, which was similar for individuals with stronger pupil
430 dilation response to negative PEs, or increased ratings of embarrassment. In contrast,
431 individuals with higher pride ratings showed relatively stronger neural responses to positive
432 as compared to negative PEs. Black arrows indicate the direction the covariates are coded in
433 the analyses; embarrassment is coded negatively as high embarrassment was supposed to be
434 associated with increased activity for negative rather than positive PE valence [neg↗ pos].
435 Clusters refers to $p < .005$, uncorrected for displaying purposes; see **Supplementary Table**
436 **S7** for FWE-corrected statistics. **e)** Pearson correlations for parameter estimates derived from
437 the whole areas of our predefined ROIs with the Valence Learning Bias, Pupil Dilation Bias,
438 embarrassment and pride are color coded. * = $p < .05$, FDR corrected.
439

440 **Neural activity in response to self-related prediction error valence is associated with** 441 **affect, learning bias, and pupil dilation**

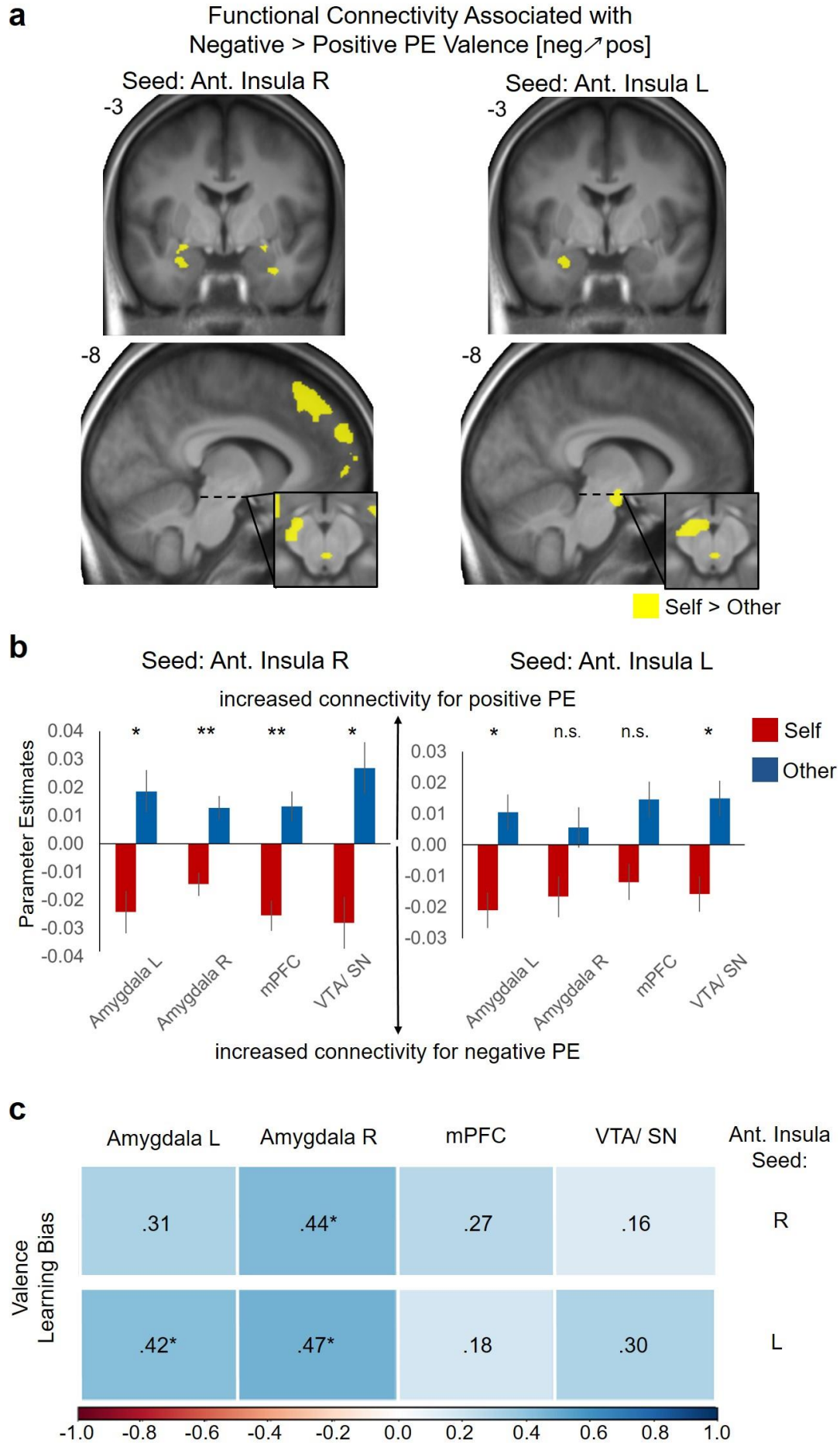
442 To assess how biases in learning as well as affective experience and pupil dilation were
443 associated with valence specific PE processing on the single trial level, multiple general linear
444 models (GLMs) were implemented. These included the Valence Learning Bias, self-conscious
445 emotions and a score representing a valence bias for pupil dilation responses towards positive
446 vs negative PEs (Pupil Dilation Bias = $\text{PupilSlope}_{\text{Self/PE}^+} - \text{PupilSlope}_{\text{Self/PE}^-}$) as between
447 subject covariates for PE valence [neg↗ pos] tracking. Analyses within our predefined regions
448 of interest (ROIs) revealed that the more negative the Valence Learning Bias was, the more
449 neural activity increased in response to negative vs positive PEs in the bilateral dAI, vAI,
450 amygdala, mPFC, and VTA/ SN (all results are $p < .05$ FWE-corrected within ROIs; see **Figure**
451 **4d** and **Supplementary Table S7**). Overall higher experience of embarrassment showed

452 similar associations with increased PE tracking for negative vs positive PEs in the right dAI,
453 bilateral amygdala, and VTA/ SN. Trendwise effects for embarrassment were found in the left
454 dAI, bilateral vAI, and mPFC. In line with this, lower experience of pride showed the same
455 association in the dAI and vAI, amygdala, VTA/ SN and mPFC. Additional analyses revealed,
456 that effects for embarrassment and pride were mainly independent (see **Supplementary**
457 **Results**). Similarly, the more negative the Pupil Dilation Bias was, the stronger the increase
458 in activation of the dAI and vAI, amygdala and VTA/ SN towards negative vs positive PEs.
459 Thus, the greater the response of this neural system to PEs with negative valence the greater
460 was the preference for negative information during learning as well as the negativity of the
461 affective experience. This gained multi-modal support by similar associations of the Valence
462 Learning bias and affect with the pupil dilation response, which reflects the activity of this
463 underlying neural system. In contrast, participants which showed a greater response of this
464 neural system towards positive PEs also had a preference for positive information during
465 learning and reported more positive affect.

466 **Functional connectivity of the dorsal anterior insula depends on prediction error valence**
467 **in line with the negativity bias**

468 Functional connectivity dynamics of the left and right dAI were assessed as these were
469 activated during feedback processing for self- and other-related feedback, independent of
470 Agent and PE valence. Here, we tested if the connectivity of the dAI differed depending on
471 the level of the valence of the PE, by using psychophysiological interaction (PPI) analyses.
472 We calculated the interaction of the continuous PE valence [neg↗pos] and the timeseries
473 extracted from the left and right dAI seed regions separately for Self and Other on the first
474 level and the two agents were contrasted against each other on the second level GLM as we
475 were specifically interested in connectivity dynamics that could reflect the differential
476 learning from negative PEs when processing self-relevant information. Contrasting the PPI
477 effects for PE valence [neg↗pos] between the Self and Other demonstrated that during self-

478 related learning, functional connectivity dynamics of the right dAI with the bilateral
479 amygdala, mPFC and VTA/ SN ($p < .05$, FWE-corrected within ROIs) more strongly aligned
480 with the negativity of the PEs. The left dAI showed a weaker but similar spatial distribution
481 with significant differences between self- and other-related PE valence [neg \nearrow pos] for the left
482 amygdala and VTA/ SN ($p < .05$, FWE corrected, see **Figure 5a/ b** and **Supplementary Table**
483 **S8**). Thus, those brain regions that preferably tracked PEs of negative valence in individuals
484 with increased negative affect and learning biases also showed connectivity dynamics with
485 the dAI in a similar direction during self-related learning. Individuals who showed more
486 pronounced differences in functional connectivity, that is, stronger functional connectivity for
487 negative PEs during Self>Other, also showed a more negative Valence Learning Bias,
488 although this pattern was not fully consistent across all ROIs (see **Figure 5c**).



490 **Figure 4.** Differences in functional connectivity of the dorsal anterior insula during prediction
491 error (PE) valence [neg↗pos] tracking for self- and other-related learning and associations
492 with the Valence Learning Bias. **a)** Increased functional connectivity of the dorsal anterior
493 insula for the negative effect of PE Valence [neg↗pos] for self- vs other-related learning in
494 the predefined ROIs (amygdala, mPFC, VTA/ SN; $p < .005$ uncorrected for displaying
495 purposes). **b)** Functional connectivity dynamics of the dorsal anterior insula plotted separately
496 for self- and other-related learning. For displaying purposes parameter estimates are plotted
497 separately for Self and Other and refer to the peak voxels of the contrast Self vs Other that are
498 reported in **Supplementary Table S8**. * = $p < .05$, ** = $p < .01$. **c)** Spearman correlations of
499 the Valence Learning Bias with the functional connectivity dynamics between the dorsal
500 anterior insula (seed region reported on the right side) and the amygdala, mPFC and VTA/ SN
501 associated with PE valence [neg↗pos] for self- vs other-related learning are color coded. * =
502 $p < .05$, FDR corrected.
503

504 Discussion

505 Belief formation is essentially biased and various studies have shown how motivations
506 shape belief formation (Elder et al., 2021; Sedikides and Gregg, 2008; Sedikides and Hepper,
507 2009; Sharot and Garrett, 2016). Our results demonstrate that the affect people experience
508 during learning is linked to the process of self-related belief formation on the level of neural
509 systems. Our computational modelling results imply that biases in the formation of self-
510 efficacy beliefs in mastering a conceptually novel task are associated with the experience of
511 the self-conscious emotions of embarrassment and pride. Critically, on the level of neural
512 systems, the valence of PEs was associated with biases in self-related learning biases and
513 negativity of the affective experience. Individual differences in the response preference
514 towards negative PEs indicated by the pupil dilation response and activation of the AI,
515 amygdala, mPFC, and VTA/ SN were associated with a more negative learning bias and
516 negative affective experience, hinting towards a neurobiological system that integrates affect
517 during learning.

518 The novel framework on the “value of beliefs” proposed by Bromberg-Martin and
519 Sharot (2020) nicely details how beliefs elicit emotions while at the same time these emotions
520 shape how beliefs are updated in a reciprocal relation. Based on this framework and prior
521 research on self-conscious emotions a negative belief about one’s abilities should elicit

522 stronger embarrassment after failures and reduced pride after successes (Müller-Pinzler et al.,
523 2015; Tangney et al., 2007). In the present data the association of the learning bias with the
524 affective experience during learning supports this notion, with individuals who experience
525 more negative affect (embarrassment) and less positive affect (pride) when receiving self-
526 efficacy feedback being also inclined to update their beliefs in a more negative way. Negative
527 emotions at the same time guide the information processing at various stages, including
528 perception, attention, and decision-making as discussed in the context of “motivated
529 cognition” (Hughes and Zaki, 2015). This reciprocal relation finally results in biased belief
530 formation and beliefs that are both driver of affect as well as affected by emotional responses
531 to incoming information. Embarrassment in particular is one relevant example illustrating the
532 recursive relation between both, as the fear of failure as often discussed in the context of social
533 anxiety (disorder) (Koban et al., 2017; Morrison and Heimberg, 2013; Müller-Pinzler et al.,
534 2015, 2019), leads to shifts of expectations and attention (threat monitoring) towards negative
535 information. At best this results in reparative behaviour and performance improvement (Darby
536 and Harris, 2010; Keltner and Potegal, 1997) and at worst results in a vicious cycle of fear and
537 pathologically increasing negative beliefs about the self (Heimberg et al., 2010) as it is
538 reflected in our results, when individuals who experience more intense embarrassment end up
539 with lower self-efficacy beliefs.

540 There are different ways on how emotions shape learning processes: First, emotions
541 can imbue how information is processed in the brain by adaptively shifting attention towards
542 salient aspects of the situation (Christianson, 2014; Kaspar and König, 2012). Second,
543 emotions entail arousal that intensifies internal rehearsal and evaluations leading to increased
544 learning (Christianson, 2014; Frijda, 1987; Storbeck and Clore, 2008), although these
545 processes often interact and are intricately related (Hughes and Zaki, 2015). The increased
546 pupil dilation in response to negative PEs in our study is in line with both increased salience
547 of and attentional shifts towards negative PEs (Koenig et al., 2018; Preuschoff, 2011) or

548 increased arousal elicited by negative PEs (Bradley et al., 2008; Müller-Pinzler et al., 2015).
549 Here, we think that the stronger impact of positive or negative information on pupil responses
550 and brain reactivity maps arousal and affect according to the valence of individual learning
551 biases and affective experiences.

552 Specifically the AI has been suggested to function as an integrative hub for motivated
553 cognition and emotional behavior (Koban and Pourtois, 2014; Wager and Feldman Barrett,
554 2017). While ventral aspects of the AI are associated with affective processing, emotions, and
555 physiological arousal (Craig, 2003; Lindquist et al., 2012; Phan et al., 2002; Wager and
556 Feldman Barrett, 2017) dorsal aspects of the AI are strongly associated with the detection of
557 salient events, allocation of attention resources, executive working memory (Menon and
558 Uddin, 2010; Touroutoglou et al., 2012) and also (absolute) PEs and uncertainty during
559 learning (Loued-Khenissi et al., 2020; Rutledge et al., 2010; Ullsperger et al., 2010). These
560 findings suggest that the functions of the AI provide a physiological basis for how emotions
561 are translated into biased, motivated, or affected beliefs (Koban and Pourtois, 2014; Wager
562 and Feldman Barrett, 2017). A similar role as a link for the attention-emotion interaction has
563 also been suggested for the amygdala (Kaspar and König, 2012; Koban and Pourtois, 2014),
564 that shows similar responses in our task. The functional connectivity dynamics of the dAI,
565 matching the modelled learning rates with a stronger impact of self-related negative PEs,
566 underline the insula's role as an integrative hub receiving and forwarding information that
567 affects information processing in other brain regions.

568 Tracking of PEs in the dopaminergically innervated VTA/ SN is influenced by
569 motivational factors during learning (Adcock et al., 2006). The subjective value of self-related
570 information significantly varies between subjects which is indicated by idiosyncratic response
571 patterns of the VTA/ SN to gains or losses (Charpentier et al., 2018). In this line, we think that
572 the present results reflect individual response tendencies at a very basic level of PE tracking.
573 On higher layers of the computational hierarchy regions in the ACC and mPFC are also

574 associated with PE tracking and value representation (Hare et al., 2008; Lockwood and
575 Wittmann, 2018; Wallis and Kennerley, 2010) and have been previously associated with
576 biases in belief updating (Korn et al., 2012; Kuzmanovic et al., 2016, 2018). Affect and arousal
577 could therefore bias learning on various stages of the computational hierarchy of PE
578 processing from more basic dopaminergic midbrain responses to more abstract value
579 representations in the neocortex (Diaconescu et al., 2017). While the directionality of the
580 effects remains to be determined the dynamics in the functional connectivity of the dAI
581 suggest a modulatory role in this process. Here, information is forwarded to and/ or integrated
582 from VTA/ SN and mPFC, the same regions, whose response to the valence of PEs was also
583 modulated by differences in learning bias and affective experience. This strengthens the idea
584 that the AI plays a role in shifting responses to negative or positive information in other brain
585 regions (e.g. by shifting attention and by affective tagging) or already receives stronger signals
586 in response to PEs of negative or positive valence from midbrain regions and mPFC.

587 The tracking of the absolute error, PE surprise (Rouhani and Niv, 2021), independent
588 of the agent, is in line with the “common currency” assumption (Izuma et al., 2008; Ruff and
589 Fehr, 2014) for the positive and negative value of one’s own and others’ performance
590 feedback. The common and valence independent coding of surprise in the insula and the
591 mPFC might therefore be sufficient to complete the learning task per se. Valence, however,
592 matters when individuals learn about themselves as indicated by an additional shift in error
593 tracking in the same regions, AI and mPFC, which also track surprise in a valence independent
594 manner. As a consequence, across individuals we observe a robust effect for surprise,
595 however, when people learn more negatively biased and experience more negative affect, this
596 signal is unbalanced and increases with more negative PEs. This pattern hints towards a
597 neurocomputational mechanism of how affect shapes the formation of beliefs as proposed
598 earlier (Bromberg-Martin and Sharot, 2020).

599 In the current study, some of the key findings emerge at the level of individual

600 differences. We observed a wide inter-individual variance in the affective experience during
601 the task and in the Valence Learning Bias, that is, the kind of information participants
602 preferably used to update self-related beliefs. While on average we find a negativity bias
603 during self-related belief formation, a little less than one third of the participants still shows a
604 positive learning bias, pointing out the importance of individual factors and meaningfulness
605 of variability. Studies do not only suggest that biases in belief formation differ between tasks
606 (Ertac, 2011; Müller-Pinzler et al., 2019; Sharot and Garrett, 2016) but also depend on
607 situational factors like stress (Czekalla et al., 2021; Garrett et al., 2018). The individual's
608 ability to adjust the current information processing strategy to the context might be adaptive
609 (Bromberg-Martin and Sharot, 2020): for example, adaptation to an increased relevance of
610 negative or threat-related information during stress (Garrett et al., 2018) or coping with a
611 negative self-concept following social stress by means of more self-beneficial belief updating
612 (Czekalla et al., 2021). It might also be adaptive for people who fear negative feedback to pay
613 more attention to failure-related information in order to learn and circumvent potential future
614 failures (Sedikides and Hepper, 2009). However, it is not always straightforward to determine
615 under which conditions a strategy is adaptive or whether the affective experience can
616 ameliorate the individual's well-being. A maladaptive consequence of biased self-efficacy
617 beliefs becomes apparent in psychiatric disorders such as depression and social anxiety, in
618 which amplified negative updating can lead to persistently distorted self-views and overly
619 negative beliefs about one's own capabilities in everyday life (Alden et al., 2008; Amir et al.,
620 2012; Koban et al., 2017; Korn et al., 2014; Taylor and Brown, 1988).

621 **Conclusions**

622 Emotions experienced during learning affect computational mechanisms and manifest
623 in distributed neural activity during belief formation. In particular, neural activity of the AI,
624 amygdala, VTA/SN, and mPFC and pupil responses map the valence of PEs in
625 correspondence to the experienced affect and learning bias people show during belief

626 formation. The more negative balancing in the functional connectivity dynamics of the dAI
627 during processing of self-related PEs within this network outline a scaffold for neural and
628 computational mechanism integrating affect during belief formation. The results of the first
629 empirical spell-out of the “value of beliefs model” (Bromberg-Martin and Sharot, 2020) have
630 broader implications concerning any context which provides personal evaluations based on
631 behavioral performance. Here, the focus on the affective experience during learning provides
632 a deeper understanding on how feedback manifests in self-related beliefs which may then
633 significantly impact developmental processes and future behavior.

634

635

Materials and Methods

636 Participants

637 The study was approved by the ethics committee of the University of Lübeck (AZ 18-
638 066), has been conducted in compliance with the ethical guidelines of the American
639 Psychological Association (APA), and all subjects gave written informed consent. Participants
640 were recruited at the University Campus of Lübeck, were fluent in German, and had normal
641 or corrected-to-normal vision. Two independent samples were recruited, one for the fMRI
642 study and the other for a behavioral study that was added to increase the sample size for the
643 behavioral data. All participants received monetary compensation for their participation in the
644 study. The final sample size for the fMRI sample was 39 participants (26 females, aged 18-28
645 years; $M=22.3$; $SD=2.65$). We initially recruited 48 participants and had to exclude six
646 participants because they did not believe the cover-story of the task and three participants
647 because they did not attentively complete the task until the end (e.g. participants reported
648 having been too tired or the ratings indicated that they stopped responding to the estimation
649 task). The additional behavioral sample consisted of 30 participants (24 females, aged 18-32
650 years; $M=23.3$; $SD=3.97$). For more details on the sample characteristics see **Supplementary**
651 **Table S9**.

652 **Learning of own performance task.**

653 The Learning of own performance (LOOP) task enables participants to incrementally
654 learn about their or another person’s alleged ability in estimating properties. The LOOP task
655 has been previously introduced and validated in a set of behavioral studies (Müller-Pinzler et
656 al., 2019). For the LOOP task all participants were invited to take part in an experiment on
657 "cognitive estimation" together with a confederate who allegedly was another participant. In
658 contrast to the fMRI study, for the behavioral study two participants were invited and tested
659 together instead of introducing a confederate. Participants were informed they would take
660 turns with the other participant/ confederate, either performing the task themselves (Self) or
661 observing the other person performing (Other). For the task participants were instructed to
662 estimate different properties (e.g. the height of houses or the weight of animals). On a trial-
663 by-trial basis participants received manipulated performance feedback in two distinct
664 estimation categories for their own estimation performance as well as for the other person’s
665 estimation performance. Unbeknownst to the participant, one of the two categories was
666 arbitrarily paired with rather positive feedback while the other was paired with rather negative
667 feedback (e.g. “height” of houses = High ability category and “weight” of animals = Low
668 ability category or vice versa; estimation categories were counterbalanced between Ability
669 conditions and Agent [Self vs Other] conditions). This resulted in four feedback conditions
670 with 20 trials each (Agent condition [Self vs Other] x Ability condition [High Ability vs Low
671 Ability]). Trials of all conditions were intermixed in a fixed order with a maximum of two
672 consecutive trials of the same condition. Performance feedback was provided after every
673 estimation trial, indicating the participant’s own or the other person’s current estimation
674 accuracy as percentiles compared to an alleged reference group of 350 university students
675 who, according to the cover-story, had been tested beforehand (e.g. "You are better than 94%
676 of the reference participants."; see **Figure 1a**). The feedback was defined by a sequence of
677 fixed PEs with respect to the participants’ “current belief” about their abilities. The “current

678 belief” was calculated as the average of the last five performance expectation ratings per
679 category, which started at 50% before participants actually rated their performance
680 expectation. This procedure led to varying feedback sequences between participants but kept
681 PEs mostly independent of the participants’ performance expectations and insured a relatively
682 equal distribution of negative and positive PEs across conditions (Self: mean positive PE =
683 13.6, SD = 1.8 (average number = 20.3); mean negative PE = -12.6, SD = 1.4 (average number
684 = 19.7); Other: mean positive PE = 13.0, SD = 1.3 (average number = 19); mean negative PE
685 = -13.1, SD = 1.1 (average number = 21)). At the beginning of each trial a cue was presented
686 indicating the estimation category (e.g. “height”) and participants were asked to state their
687 expected performance for this trial on the same percentile scale used for feedback. As part of
688 the cover story, participants were informed that accurate expected performance ratings would
689 be rewarded with up to 6 cents per trial, that is, the better their expected performance rating
690 matched their actual feedback percentile the more money they would receive, to increase
691 motivation and encourage honest response behavior. Following each performance expectation
692 rating, the estimation question was presented for 10 seconds. During the estimation period,
693 continuous response scales below the pictures determined a range of plausible answers for
694 each question. Participants indicated their responses by navigating a pointer on the response
695 scale with an MRI compatible computer mouse. Subsequently, feedback was presented for 3
696 seconds (see **Figure 1a**). Jittered inter-stimulus intervals were presented following the cue (2-
697 $6 * TR$ (0.992 secs)), estimation ($2.5 - 6.5 * TR$) and feedback phase ($4-8 * TR$) for the fMRI
698 task. All stimuli were presented using MATLAB Release 2015b (The MathWorks, Inc.) and
699 the Psychophysics Toolbox (Brainard, 1997). The fMRI task was completed in two separate
700 sessions of each 20 min with a short break in between.

701 Before starting the experiment all participants answered several questions about their
702 self-related beliefs and filled in a self-esteem personality questionnaire (SDQ-III; Marsh &
703 O’Neill, 1984). During the LOOP task participants were also asked to rate their current levels

704 of embarrassment, pride, happiness and stress/ arousal on a continuous scale ranging from not
705 at all (coded as 0) to very strong (coded as 100). Two emotion rating phases followed self-
706 related feedback and two rating phases followed other-related feedback. The two emotion
707 rating phases following self-related feedback were averaged to receive a rating for the
708 experience of self-conscious affect (embarrassment and pride) during self-related learning.
709 After the task participants completed an interview including ratings about self-related beliefs,
710 were debriefed about the cover-story, and reimbursed for their time before leaving. The whole
711 procedure took approximately 2 h.

712 **Behavioral Data Analysis and Modeling**

713 A model free analysis was performed on the participants' expected performance
714 ratings for each trial to illustrate the basic effects we see in our behavioral data. A repeated-
715 measures ANOVA was calculated with the factors Trial (20 Trials) x Ability condition (High
716 ability vs Low ability) x Agent condition (Self vs Other) as well as Group as a between-subject
717 factor to control for potential differences between the two samples. All statistical analyses on
718 the behavioral data apart from the modeling procedure were performed using *jamovi* (Version
719 1.2.27, The jamovi project (2020). Retrieved from <https://www.jamovi.org>).

720 Dynamic changes in self-related efficacy beliefs, that is, performance expectation
721 ratings, were then modeled using PE delta-rule update equations (adapted Rescorla-Wagner
722 model; Rescorla & Wagner, 1972). The model space contained three main models varying
723 with regards to their assumptions about biased updating behavior when learning about the self
724 (see **Supplementary Figure S1**). The simplest learning model used one single learning rate
725 for all conditions for each participant, thus not assuming any learning biases (Unity Model).
726 The second model, the Valence Model, included separate learning rates for positive PEs (α_{PE+})
727 vs negative PEs (α_{PE-}) across both ability conditions, thus suggesting that the valence (positive
728 vs negative) of the PE biases self-related learning. The third model, the Ability Model,
729 contained a separate learning rate for each of the ability conditions indicating context specific

730 learning. In addition, learning rates were either estimated separately for Self vs Other or across
731 Agent conditions. The Valence Model with separate learning rates for Self vs Other (Model
732 5), winning model in our previous studies (Czekalla et al., 2021; Müller-Pinzler et al., 2019),
733 was further extended by adding a weighting factor reducing learning rates towards the ends
734 of the feedback scale (percentiles close to 0 % or 100 %), assuming that participants
735 experienced extreme feedback values as less likely than more average feedback (Kube et al.,
736 2021). In the first model of these (Model 7) a linear decrease of the learning rates was assumed
737 beginning at 50 % and ending at 0 % and 100 %. A weighting factor w was fitted for each
738 participant defining how strong the linear decrease was present for each individual. Since
739 many variables people encounter in every-day life (e.g. many test results) approximately
740 follow a normal distribution with extreme values being less likely, for the second model of
741 this kind (Model 8) we assigned the relative probability density of the normal distribution to
742 each feedback percentile value. Again, a weighting factor w was fitted for each individual
743 indicating how strongly the relative probability density reduced the learning rates for feedback
744 further away from the mean. In contrast to our previous studies implementing the LOOP task
745 with fixed feedback sequences, here, feedback depended on the participants' current
746 expectations and thus differed between participants and conditions. Reduced learning rates
747 towards the ends of the feedback scale which could systematically confound learning rates
748 between participants and conditions could thus be accounted for in Models 7 and 8. To test if
749 the participants' performance expectation ratings can be better explained in terms of PE
750 learning as compared to stable assumptions in each Ability condition, we included a simple
751 Mean Model with a mean value for each task condition (Model 9; for more details see
752 **Supplementary Methods**).

753 **Model Fitting.**

754 For model fitting we used the RStan package (Stan Development Team, 2016. RStan:
755 the R interface to Stan. R package version 2.14.1.), which uses Markov chain Monte Carlo

756 (MCMC) sampling algorithms. All of the learning models in the model space were fitted for
757 each participant individually and posterior parameter distributions were sampled for each
758 participant. A total of 2400 samples were drawn after 1000 burn-in samples (overall 3400
759 samples; thinned with a factor of 3) in three MCMC chains. We assessed if MCMC chains
760 converged to the target distributions by inspecting \hat{R} values for all model parameters (Gelman
761 and Rubin, 1992). Effective sample sizes (n_{eff}) of model parameters, which are estimates of
762 the effective number of independent draws from the posterior distribution, were typically
763 greater than 1500 (for most parameters and subjects). Posterior distributions for all parameters
764 for each of the participants were summarized by their mean as the central tendency resulting
765 in a single parameter value per participant that we used in order to calculate group statistics.

766 **Bayesian Model Selection and Family Inference.**

767 For model selection we estimated pointwise out-of-sample prediction accuracy for all
768 fitted models separately for each participant by approximating leave-one-out cross-validation
769 (LOO; corresponding to leave-one-trial-out per subject; Acerbi et al., 2018; Vehtari et al.,
770 2016). To do so, we applied Pareto-smoothed importance sampling (PSIS) using the log-
771 likelihood calculated from the posterior simulations of the parameter values as implemented
772 by Vehtari et al. (2016). Sum PSIS-LOO scores for each model as well as information about
773 \hat{k} values – the estimated shape parameters of the generalized Pareto distribution – indicating
774 the reliability of the PSIS-LOO estimate are depicted in **Supplementary Table S1**. As
775 summarized in **Supplementary Table S1** very few trials resulted in insufficient parameter
776 values for \hat{k} and thus potentially unreliable PSIS-LOO scores (on average 1.1 trials per subject
777 with $\hat{k} > 0.7$ for the winning model; Vehtari et al., 2016). BMS on PSIS-LOO scores was
778 performed on the group level accounting for group heterogeneity in the model that best
779 describes learning behavior (Rigoux et al., 2014). This procedure provides the protected
780 exceedance probability for each model (p_{xp}), indicating how likely a given model has a higher
781 probability explaining the data than all other models in the comparison set. The Bayesian

782 omnibus risk (*BOR*) indicates the posterior probability that model frequencies for all models
783 are all equal to each other (Rigoux et al., 2014). We also provide difference scores of PSIS-
784 LOO in contrast to the model that won the BMS that can be interpreted as a simple ‘fixed-
785 effect’ model comparison (see **Supplementary Table S1**; Acerbi et al., 2018; Vehtari et al.,
786 2016). Model comparisons according to PSIS-LOO difference scores were qualitatively
787 comparable to the BMS analyses for our data.

788 **Posterior Predictive Checks and Statistical Analyses of Learning Parameters**

789 First, posterior predictive checks were conducted by quantifying if the predicted data
790 could capture the variance in performance expectation ratings for each subject within each of
791 the experimental conditions using regression analyses. Additionally, we repeated the model
792 free analysis we had conducted on the behavioral data with the data predicted by the winning
793 model to assess if the winning model captured the core effects in the behavioral data (see
794 **Supplementary Results**).

795 Model parameters, i.e. learning rates, of the winning models for all experiments were
796 analyzed on the group level. A repeated-measures ANOVA was calculated on the learning
797 rates with the factor Agent (Self [$\alpha_{\text{Self/PE+}}$, $\alpha_{\text{Self/PE-}}$] vs Other [$\alpha_{\text{Other/PE+}}$, $\alpha_{\text{Other/PE-}}$]) and factor PE
798 valence [pos| neg] (PE+ [$\alpha_{\text{Self/PE+}}$, $\alpha_{\text{Other/PE+}}$] vs PE- [$\alpha_{\text{Self/PE-}}$, $\alpha_{\text{Other/PE-}}$]) as well as Group as a
799 between-subject factor testing if learning about one’s own performance was more valence
800 specific as compared to learning about the other person’s performance.

801 To associate learning biases with self-conscious affect, that is, embarrassment and
802 pride, and self-esteem (SDQ-III subscale scores) we calculated a normalized learning rate
803 valence bias score for self-related learning (Valence Learning Bias= $(\alpha_{\text{PE+(S)}} - \alpha_{\text{PE-(S)}})/(\alpha_{\text{PE+(S)}} + \alpha_{\text{PE-(S)}})$) (Müller-Pinzler et al., 2019; Niv et al., 2012; Palminteri et al., 2017). Spearman
805 correlations were calculated between Valence Learning Bias, affect ratings, and self-esteem
806 scores.

807 **Pupil Data Analysis**

808 For the fMRI sample, eye-tracking data were assessed during scanning. Pupil diameter
809 and gaze behavior were recorded non-invasively in one eye at 500 Hz using an MRI-
810 compatible Eyelink-1000 plus device (SR Research, Kanata, ON, Canada) with manufacturer-
811 recommended settings for calibration and blink detection. Due to insufficient pupillometry
812 data quality three participants had to be excluded from analyses (final sample n=36). Pupil
813 data were preprocessed by cutting out periods of blinks and values in this gap were
814 interpolated by piecewise cubic interpolation. The pupil trace was subsequently z-normalized
815 over the whole session. To characterize the pupil dilation for each trial by a single value, we
816 calculated a linear slope for each feedback phase of three seconds. Pupil traces were only
817 analyzed for the Self condition as onsets during feedback strongly differed between Agent
818 conditions, which made a meaningful comparison between pupil slopes impossible. Pupil
819 slopes during self-related feedback phases for each trial were then entered in linear mixed
820 models fitted by restricted maximum likelihood including PE valence [neg↗pos] (continuous
821 signed PE values) and PE surprise (continuous unsigned/ absolute PE values) as fixed effects
822 and participant and PE valence [neg↗pos] as random effects. Additionally, separate linear
823 mixed models including embarrassment ratings, pride ratings or the Valence Learning Bias as
824 well as their interaction with PE valence [neg↗pos] were implemented to assess if variance
825 in individual pupil responses to positive and negative PEs (random PE valence [neg↗pos]
826 slopes) was explained by with different emotional reactions and learning behavior (see **Figure**
827 **2c**).

828 **fMRI Data**

829 **fMRI Image Acquisition**

830 Participants were scanned using a 3T Siemens MAGENTOM Skyra scanner (Siemens,
831 München, Germany) at the Center of Brain, Behavior, and Metabolism (CBBM) at the
832 University of Lübeck, Germany with 60 near-axial slices. An echo planar imaging (EPI)
833 sequence was used for acquisition of on average 1520 functional volumes (min=1395, max=

834 1672) during each of the two sessions of the experiment, resulting in a total of on average
835 3040 functional volumes (TR=0.992s, TE=28ms, flip angle=60°, voxel size=3 × 3 × 3mm³,
836 simultaneous multi-slice factor 4). In addition, a high-resolution anatomical T1 image was
837 acquired that was used for normalization (voxel size=1 × 1 × 1mm³, 192×320×320mm³ field
838 of view, TR= 2.300s, TE = 2.94ms, TI = 900ms; flip angle=9°; GRAPPA factor 2; acquisition
839 time 6.55 min).

840 **FMRI data analyses**

841 FMRI data were analyzed using SPM12 (www.fil.ion.ucl.ac.uk/spm). Field maps were
842 reconstructed to obtain voxel displacement maps (VDMs). EPIs were corrected for timing
843 differences of the slice acquisition, motion-corrected and unwrapped using the corresponding
844 VDM to correct for geometric distortions and normalized using the forward deformation fields
845 as obtained from the unified segmentation of the anatomical T1 image. The normalized
846 volumes were resliced with a voxel size of 2 × 2 × 2 mm and smoothed with an 8 mm full-
847 width-at-half-maximum isotropic Gaussian kernel. To remove low-frequency drifts,
848 functional images were high-pass filtered at 1/384.

849 Statistical analyses were performed in a two-level, mixed-effects procedure. Three
850 main GLMs were implemented on the first level. The first fixed-effects GLM included four
851 epoch regressors modeling the hemodynamic responses to the different cue conditions
852 (Ability: High vs Low × Agent: Self vs Other), weighted with the performance expectation
853 ratings per trial as parametric modulator for each condition. Four regressors modeled the four
854 feedback conditions (PE valence [pos| neg]: Positive vs Negative × Agent: Self vs Other),
855 each weighted with the PE value for each trial. The estimation periods for Self and Other as
856 two regressors, and emotion ratings phase and the instruction phase as separate regressors. To
857 account for noise due to head movement, six additional regressors modeling head movement
858 parameters were introduced and a constant term was included for each of the two sessions.
859 The second first-level GLM differed only with respect to the feedback regressors. Here, only

860 two regressors modeled feedback separately for Self and Other and two parametric modulators
861 were included per condition weighting feedback trials with PE valence [neg↗pos] (continuous
862 effect of the signed PE values) and PE surprise (continuous effect of the unsigned PE values).
863 The third first-level GLM was set-up to show that activation found in response to PEs was
864 actually related to PEs and not only to feedback level alone. Therefore, the parametric weights
865 of the two feedback conditions in the second GLM were replaced by feedback level and
866 performance expectation ratings, allowing us to assess if neural activity goes up with feedback
867 level and down with performance expectation ratings confirming a potential interpretation in
868 terms of PE tracking (Zhang et al., 2020).

869 On the second level for the first GLM model beta images for the four feedback
870 conditions were included in a flexible factorial design with two repeated-measurement factors
871 (PE valence [pos| neg] and Agent). Beta images for the parametric weights of feedback were
872 extracted from the second and third first-level model for Self and Other. Separate repeated-
873 measures ANOVAs and one sample t-tests (for baseline contrasts) were implemented for PE
874 valence [neg↗pos] and PE surprise as well as feedback level and performance expectation
875 level. Additional second level models for the PE valence [neg↗pos] contrast included the
876 Valence Learning Bias, embarrassment and pride ratings as between subject covariates,
877 assessing differential tracking of PEs depending on biased learning and self-conscious affect.
878 A self-related Pupil Dilation Bias (average slope for positive PEs - average slope for negative
879 PEs; higher scores indicate stronger pupil dilation for positive PEs) was also included as
880 covariate in another second level model to assess if the neural response towards negative vs
881 positive PEs was associated with the pupil dilation response.

882 We additionally performed psychophysiological interaction (PPI) analyses on the first
883 level, investigating whether functional connectivity of the dAI, that is commonly activated
884 during feedback processing independent of agent and feedback valence (conjunction of
885 baseline contrasts: feedback Self ^ feedback Other) would differ depending on the PE valence

886 [neg↗ pos]. PPI analyses were computed separately for Self and Other and the resulting
887 contrast images for the PPI effects were aggregated on the second level using two-sample t-
888 tests contrasting PPI effects for Self vs Other. For each participant, we defined 6-mm radius
889 spherical ROIs, centered at the nearest local maximum for the conjunction contrast feedback
890 Self ^ feedback Other and located within 10 mm from the group maximum within the dAI,
891 separately for the left dAI (x, y, z: -33 20 -4) and right dAI (x, y, z: 36 20 -7). By computing
892 the first eigenvariate for all voxels within these ROIs that showed a positive effect for the
893 conjunction ($p < .500$), we extracted the time course of activations and constructed PPI terms
894 using the contrast for the parametric weights of PE valence [neg↗ pos] for Self or Other,
895 respectively, resulting in four distinct PPI first level GLMs. One participant was excluded
896 from the PPI analyses for the right dAI, because no voxels survived the predefined threshold
897 for eigenvariate extraction. The PPI term, along with the activation time course from the (left
898 or right) dAI was included in a new GLM for each participant that also included all the
899 regressors in the initial first level GLM (four regressors for the different cue conditions, each
900 weighted with the expected performance ratings; two feedback regressors for Self and Other
901 with each two parametric modulators for PE valence [neg↗ pos] and PE surprise; two
902 regressors for the estimation periods for self and other; one regressor for the emotion ratings
903 phase; one regressor for the instruction phase; six regressors modeling head movement
904 parameters; a constant term for each session). On the second level we assessed if there was a
905 stronger functional coupling of the dAI with our predefined ROIs (Amygdala, mPFC, VTA/
906 SN) for the Self in contrast to the Other when PE valence [neg↗ pos] was more negative.
907 Functional connectivity dynamics were also associated with learning behavior by calculating
908 Spearman correlations for the Valence Learning Bias and the parameter estimates for the PPI
909 effect of Self > Other derived from a sphere of 6mm around the peak voxels within our
910 predefined ROIs.

911 **Thresholding procedure and regions of interest**

912 According to its suggested role as an integrative hub for motivated cognition and
913 emotional behavior the AI was defined as one of the regions of interest (ROIs) (Koban and
914 Pourtois, 2014; Wager and Feldman Barrett, 2017). Due to their specific functional
915 associations, a bilateral ventral and a bilateral dorsal AI ROI was defined according the three
916 cluster solution by Kelly and colleagues (2012). The bilateral amygdala was defined as another
917 ROI and derived from the AAL atlas definition in the WFU PickAtlas (Tzourio-Mazoyer et
918 al., 2002) due to its similar role for the attention-emotion interaction (Kaspar and König, 2012;
919 Koban and Pourtois, 2014). The mPFC ROI was also derived from the AAL atlas in the WFU
920 PickAtlas (label: bilateral frontal superior medial) due to its specific role during social learning
921 and for biases in self-related belief updating in previous studies (Kuzmanovic et al., 2018;
922 Sharot, 2011). Additionally, a VTA/ SN ROI, dopaminergic nuclei in the midbrain, was
923 included (Ballard et al., 2011; Murty et al., 2014) as dopamine signals motivationally important
924 events, e.g. during reward learning (Schultz, 1998), and has been associated with biases in
925 memory towards events that are of motivational significance (Adcock et al., 2006).

926 FMRI results were family-wise-error (FWE) corrected on the whole brain level if not
927 mentioned otherwise and all coordinates are reported in MNI space. As our predefined ROIs
928 were chosen with respect to their involvement with the emotion-cognition link, we tested the
929 effects of our covariates on PE valence [neg/pos] tracking and PPI effects within the ROIs.
930 Anatomical labels of all resulting clusters were derived from the Automated Labeling Atlas
931 Version 3.0 (Eickhoff et al., 2005).

932 **Acknowledgments**

933 We would like to thank Prof. Christoph W Korn for his very helpful comments and
934 discussions on the manuscript. We are also grateful to Clara Gunzelmann and Rebecca
935 Rocksien for their help with data collection. The research was funded by the German Research
936 Foundation (Temporary Positions for Principal Investigators: MU 4373/1-1; Sachbeihilfe KR
937 3803/11-1) and the Medical Department of the University of Lübeck (J21-2018).

938

References

- 939 Acerbi, L., Dokka, K., Angelaki, D.E., and Ma, W.J. (2018). Bayesian comparison of
940 explicit and implicit causal inference strategies in multisensory heading perception. *PLOS*
941 *Comput. Biol.* *14*, e1006110.
- 942 Adcock, R.A., Thangavel, A., Whitfield-Gabrieli, S., Knutson, B., and Gabrieli, J.D.E.
943 (2006). Reward-Motivated Learning: Mesolimbic Activation Precedes Memory Formation.
944 *Neuron* *50*, 507–517.
- 945 Alden, L.E., Taylor, C.T., Mellings, T.M.J.B., and Laposa, J.M. (2008). Social anxiety and
946 the interpretation of positive social events. *J. Anxiety Disord.* *22*, 577–590.
- 947 Amir, N., Prouvost, C., and Kuckertz, J.M. (2012). Lack of a Benign Interpretation Bias in
948 Social Anxiety Disorder. *Cogn. Behav. Ther.* *41*, 119–129.
- 949 Apsler, R. (1975). Effects of embarrassment on behavior toward others. *J. Pers. Soc.*
950 *Psychol.* *32*, 145–153.
- 951 Ballard, I.C., Murty, V.P., Carter, R.M., Macinnes, J.J., Huettel, S.A., and Adcock, R.A.
952 (2011). Dorsolateral Prefrontal Cortex Drives Mesolimbic Dopaminergic Regions to Initiate
953 Motivated Behavior. *J. Neurosci.* *31*, 10340–10346.
- 954 Bandura, A. (1977). Self-efficacy: Toward a unifying theory of behavioral change. *Psychol.*
955 *Rev.* *84*, 191–215.
- 956 Bradley, M.M., Miccoli, L., Escrig, M.A., and Lang, P.J. (2008). The pupil as a measure of
957 emotional arousal and autonomic activation. *Psychophysiology* *45*, 602–607.
- 958 Brainard, D.H. (1997). The Psychophysics Toolbox. *Spat. Vis.* *10*, 433–436.
- 959 Bromberg-Martin, E.S., and Sharot, T. (2020). The value of beliefs. *Neuron* *106*, 561–565.
- 960 Charpentier, C.J., De Martino, B., Sim, A.L., Sharot, T., and Roiser, J.P. (2016a). Emotion-
961 induced loss aversion and striatal-amygdala coupling in low-anxious individuals. *Soc. Cogn.*
962 *Affect. Neurosci.* *11*, 569–579.
- 963 Charpentier, C.J., De Neve, J.E., Li, X., Roiser, J.P., and Sharot, T. (2016b). Models of

- 964 Affective Decision Making: How Do Feelings Predict Choice? *Psychol. Sci.* *27*, 763–775.
- 965 Charpentier, C.J., Bromberg-Martin, E.S., and Sharot, T. (2018). Valuation of knowledge
966 and ignorance in mesolimbic reward circuitry. *Proc. Natl. Acad. Sci. U. S. A.* *115*, E7255–
967 E7264.
- 968 Chib, V.S., Rangel, A., Shimojo, S., and O’Doherty, J.P. (2009). Evidence for a common
969 representation of decision values for dissimilar goods in human ventromedial prefrontal
970 cortex. *J. Neurosci.* *29*, 12315–12320.
- 971 Christianson, S.A. (2014). *The Handbook of Emotion and Memory: Research and Theory*
972 (Taylor & Francis).
- 973 Craig, A.D.B. (2003). Interoception: the sense of the physiological condition of the body.
974 *Curr. Opin. Neurobiol.* *13*, 500–505.
- 975 Craig, A.D.B. (2009). How do you feel — now? The anterior insula and human awareness.
976 *Nat. Rev. Neurosci.* *10*, 59–70.
- 977 Czekalla, N., Stierand, J., Stolz, D.S., Mayer, A. V., Voges, J.F., Rademacher, L., Paulus,
978 F.M., Krach, S., and Müller-Pinzler, L. (2021). Self-beneficial belief updating as a coping
979 mechanism for stress-induced negative affect. *Sci. Rep.* *11*, 17096.
- 980 Darby, R.S., and Harris, C.R. (2010). Embarrassment’s effect on facial processing. *Cogn.*
981 *Emot.* *24*, 1250–1258.
- 982 Diaconescu, A.O., Mathys, C., Weber, L.A.E., Kasper, L., Mauer, J., and Stephan, K.E.
983 (2017). Hierarchical prediction errors in midbrain and septum during social learning. *Soc.*
984 *Cogn. Affect. Neurosci.* *12*, 618–634.
- 985 Eickhoff, S.B., Stephan, K.E., Mohlberg, H., Grefkes, C., Fink, G.R., Amunts, K., and
986 Zilles, K. (2005). A new SPM toolbox for combining probabilistic cytoarchitectonic maps
987 and functional imaging data. *Neuroimage* *25*, 1325–1335.
- 988 Elder, J., Davis, T., and Hughes, B.L. (2021). Learning About the Self: Motives for
989 coherence and positivity constrain learning from self- relevant feedback. *PsyArXiv*.

- 990 Ertac, S. (2011). Does self-relevance affect information processing? Experimental evidence
991 on the response to performance and non-performance feedback. *J. Econ. Behav. Organ.* *80*,
992 532–545.
- 993 Feinberg, M., Willer, R., and Keltner, D. (2012). Flustered and faithful: Embarrassment as a
994 signal of prosociality. *J. Pers. Soc. Psychol.* *102*, 81–97.
- 995 Frijda, N.H. (1987). Emotion, cognitive structure, and action tendency. *Cogn. Emot.* *1*, 115–
996 143.
- 997 Garrett, N., González-Garzón, A.M., Foulkes, L., Levita, L., and Sharot, T. (2018). Updating
998 beliefs under perceived threat. *J. Neurosci.* *38*, 7901–7911.
- 999 de Gee, J.W., Knapen, T., and Donner, T.H. (2014). Decision-related pupil dilation reflects
1000 upcoming choice and individual bias. *Proc. Natl. Acad. Sci.* *111*, E618–E625.
- 1001 Gelman, A., and Rubin, D.B. (1992). Inference from iterative simulation using multiple
1002 sequences. *Stat. Sci.* *7*, 457–472.
- 1003 Hare, T.A., O’Doherty, J., Camerer, C.F., Schultz, W., and Rangel, A. (2008). Dissociating
1004 the role of the orbitofrontal cortex and the striatum in the computation of goal values and
1005 prediction errors. *J. Neurosci.* *28*, 5623–5630.
- 1006 Heimberg, R.G., Brozovich, F.A., and Rapee, R.M. (2010). A cognitive- behavioral model
1007 of social anxiety disorder: Update and extension. In *Social Anxiety: Clinical,*
1008 *Developmental, and Social Perspectives*, S.G. Hofmann, and P.M. DiBartolo, eds. (New
1009 York: NY: Elsevier), pp. 395–422.
- 1010 Hopkins, A.K., Dolan, R., Button, K.S., and Moutoussis, M. (2021). A Reduced Self-
1011 Positive Belief Underpins Greater Sensitivity to Negative Evaluation in Socially Anxious
1012 Individuals. *Comput. Psychiatry* *5*, 21.
- 1013 Hughes, B.L., and Zaki, J. (2015). The neuroscience of motivated cognition. *Trends Cogn.*
1014 *Sci.* *19*, 62–64.
- 1015 Izuma, K., Saito, D.N., and Sadato, N. (2008). Processing of Social and Monetary Rewards

1016 in the Human Striatum. *Neuron* 58, 284–294.

1017 Izuma, K., Saito, D.N., and Sadato, N. (2010). Processing of the incentive for social
1018 approval in the ventral striatum during charitable donation. *J. Cogn. Neurosci.* 22, 621–631.

1019 Kanske, P., Böckler, A., Trautwein, F.M., and Singer, T. (2015). Dissecting the social brain:
1020 Introducing the EmpaToM to reveal distinct neural networks and brain-behavior relations
1021 for empathy and Theory of Mind. *Neuroimage* 122, 6–19.

1022 Kaspar, K., and König, P. (2012). Emotions and personality traits as high-level factors in
1023 visual attention: a review. *Front. Hum. Neurosci.* 6, 1–14.

1024 Kelly, C., Toro, R., Di Martino, A., Cox, C.L., Bellec, P., Castellanos, F.X., and Milham,
1025 M.P. (2012). A convergent functional architecture of the insula emerges across imaging
1026 modalities. *Neuroimage* 61, 1129–1142.

1027 Keltner, D., and Potegal, M. (1997). Appeasement and reconciliation: Introduction to an
1028 aggressive behavior special issue. *Aggress. Behav.* 23, 309–314.

1029 King-Casas, B., Tomlin, D., Anen, C., Camerer, C.F., Quartz, S.R., and Montague, P.R.
1030 (2005). Getting to Know You: Reputation and Trust in a Two-Person Economic Exchange.
1031 *Science.* 308, 78–83.

1032 Koban, L., and Pourtois, G. (2014). Brain systems underlying the affective and social
1033 monitoring of actions: An integrative review. *Neurosci. Biobehav. Rev.* 46, 71–84.

1034 Koban, L., Schneider, R., Ashar, Y.K., Andrews-Hanna, J.R., Landy, L., Moscovitch, D.A.,
1035 Wager, T.D., and Arch, J.J. (2017). Social anxiety is characterized by biased learning about
1036 performance and the self. *Emotion* 17, 1144–1155.

1037 Koenig, S., Uengoer, M., and Lachnit, H. (2018). Pupil dilation indicates the coding of past
1038 prediction errors: Evidence for attentional learning theory. *Psychophysiology* 55, 1–12.

1039 Korn, C.W., Prehn, K., Park, S.Q., Walter, H., and Heekeren, H.R. (2012). Positively biased
1040 processing of self-relevant social feedback. *J. Neurosci.* 32, 16832–16844.

1041 Korn, C.W., Sharot, T., Walter, H., Heekeren, H.R., and Dolan, R.J. (2014). Depression is

1042 related to an absence of optimistically biased belief updating about future life events.
1043 *Psychol. Med.* *44*, 579–592.

1044 Kube, T., Kirchner, L., Lemmer, G., and Glombiewski, J.A. (2021). How the Discrepancy
1045 Between Prior Expectations and New Information Influences Expectation Updating in
1046 Depression—The Greater, the Better? *Clin. Psychol. Sci.*

1047 Kuzmanovic, B., Jefferson, A., and Vogeley, K. (2016). The role of the neural reward
1048 circuitry in self-referential optimistic belief updates. *Neuroimage* *133*, 151–162.

1049 Kuzmanovic, B., Rigoux, L., and Tittgemeyer, M. (2018). Influence of vmPFC on dmPFC
1050 predicts valence-guided belief formation. *J. Neurosci.* *38*, 7996–8010.

1051 Lindquist, K.A., Wager, T.D., Kober, H., Bliss-Moreau, E., and Barrett, L.F. (2012). The
1052 brain basis of emotion: a meta-analytic review. *Behav. Brain Sci.* *35*, 121–143.

1053 Lockwood, P.L., and Wittmann, M.K. (2018). Ventral anterior cingulate cortex and social
1054 decision-making. *Neurosci. Biobehav. Rev.* *92*, 187–191.

1055 Lockwood, P.L., Apps, M.A.J., Valton, V., Viding, E., and Roiser, J.P. (2016).
1056 Neurocomputational mechanisms of prosocial learning and links to Empathy. *Proc. Natl.*
1057 *Acad. Sci. U. S. A.* *113*, 9763–9768.

1058 Loewenstein, G. (2006). The Pleasures and Pains of Information. *Science.* *312*, 704–706.

1059 Loued-Khenissi, L., Pfeuffer, A., Einhäuser, W., and Preuschoff, K. (2020). Anterior insula
1060 reflects surprise in value-based decision-making and perception. *Neuroimage* *210*, 116549.

1061 Markus, H., and Wurf, E. (1987). The Dynamic Self-Concept: A Social Psychological
1062 Perspective. *Annu. Rev. Psychol.* *38*, 299–337.

1063 Marsh, H.W., and O’Neill, R. (1984). Self Description Questionnaire III: The Construct
1064 Validity of Multidimensional Self-Concept Ratings by Late Adolescents. *J. Educ. Meas.* *21*,
1065 153–174.

1066 Menon, V., and Uddin, L.Q. (2010). Saliency, switching, attention and control: a network
1067 model of insula function. *Brain Struct. Funct.* *214*, 655–667.

1068 Miller, R.S. (1996). *Embarrassment: Poise and Peril in Everyday Life* (New York: The
1069 Guilford Press).

1070 Morrison, A.S., and Heimberg, R.G. (2013). Social anxiety and social anxiety disorder.
1071 *Annu. Rev. Clin. Psychol.* *9*, 249–274.

1072 Müller-Pinzler, L., Gazzola, V., Keysers, C., Sommer, J., Jansen, A., Frässle, S., Einhäuser,
1073 W., Paulus, F.M., and Krach, S. (2015). Neural pathways of embarrassment and their
1074 modulation by social anxiety. *Neuroimage* *119*, 252–261.

1075 Müller-Pinzler, L., Czekalla, N., Mayer, A. V, Stolz, D.S., Gazzola, V., Keysers, C., Paulus,
1076 F.M., and Krach, S. (2019). Negativity-bias in forming beliefs about own abilities. *Sci. Rep.*
1077 *9*, 14416.

1078 Murty, V.P., Shermohammed, M., Smith, D. V., Carter, R.M., Huettel, S.A., and Adcock,
1079 R.A. (2014). Resting state networks distinguish human ventral tegmental area from
1080 substantia nigra. *Neuroimage* *100*, 580–589.

1081 Niv, Y., Edlund, J.A., Dayan, P., and O’Doherty, J.P. (2012). Neural prediction errors reveal
1082 a risk-sensitive reinforcement-learning process in the human brain. *J. Neurosci.* *32*, 551–
1083 562.

1084 O’Doherty, J.P. (2004). Reward representations and reward-related learning in the human
1085 brain: Insights from neuroimaging. *Curr. Opin. Neurobiol.* *14*, 769–776.

1086 Palminteri, S., Lefebvre, G., Kilford, E.J., and Blakemore, S.J. (2017). Confirmation bias in
1087 human reinforcement learning: Evidence from counterfactual feedback processing. *PLoS*
1088 *Comput. Biol.* *13*, e1005684.

1089 Phan, K.L., Wager, T., Taylor, S.F., and Liberzon, I. (2002). Functional neuroanatomy of
1090 emotion: A meta-analysis of emotion activation studies in PET and fMRI. *Neuroimage* *16*,
1091 331–348.

1092 Preuschoff, K. (2011). Pupil dilation signals surprise: evidence for noradrenaline’s role in
1093 decision making. *Front. Neurosci.* *5*, 1–12.

- 1094 Rescorla, R.A., and Wagner, A.R. (1972). A theory of Pavlovian conditioning: variations in
1095 the effectiveness of reinforcement and non reinforcement. In *Classical Conditioning II:
1096 Current Research and Theory*, A. Black, and W.F. Prokasy, eds. (New York: Appleton-
1097 Century-Crofts), pp. 64–99.
- 1098 Rigoux, L., Stephan, K.E., Friston, K.J., and Daunizeau, J. (2014). Bayesian model selection
1099 for group studies - revisited. *Neuroimage* *84*, 971–985.
- 1100 Rouault, M., Dayan, P., and Fleming, S.M. (2019). Forming global estimates of self-
1101 performance from local confidence. *Nat. Commun.* *10*, 1–11.
- 1102 Rouhani, N., and Niv, Y. (2021). Signed and unsigned reward prediction errors dynamically
1103 enhance learning and memory. *Elife* *10*, e61077.
- 1104 Ruff, C.C., and Fehr, E. (2014). The neurobiology of rewards and values in social decision
1105 making. *Nat. Rev. Neurosci.* *15*, 549–562.
- 1106 Rutledge, R.B., Dean, M., Caplin, A., and Glimcher, P.W. (2010). Testing the Reward
1107 Prediction Error Hypothesis with an Axiomatic Model. *J. Neurosci.* *30*, 13525–13536.
- 1108 Rutledge, R.B., Skandali, N., Dayan, P., and Dolan, R.J. (2014). A computational and neural
1109 model of momentary subjective well-being. *Proc. Natl. Acad. Sci.* *111*, 12252–12257.
- 1110 Rutledge, R.B., De Berker, A.O., Espenhahn, S., Dayan, P., and Dolan, R.J. (2016). The
1111 social contingency of momentary subjective well-being. *Nat. Commun.* *7*, 1–8.
- 1112 Schultz, W. (1998). Predictive Reward Signal of Dopamine Neurons. *J. Neurophysiol.* *40*.
- 1113 Schultz, W., Dayan, P., and Montague, P.R. (1997). A Neural Substrate of Prediction and
1114 Reward on JSTOR. *275*.
- 1115 Sedikides, C., and Gregg, A.P. (2008). Self-enhancement: Food for thought. *Perspect.
1116 Psychol. Sci.* *3*, 102–116.
- 1117 Sedikides, C., and Hepper, E.G.D. (2009). Self-Improvement. *Soc. Personal. Psychol.
1118 Compass* *3*, 899–917.
- 1119 Sharot, T. (2011). The optimism bias. *Curr. Biol.* *21*, R941–R945.

- 1120 Sharot, T., and Garrett, N. (2016). Forming beliefs: Why valence matters. *Trends Cogn. Sci.*
1121 *20*, 25–33.
- 1122 Sharot, T., Korn, C.W., and Dolan, R.J. (2011). How unrealistic optimism is maintained in
1123 the face of reality. *Nat. Neurosci.* *14*, 1475–1479.
- 1124 Späti, J., Chumbley, J., Brakowski, J., Dörig, N., Grosse Holtforth, M., Seifritz, E., and
1125 Spinelli, S. (2014). Functional lateralization of the anterior insula during feedback
1126 processing. *Hum. Brain Mapp.* *35*, 4428–4439.
- 1127 Sperduti, M., Delaveau, P., Fossati, P., and Nadel, J. (2011). Different brain structures
1128 related to self- and external-agency attribution: A brief review and meta-analysis. *Brain*
1129 *Struct. Funct.* *216*, 151–157.
- 1130 Stolz, D.S., Müller-Pinzler, L., Krach, S., and Paulus, F.M. (2020). Internal control beliefs
1131 shape positive affect and associated neural dynamics during outcome valuation. *Nat.*
1132 *Commun.* *11*, 1–13.
- 1133 Storbeck, J., and Clore, G.L. (2008). Affective Arousal as Information: How Affective
1134 Arousal Influences Judgments, Learning, and Memory. *Soc. Personal. Psychol. Compass* *2*,
1135 1824–1843.
- 1136 Tangney, J.P., Stuewig, J., and Mashek, D.J. (2007). Moral Emotions and Moral Behavior.
1137 *Annu. Rev. Psychol.* *58*, 345–372.
- 1138 Taylor, S.E., and Brown, J.D. (1988). Illusion and well-being: a social psychological
1139 perspective on mental health. *Psychol. Bull.* *103*, 193–210.
- 1140 Touroutoglou, A., Hollenbeck, M., Dickerson, B.C., and Feldman Barrett, L. (2012).
1141 Dissociable large-scale networks anchored in the right anterior insula subserve affective
1142 experience and attention. *Neuroimage* *60*, 1947–1958.
- 1143 Tzourio-Mazoyer, N., Landeau, B., Papathanassiou, D., Crivello, F., Etard, O., Delcroix, N.,
1144 Mazoyer, B., and Joliot, M. (2002). Automated anatomical labeling of activations in SPM
1145 using a macroscopic anatomical parcellation of the MNI MRI single-subject brain.

1146 Neuroimage *15*, 273–289.

1147 Ullsperger, M., Harsay, H.A., Wessel, J.R., and Ridderinkhof, K.R. (2010). Conscious
1148 perception of errors and its relation to the anterior insula. *Brain Struct. Funct.* *214*, 629–643.

1149 Vehtari, A., Mononen, T., Tolvanen, V., Sivula, T., and Winther, O. (2016). Bayesian leave-
1150 one-out cross-validation approximations for Gaussian latent variable models. *J. Mach.*
1151 *Learn. Res.* *17*, 3581–3618.

1152 Wager, T.D., and Feldman Barrett, L. (2017). From affect to control: Functional
1153 specialization of the insula in motivation and regulation. *BioRxiv*.

1154 Wallis, J.D., and Kennerley, S.W. (2010). Heterogeneous reward signals in prefrontal
1155 cortex. *Curr. Opin. Neurobiol.* *20*, 191–198.

1156 Will, G.-J., Moutoussis, M., Womack, P.M., Bullmore, E.T., Goodyer, I.M., Fonagy, P.,
1157 Jones, P.B., Rutledge, R.B., and Dolan, R.J. (2020). Neurocomputational mechanisms
1158 underpinning aberrant social learning in young adults with low self-esteem. *Transl.*
1159 *Psychiatry* *10*, 96.

1160 Williams, L.A., and DeSteno, D. (2008). Pride and Perseverance: The Motivational Role of
1161 Pride. *J. Pers. Soc. Psychol.* *94*, 1007–1017.

1162 Zhang, L., Lengersdorff, L., Mikus, N., Gläscher, J., and Lamm, C. (2020). Using
1163 reinforcement learning models in social neuroscience: Frameworks, pitfalls and suggestions
1164 of best practices. *Soc. Cogn. Affect. Neurosci.* *15*, 695–707.

1165