1    **Asymmetric effects of acute stress on cost and benefit learning**

2    Stella Voulgaropoulou[1*], Fasya Fauzani[1], Janine Pfirrmann[1], Claudia Vingerhoets[1, 2], Thérèse

3    van Amelsvoort[1], Dennis Hernaus[1]

4

5    [1] Department of Psychiatry & Neuropsychology, Maastricht University, Minderbroedersberg

6    4-6, 6211 LK Maastricht, The Netherlands

7    [2] Department of Radiology & Nuclear Medicine, Amsterdam University Medical Centre,

8    Location AMC, Meibergdreef 5, 1105 AZ Amsterdam, The Netherlands.

9

10    **Abstract**

11    Stressful events trigger a complex physiological reaction – the *fight-or-flight* response – that

12    can hamper flexible decision-making. Inspired by key neural and peripheral characteristics of

13    the fight-or-flight response, here we ask whether acute stress changes how humans learn

14    about costs and benefits. Participants were randomly exposed to an acute stress or no-stress

15    control condition after which they completed a cost-benefit reinforcement learning task.

16    Acute stress improved learning to maximize benefits (monetary rewards) relative to

17    minimising energy expenditure (grip force). Using computational modelling, we demonstrate

18    that costs and benefits can exert asymmetric effects on decisions when prediction errors that

19    convey information about the reward value and cost of actions receive inappropriate

20    importance; a process associated with distinct alterations in pupil size fluctuations. These

21    results provide new insights into learning strategies under acute stress – which,  depending on

22    the context, may be maladaptive or beneficial - and candidate neuromodulatory mechanisms

23    that could underlie such behaviour.

**Introduction**

24

25    Stress is ubiquitous in everyday life. From recurrent, brief, events (a work meeting, moving

26    to a new house) to major life events (armed combat, pandemic, financial crisis), humans are

27    continuously exposed to challenges in their daily environment. The immediate central and

28    peripheral physiological cascade triggered by such events, collectively termed the *fight-or-*

29    *flight (or acute stress) response* (Cannon, 1915), serves an allostatic role that enables

30    organisms to adequately respond to environmental demands (de Kloet, Joëls, & Holsboer,

31    2005). Although beneficial for survival, this allostatic process comes at a cost: stress-induced

32    redistributions of neural resources - e.g., towards vigilance or threat detection - may hamper

33    the deployment of strategies that support adaptive and optimal decision-making (Hermans,

34    Henckens, Joëls, & Fernández, 2014).

35         Optimal decisions essentially depend on the ability to rapidly learn from the positive

36    and negative outcomes of previous actions, also known as *reinforcement learning* (Niv,

37    2009). Considerable evidence now suggests that acute stress impairs aspects of reinforcement

38    learning (Carvalheiro, Conceição, Mesquita, & Seara-Cardoso, 2020; de Berker et al., 2016;

39    Raio, Hartley, Orederu, Li, & Phelps, 2017). Acute stress, among others, modulates the

40    impact of positive outcomes on future decisions - both positively and negatively - (Berghorst,

41    Bogdan, Frank, & Pizzagalli, 2013; Carvalheiro et al., 2020; Lighthall, Gorlick, Schoeke,

42    Frank, & Mather, 2013; Petzold, Plessow, Goschke, & Kirschbaum, 2010), likely driven by

43    changes in reward sensitivity and the signalling of reward prediction errors (RPEs)

44    (Berghorst et al., 2013; Carvalheiro et al., 2020; Huys, Pizzagalli, Bogdan, & Dayan, 2013);

45    putatively dopaminergic teaching signals that represent the mismatch between actual and

46    expected outcomes, which are used to flexibly adjust behaviour (Niv, 2009; Rescorla, 1972).

47    Alterations in the influence of RPEs on future decisions play a key role in the development of

48    motivational impairments, which are frequently observed in behavioural disorders associated

49    with repeated and/or prolonged stress exposure (Huys et al., 2013).

50          Intuitive as it is, the notion that the impact of acute stress on (potentially maladaptive)

51    decisions *primarily* involves changes in how reward value influences action may be

52    oversimplified. Decisions are not only motivated by appetitive properties; they equally

53    depend on the – cognitive (e.g., mental effort) or physical (e.g., energy) – *cost* associated

54    with actions (Hauser, Eldar, & Dolan, 2017; Pessiglione, Vinckier, Bouret, Daunizeau, & Le

55    Bouc, 2017; Schmidt, Lebreton, Cléry-Melin, Daunizeau, & Pessiglione, 2012). Expectations

56    about action costs are also updated according to a prediction error rule (Skvortsova, Degos,

57    Welter, Vidailhet, & Pessiglione, 2017; Skvortsova, Palminteri, & Pessiglione, 2014)

58    (henceforth "effort" prediction errors; EPEs), which due to the aversive and resource-

59    consuming nature of effort, optimal learners should utilize to minimize effort expenditure.

60    When decisions involve a potential cost *and* benefit, the former is subtracted from the latter

61    to compute a "net" or subjective decision value (i.e., effort-discounted reward value) (Klein-

62    Flügge, Kennerley, Friston, & Bestmann, 2016; Skvortsova et al., 2017; Skvortsova et al.,

63    2014). Notably, stress exposure impairs cost-benefit decisions in rodents when learning is not

64    explicitly required (Friedman et al., 2017; Shafiei, Gray, Viau, & Floresco, 2012). Moreover,

65    in a reinforcement learning context, acute stress blocks the flexible updating of aversive

66    value (Raio et al., 2017), an inherent property of costly actions. These results suggest that

67    decisions during acute stress may involve a complex shift in reinforcement learning strategies

68    that serve to balance the cost versus benefits of decisions; a hypothesis that hitherto has

69    remained unexplored.

70          Although computationally similar in nature, distinct neural correlates of RPEs (e.g.,

71    striatal subdivisions, ventromedial prefrontal cortex [vmPFC]) and EPEs (e.g., parietal

72    cortex, insula, dorsomedial PFC) can be observed in cost-benefit reinforcement learning

73  paradigms (Hauser et al., 2017; Skvortsova et al., 2014). The ascending dopaminergic (e.g.,

74  RPEs, action cost, reward value) (Schultz, Dayan, & Montague, 1997; Skvortsova et al.,

75  2017; Yohn et al., 2016), noradrenergic (e.g., mobilizing energy) (Pessiglione et al., 2017;

76  Varazzani, San-Galli, Gilardeau, & Bouret, 2015) and serotonergic (e.g., aversive value,

77  overcoming action costs) (H. E. den Ouden et al., 2015; Meyniel et al., 2016)

78  neuromodulatory systems, moreover, encode partly dissociable aspects of goal-directed

79  actions that involve learning about costs and benefits, which together support optimal

80  decision-making. These observations are noteworthy because the initial fight-or-flight

81  response triggers a large-scale reorganization of brain networks that is driven by alterations in

82  the firing mode of midbrain dopaminergic ventral tegmental area and noradrenergic *locus*

83  *coeruleus* neurons (Arnsten, 2015; Hermans et al., 2014); neurons that signal prediction

84  errors (Steinberg et al., 2013) and that are also responsive to reward value, action cost and

85  energy expenditure (Del Arco, Park, & Moghaddam, 2020; Varazzani et al., 2015). Thus,

86  catecholaminergic mechanisms that are recruited by the fight-or-flight response may

87  differentially impact cost and benefit reinforcement learning, resulting in a potential scenario

88  in which costs and benefits exert asymmetric influences on decisions.

89      As mentioned above, the central (i.e., neural) effects of acute stress trigger a shift in

90  cognitive strategies, including reinforcement learning. The peripheral counterpart of the acute

91  stress response, however, mobilizes the energy (i.e., adrenaline-mediated glucose release (de

92  Kloet et al., 2005; Russell & Lightman, 2019)) that is required to exert effortful actions

93  aimed at preserving homeostasis (Cannon, 1915). Therefore, decision-making and learning

94  policies regarding physical costs may be *especially* susceptible to stress: both via

95  computational (neural) mechanisms that support learning about and representation of action

96  cost, as well as peripheral mechanisms that co-determine the amount of available energy that

97  can be directed towards effortful actions. Indeed, preliminary evidence suggests that acute

98    stress alters the willingness to exert physical effort for rewards (Bryce & Floresco, 2016) and

99    reward-associated cues in a Pavlovian-instrumental transfer context (Pool, Brosch,

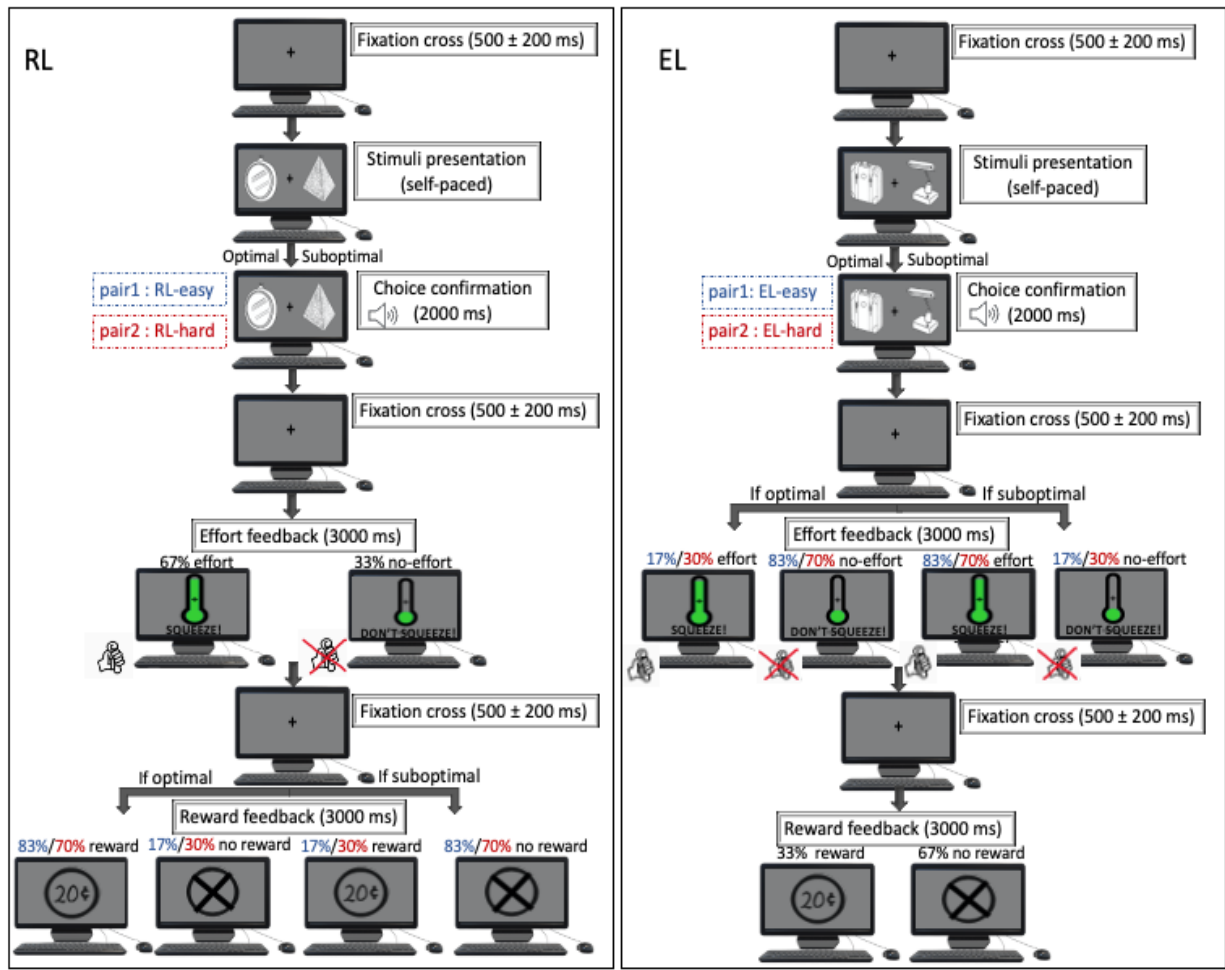100   Delplanque, & Sander, 2015).

101        How acute stress impacts reinforcement learning involving costs *and* benefits has not

102   been investigated to date in humans. Based on the above considerations, we expect that

103   computationally frugal learning strategies, in concert with increased energy availability,

104   during acute stress should asymmetrically impact cost versus benefit learning. Using an acute

105   stress-induction paradigm, a cost-benefit learning paradigm and computational model of cost-

106   benefit reinforcement learning (Skvortsova et al., 2017; Skvortsova et al., 2014), we

107   demonstrate that acute stress asymmetrically prioritizes reward (maximization) learning over

108   physical effort (minimization) learning. Better benefit versus cost learning results from a

109   stress-induced change in the influence of RPEs versus EPEs on future decisions, and is

110   associated with altered pupil encoding of RPEs, EPEs, and subjective decision value. These

111   results reveal how neural and peripheral mechanisms that support the fight-or-flight response

112   may facilitate a shift in reinforcement learning strategies that confers strategic benefits during

113   acutely stressful situations (e.g., ignoring high action costs to achieve a desirable outcome),

114   yet might also give rise to maladaptive behaviour (e.g., stress-induced relapse in substance

115   users).

116 **Results**

117 **Experiment design**

118 Healthy human participants were randomly assigned to the acute stress (19 males/21 females;

119 age M=23.48, SD=3.94) or no-stress control condition (18 males/22 females; age M=23.80,

120 SD=4.23) of the Maastricht Acute Stress Task (MAST) (Smeets et al., 2012), a validated

121 psychological and physical stress-induction paradigm (see Materials and Methods).

122 Immediately post-MAST and within the confines of the acute stress response (Hermans et al.,

123 2014), all participants completed a ~40 minute probabilistic cost-benefit reinforcement

124 learning paradigm, adapted from Skvortsova et al. (Skvortsova et al., 2017; Skvortsova et al.,

125 2014), in which they learned to select stimuli with high reward value (20 Eurocents) and

126 avoid stimuli with high action cost (exerting grip force above a pre-calibrated individual

127 threshold of 50% maximum voluntary contraction for 3000ms), followed by a surprise test

128 phase. A detailed overview of the paradigm is provided in Figure 1 and the Materials and

129 Methods. Pupil size was continuously recorded while participants performed the task (see

130 Materials and Methods).

131    **Figure 1**



132

133    **Reward maximisation/action cost minimization reinforcement learning task**.

134    Visual depiction of the learning phase. Participants were presented with four distinct stimulus

135    pairs, and all stimuli were associated with a predetermined chance of a €0.20 monetary

136    reward (versus no reward) and a chance of having to exert physical effort (grip force) using a

137    dynamometer (versus no grip force required). Stimulus-outcome probabilities were yoked in

138    such a way that, for a given pair, two stimuli *only* differed in the probability of earning a

139    reward ("reward learning", RL, left) or the probability of having to exert effort ("effort

140    learning", EL, right). That is, for RL (left)/EL (right) pairs, reward/effort outcomes were

141    choice-dependent, respectively (see "Reward feedback" for RL and "Effort feedback" for EL

142    for outcome contingencies). For RL pairs, effort outcomes were independent of choice and

143    fixed (see "Effort feedback" for RL), while for EL pairs, reward outcomes were independent

144    of choice and fixed (see "Reward feedback" for EL). Percentages in blue and red refer to

145    outcomes for the Easy RL/EL and Hard RL/EL pair, respectively.

146
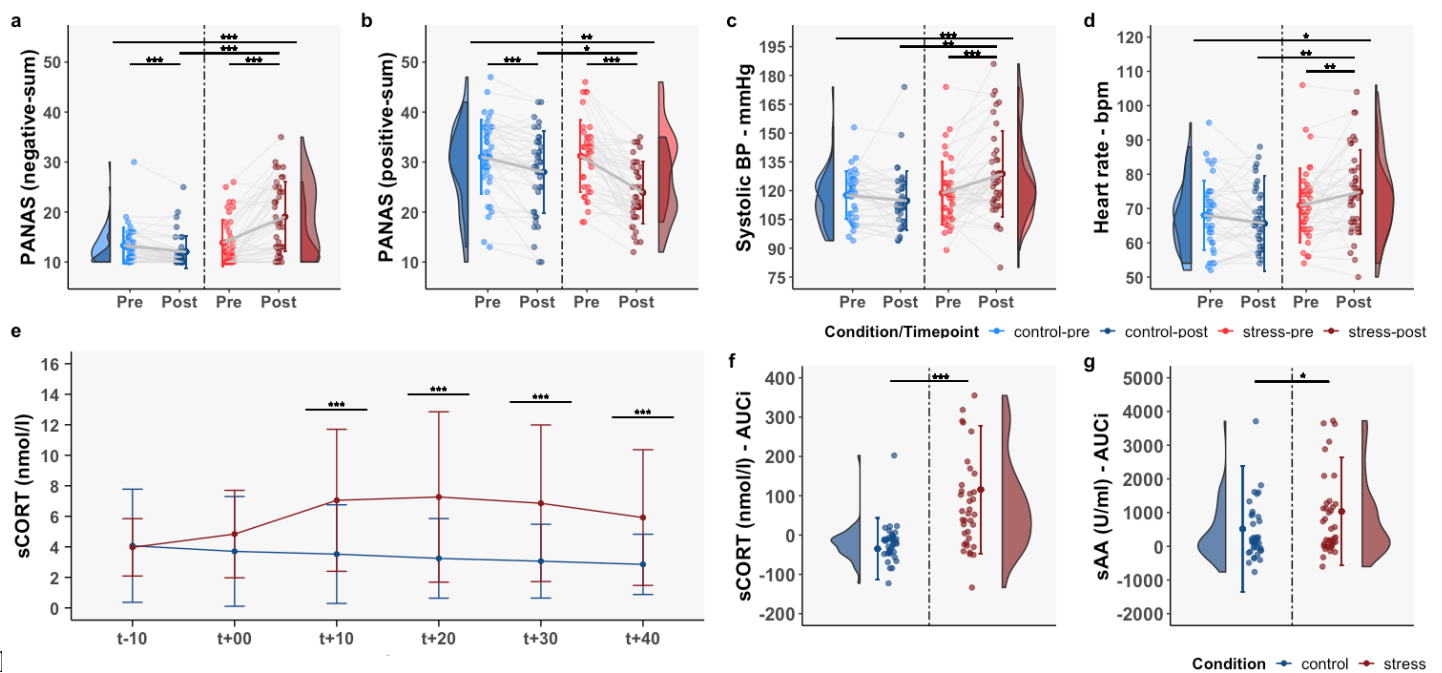
147    **Acute stress manipulation**

148    We first ascertained whether the acute stress manipulation was successful. Subjective stress,

149    physiological and neuroendocrine measurements are displayed in Figure 2. Acute stress and

150    no-stress control groups did *not* differ on physiological, subjective stress, or neuroendocrine

151    measurements pre-MAST (all *p-values*>0.05). We observed significant Condition-by-Time

152    interactions for subjective stress ratings [PANAS negative: $F(1,78)=52.66$, $p<0.001$,

153    $n^2_G=0.10$; PANAS positive: $F(1,78)=9.82$ $p=0.002$, $n^2_G=0.02$] and physiological measures

154    [systolic blood pressure (SBP): $F(1,78)=15.50$, $p<0.001$, $n^2_G= 0.04$; heart rate: $F(1,78)=6.83$,

155    $p=0.011$, $n^2_G= 0.02$]. Simple main effect analyses revealed that only the acute stress group

156    exhibited pre-to-post *increases* in negative affect [control pre-post: $t(39)=4.21$, $p<0.001$;

157    stress pre-post: $t(39)=-6.17$, $p<0.001$; control-stress post-MAST: $t(55.1)=-5.78$, $p<0.001$],

158    and greater pre-to-post *decreases* in positive affect [control pre-post: $t(39)=4.09$, $p<0.001$;

159    stress pre-post: $t(39)=6.45$, $p<0.001$; control-stress post-MAST: $t(72.8)=2.53$, $p=0.014$]

160    (Figure 2). Similarly, only the acute stress group exhibited stress-induced *increases* in SBP

161    [control pre-post: $t(39)=1.60$, $p<0.117$; stress pre-post: $t(39)=-3.66$, $p<0.001$; control-stress

162    post-MAST: $t(69.1)=-3.27$, $p=0.002$] and heart rate [control pre-post: $t(39)=1.21$, $p=0.234$;

163    stress pre-post: $t(39)=-2.78$, $p=0.008$; control-stress post-MAST: $t(76.9)=-3.14$, $p=0.002$]

164    (Figure 2a-d).

165         An expected Condition-by-Time interaction was found for salivary cortisol (sCORT)

166    responses [$F(5,390)=18.05$, $p<0.001$, $n^2_G= 0.04$], with the acute stress group displaying

167    greater sCORT levels 10 min post-MAST and onwards (all *p-values*<0.01). We additionally

168    observed a main effect of Condition on sCORT area-under-the-curve with respect to increase:

8

169     (AUCi) (Pruessner, Kirschbaum, Meinlschmid, & Hellhammer, 2003) ($t(56.32)$=-5.28,

170     $p<0.001$) and salivary alpha-amylase (sAA) AUCi ($t(67.45)$=-2.50, $p$=0.015; after excluding

171     one extreme outlier from the control group), suggesting greater sCORT and sAA levels in

172     response to acute stress (Figure 2e-g). These results confirm that the MAST robustly induced

173     stress on all levels of inquiry.

174

175     **Figure 2**



177     **Neuroendocrine, physiological and subjective stress ratings.**

178     PANAS negative (**a**) and positive (**b**) subscale sum scores, systolic blood pressure (mmHg:

179     millimetres of mercury; **c**) and heart rate (bpm: beats per minute; **d**) are displayed for no-

180     stress control (blue) and acute stress (red) groups separately for pre (light blue/red) and post

181     (dark blue/red) MAST time points. SCORT responses for both conditions across the 6

182     timepoints are displayed in panel **e** ("$t_{+00}$" represents the first post-MAST measurement, and

183     the start of the reward maximization/action cost reinforcement learning paradigm; "$t_{-10}$"

184     represent a baseline sample). Panel **f** and **g** show AUCi for sCORT (nmol/l: nanomoles per

185     litre) and sAA (U/mL: Units per millilitre) responses for both MAST conditions. Significant

9

186    differences are denoted by asterisks (*: $p < 0.05$, **: $p < 0.01$, ***: $p < 0.001$). In the upper

187    panel, the top line denotes a significant Condition-by-Time interaction; lower lines represent

188    simple main effects of Condition or Time.

189

190    **Participants use reinforcement learning to optimize decisions**

191    Next, we investigated whether participants in both conditions exhibited evidence of

192    reinforcement learning to optimize actions, which in this paradigm should be reflected by an

193    increased tendency to select stimuli with high reward value and avoid stimuli with high

194    action cost as a function of increasing number of stimulus pair presentations (i.e., "time").

195    This intuition was confirmed by a main effect of Time on two distinct trial types: trials on

196    which participants could learn to accumulate frequent rewards while the probability of effort

197    (action cost) was kept constant (RL; selecting the stimulus more frequently associated with

198    €0.20) [control: $F(2,78)=10.16$, $p<0.001$, $n^2{}_G= 0.06$; stress: $F(2,78)=20.44$, $p<0.001$, $n^2{}_G=$

199    0.17] and trials on which participants could learn to frequently avoid effort while the

200    probability of reward was kept constant (EL; selecting the stimulus more frequently

201    associated with avoidance of physical energy expenditure) [control: $F(2,78)=12.35$, $p<0.001$,

202    $n^2{}_G= 0.07$; stress: $F(2,78)=9.76$, $p<0.001$, $n^2{}_G= 0.05$]. We additionally observed greater than

203    chance-level performance ($\geq0.5$) on both trial types during the final part of the task

204    (presentation 21-30; all $p$-values$<0.001$ Figure Supplement 1).

205

206    **Asymmetric cost-benefit reinforcement learning during acute stress**

207    After having observed evidence for reward (maximization) learning and action cost

208    (minimization) learning, we tested our key assumption; that acute stress would induce a

209    reprioritization in learning to maximize reward value versus learning to minimize action cost.

210    Crucially, we observed a significant Condition-by-Trial Type interaction [$F(1,78)=6.53$,

10

211    $p$=0.013, $n^2{}_G$= 0.039] (Figure 3a) with pairwise comparisons indicating that the acute stress

212    group performed significantly better on RL than EL trials [$t$(39)=5.40, $p$<0.001], while the

213    no-stress control group performed similarly on both trial types [$t$(39)=1.01, $p$=0.320]. A main

214    effect of Condition on RL-EL accuracy difference scores [$t$(74.02)=-2.55, $p$=0.013] (Figure

215    3b) and one-sample t-tests revealed that RL-EL accuracy difference scores were significantly

216    greater than zero in the acute stress group but not in the no-stress controls. Simple main

217    effects of Condition on RL [$t$(65.9)=-1.75, $p$=0.085] and EL [$t$(77.5)=1.80, $p$=0.076]

218    performance showed numerical trends for group differences that failed to reach significance.

219            When we included participants that still performed at chance level at the end of the

220    learning phase (see Participants) in the Condition-by-Trial Type interaction analysis, the

221    interaction remained significant [$F$(1,91)=7.30, $p$=0.035, $n^2{}_G$= 0.04], with the acute stress

222    group displaying better RL vs. EL performance ($t$(46)=5.83, $p$<0.001), while no-stress

223    controls performed similarly on both trial types ($t$(45)=1.24, $p$=0.22). Participants in the acute

224    stress group outperformed participants in the no-stress control group on RL ($t$(91)=2.04,

225    $p$=0.04), but no simple main effect of Condition was observed for EL ($t$(91)=-1.67, $p$=0.1).

226    These participants were excluded for all other analyses reported below.
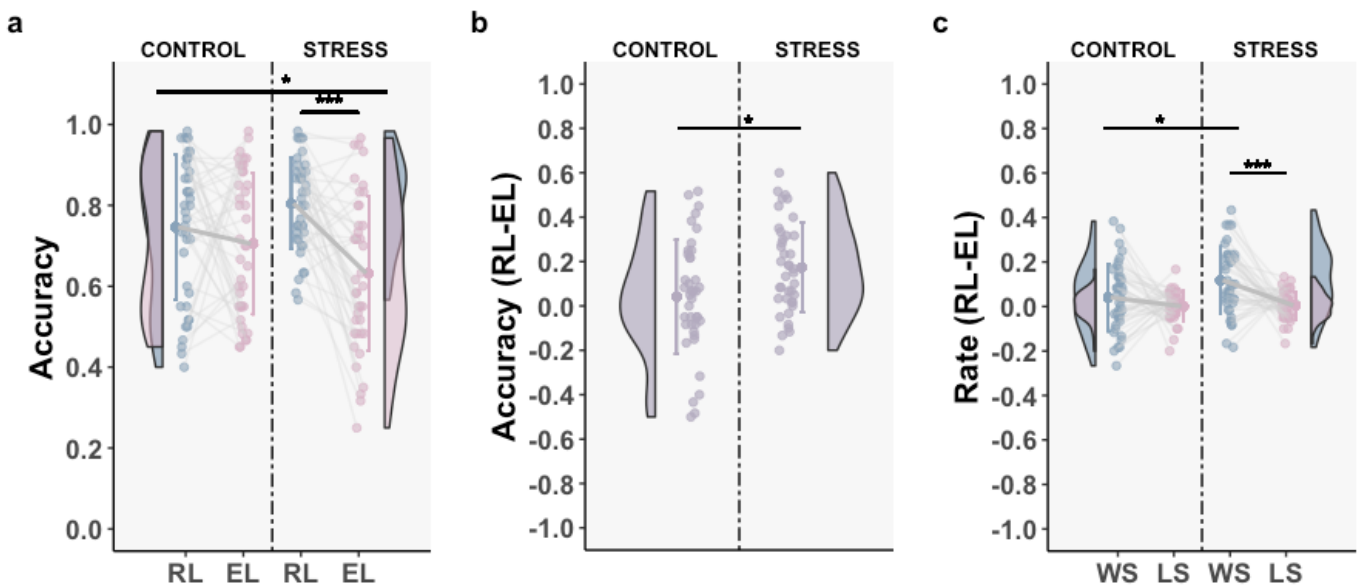
227            The use of different reinforcement probabilities for each RL and EL pair (Figure 1)

228    allowed us to discern whether the observed pattern of results reflected a *specific* change in

229    cost (EL) versus benefit (RL) reinforcement learning, or a more general impairment in

230    reinforcement learning for more difficult stimulus-response associations. We found no

231    evidence for the latter scenario in Condition-by-Trial Type-by-Difficulty [$F$(1,78)=1.05,

232    $p$=0.310] or Condition-by-Difficulty [$F$(1,78)=0.36, $p$=0.549; Figure Supplement 2]

233    interaction analyses.

234            When we investigated the use of win-stay (resampling a stimulus following a positive

235    outcome) and lose-shift (switching to the other stimulus following a negative outcome)

236    strategies (Hanneke E. M. den Ouden et al., 2013), we observed no significant Condition-by-

237    Strategy interaction [$F(1,78)=2.99$, $p=0.087$, $n^2_G= 0.03$]. However, *post hoc* comparisons

238    revealed that participants in the acute stress condition [$t(39)=-3.73$, $p<0.001$] but not those in

239    the no-stress control condition [$t(39)=-1.30$, $p=0.200$], exhibited different win-stay compared

240    to lose-shift (difference) rates. Separate Condition (main effect) analyses indicated that acute

241    stress participants compared to no-stress controls were more likely to win-stay for rewards

242    (RL trials) than for avoidance of action cost (EL trials) [$t(78)=-2.28$, $p=0.025$], but not for

243    lose-shifting for reward omissions (RL trials) compared to exerting effort (EL trials)

244    [$t(77.7)=-0.23$, $p=0.820$]. Differences in win-stay rates for RL compared to EL trials (one-

245    sample t-test) were greater than zero for the acute stress [$t(39)=4.91$, $p<0.001$] but not the no-

246    stress control group [$t(39)=1.68$, $p=0.101$] (Figure 3c).

247        Taken together, our model-free results indicate that acute stress leads to a

248    reinforcement learning strategy that favours learning to maximise reward value over

249    minimisation of action cost, which based on analyses of win-stay/lose-shift rates, could be

250    attributed to increased sensitivity to positive reinforcement (i.e., reward delivery) compared

251    to negative reinforcement (i.e., avoidance of physical effort).

**Figure 3**



**Acute stress leads to improved benefit versus cost learning.**

Panel **a**: Average accuracy (choices of the optimal stimulus) for RL and EL trials, for each

condition separately. Panel **b**: RL-EL average accuracy difference scores. Panel **c**:

Win-stay ($WS_{RL}$-$WS_{EL}$) and lose-shift ($LS_{RL}$-$LS_{EL}$) difference scores for each condition

separately. Means ± SD, individual data points, distribution and frequency of the data are

displayed. In panel **a**, the top line indicates a significant Condition-by-Trial type interaction.

Significant differences are denoted by asterisks (*: $p < 0.05$, **: $p < 0.01$, ***: $p < 0.001$).

Source files of task performance data used for the analyses are available in the Figure 3 –

Source Data 1.

**Asymmetric cost-benefit reinforcement learning biases actions in acute stress subjects**

During a surprise 64-trial test phase (D. Hernaus, Gold, Waltz, & Frank, 2018), we asked

participants to discriminate original and novel combinations of stimuli on the basis of reward

value or action cost without receiving feedback (n=16 trials for original combinations; n=48

for novel combinations; see Materials and Methods). The surprise test phase allowed us to

269 assess learned choice tendencies without having to arbitrarily choose a given number of final

270 learning phase trials, during which participants may still learn. This approach also allowed us

271 to assess the degree to which learned tendencies would carry over to novel contexts.

272   First, both groups chose the optimal (most rewarding/effort avoiding) stimulus on

273 surprise test phase trials involving the original four pairs [one-sample t-test against chance;

274 control$_{RL}$: $t(39)=8.73$, $p<0.001$; control$_{EL}$: $t(39)=3.72$, $p=0.002$; stress$_{RL}$: $t(39)=13.54$,

275 $p<0.001$; stress$_{EL}$: $t(39)=4.47$, $p<0.001$], confirming that both groups had developed a

276 preference for the optimal stimulus.

277   Although we observed no Condition-by-Trial type (reward value, action cost

278 discrimination) interaction or main effects of Condition for novel stimulus combinations

279 [$F(1,78)=1.10$, $p=0.298$, $n^2{}_G= 0.01$; stress vs. controls reward value discrimination: $t(75.9)=$

280 $0.15$, $p=0.878$; stress vs. controls action cost discrimination: $t(78)= 1.77$, $p=0.080$], pairwise

281 comparisons revealed that the acute stress group performed better on reward discrimination

282 compared to action cost discrimination trials [$t(39)=-2.23$, $p=0.032$], while no-stress controls

283 performed similarly on both trial types [$t(39)=-0.87$, $p=0.387$]. These results provide some

284 evidence that a reward maximisation-over-action cost minimisation reinforcement learning

285 policy might bias future actions in novel contexts (Figure Supplement 3).

286

287 **Computational cost-benefit reinforcement learning model: Model Fitting, Selection,**

288 **Demonstrations, and Simulations**

289 To uncover latent mechanisms by which acute stress affects cost-benefit reinforcement

290 learning, we turned to computational cognitive modelling. Trial-by-trial choices of

291 participants were fit to all candidate models described in the Materials and Methods (see

292 Model Space). To calculate model fit, log likelihood was updated trial-wise by the log of the

293 probability of the observed choice, calculated via a softmax rule (see Materials and Methods,

294    equation 6), and best-fitting parameters were identified using fmincon in MATLAB  v.2019B

295    (Mathworks, Natick, MA, USA).

296         Bayesian Model Selection (BMS; spm_BMS function in SPM12,

297    http://www.fil.ion.ucl.ac.uk/spm/software/spm12/) using the Akaike Information Criterion

298    (AIC) as a fit statistic that penalizes for the number of model parameters (Myung, Tang, &

299    Pitt, 2009), suggested that the 2LR_γ model was the most likely model, as indicated by the

300    protected exceedance probability (pxp, $\varphi = 0.99$) (Rigoux, Stephan, Friston, & Daunizeau,

301    2014) and expectation of the posterior [p(r|y), 0.70] (Figure 4a for p(r|y) of all candidate

302    models). We note that 2LR_γ remained the most likely model when we considered additional

303    models with greater redundancy and/or lesser biological plausibility (e.g., models with all

304    combinations of reward value/action cost discounting *and* weight parameters).

305         The 2LR_γ model contains separate learning rates that weight the importance of RPEs

306    and EPEs ($\alpha_R$, $\alpha_E$), an action cost discounting parameter ($\gamma$), and an inverse temperature

307    parameter ($\beta$), which in previous work could account for performance on a conceptually

308    similar cost-benefit learning task (Skvortsova et al., 2014). To demonstrate the effect of

309    changes in parameters values on choice preferences within the 2LR_γ architecture, we first

310    simulated choices from 50 artificial agents (averaged across 10 repetitions) performing the

311    reward maximization/action cost minimization reinforcement learning task using a range of

312    parameter values. As expected, greater values of $\alpha_R$ and $\alpha_E$ primarily impacted the speed of

313    RL and EL choice preferences, while low values of $\gamma$ lead to asymmetric choice preferences

314    through discounting of action cost, and lower values of $\beta$ lead to non-selective increases in

315    random sampling (Figure 4b).

316         In *post hoc* simulations, i.e., generating participant choices using the obtained

317    parameters, we additionally observed moderate-to-high correlations between simulated and

318    empirical RL/EL for the acute stress and no-stress control group [$\rho_{RL\_control} = 0.55$, $p < 0.01$;

15

319    $\rho_{RL\_stress}$ = 0.84, $p$ < 0.01; $\rho_{EL\_control}$ = 0.56, $p$ < 0.01; $\rho_{EL\_stress}$ = 0.77, $p$ < 0.01; see Figure

320    Supplement 4], although the canonical performance difference in RL versus EL accuracy was

321    not selective to the acute stress group [$t_{control}(39)$=-6.72, $p$<0.001; $t_{stress}(39)$=-6.01, $p$<0.001].

322    However, after we fixed $\beta$ and $\gamma$ to group-level averages, to better demonstrate the effect of

323    group differences in the learning rate parameters, we recovered a small but significant

324    simulated difference in RL versus EL performance for the acute stress group [$t(39)$=2.27,

325    $p$=0.029], which was not predicted in the no-stress control group [$t(39)$=0.91, $p$=0.367]

326    (Condition-by-Trial Type interaction: [$F(1,78)$= 0.77, $p$=0.38, $n^2_G$= 0.006]) (Figure 4c for

327    empirical versus simulated data, averaged across 100 repetitions per subject).

328         Importantly, even if a given model is the most likely one based on model fitting and

329    *post hoc* simulation results from the entire sample, there is still the possibility that different

330    models can better explain task performance in the no-stress control and acute stress

331    condition. When repeating BMS for each condition separately, 2LR_$\gamma$ was the most likely

332    model in the no-stress control group [$\varphi$=0.99, p(r|y)=0.83], while for acute stress subjects

333    2LR_$\gamma$ was not convincingly the most likely model [$\varphi$=0.47 p(r|y)=0.46]. Here, the 2LR

334    model (containing $\alpha_R$, $\alpha_E$, and $\beta$ parameters) was equally likely to be the optimal model

335    [$\varphi$=0.53 p(r|y)=0.47]. Post-hoc simulations from the 2LR model also correlated with actual

336    data, both for no-stress control [$\rho_{RL}$ = 0.65, $p$ < 0.001; $\rho_{EL}$ = 0.60, $p$ < 0.001] and acute stress

337    participants [$\rho_{RL}$ = 0.81, $p$ < 0.001; $\rho_{EL}$ = 0.75, $p$ < 0.001].

338         Similar to the 2LR_$\gamma$ model (results discussed in next section), the 2LR model

339    seemingly also explained stress-induced changes in cost-benefit reinforcement learning via

340    changes in learning rates; in the 2LR model, the acute stress group exhibited greater values of

341    $\alpha_R$ versus $\alpha_E$ (t(39)=2.65, $p$=0.01), while no-stress control subjects did not ($t(39)$=0.69,

342    $p$=0.50) (Condition-by-Learning Rate interaction: [$F(1,78)$=2.88, $p$=0.094, $n^2_G$= 0.01]. The

343    difference in learning rates between 2LR_$\gamma$ (where $\alpha_R$ and $\alpha_E$ are similar for the acute stress

16

344     group, see next section) and 2LR (where $\alpha_R > \alpha_E$ for the acute stress group) can be explained

345     by the absence of discounting parameter γ: 2LR is a special case of 2LR_γ, where γ=1, and

346     thus asymmetric effects of acute stress on reward value maximization and action cost

347     minimization can only be explained by *dissimilarity* in learning rates.

348         Although the effects of acute stress on reward value and action cost learning rates are

349     opposite in 2LR_ γ versus 2LR architectures, these results bolster our confidence in the

350     overall model space, as well as the interpretation that acute stress primarily impacts reward

351     value and action cost learning rates, and *not* discounting. The observations that I) 2LR_ γ fit

352     better in the entire group of participants, II) 2LR is fully contained within the 2LR_ γ model,

353     and III) 2LR_ γ displayed good recoverability (see *below)* motivated our choice to focus on

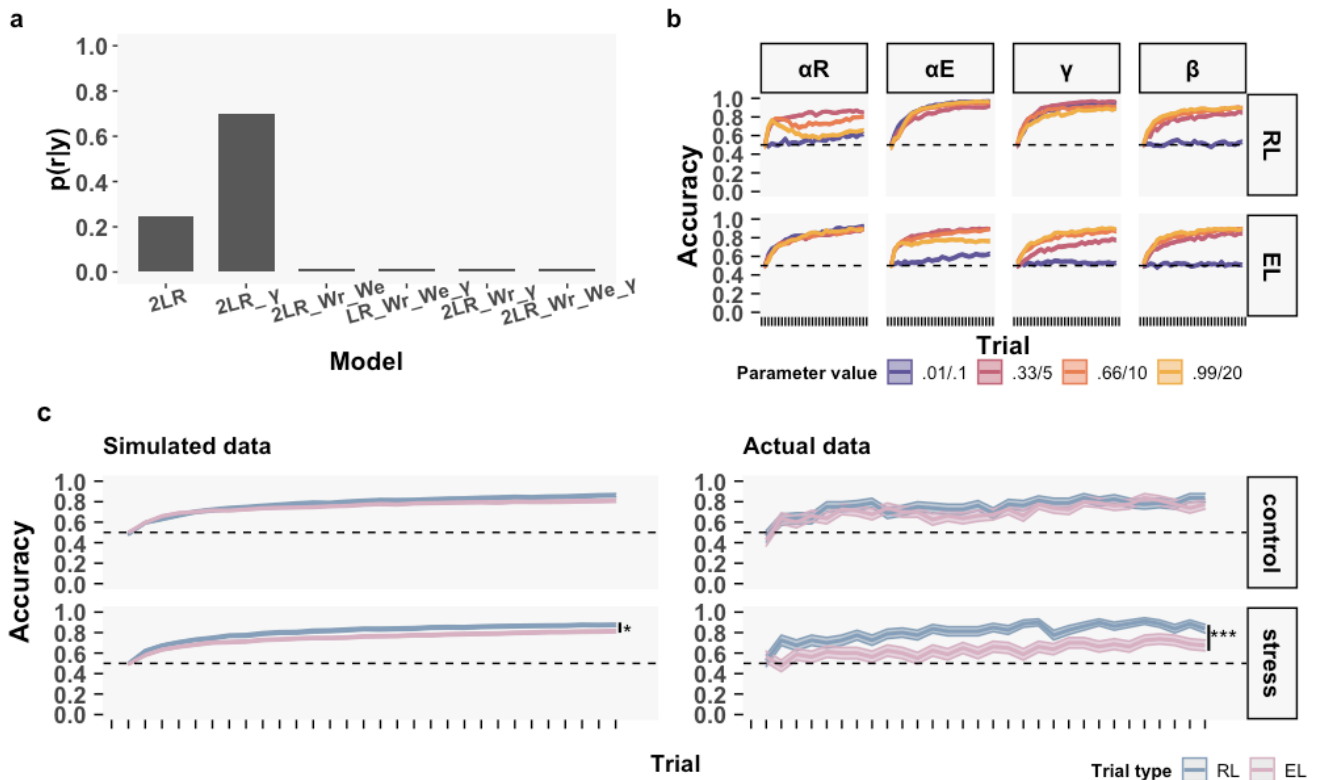354     the 2LR_γ model.

355         In model recoverability analyses i.e., re-fitting the simulated data from the model to

356     all candidate models (Wilson & Collins, 2019), BMS confirmed that the simulated 2LR_γ

357     data (that is, simulations *without* fixed parameters) were most likely to be generated from

358     2LR_γ [φ=0.99, p(r|y)=0.71].

359         To assess the stability of 2LR_γ parameters, we repeated model fitting using a

360     Bayesian hierarchical model fitting approach consisting of two steps, as described previously

361     (Daw, 2011; Frey, Frank, & McCabe, 2019). In the first step we fit the 2LR_γ model to trial-

362     wise choices to obtain subject-specific parameters; in a second step we again fit the model to

363     trial-wise choices, but this time we used the group-level average and covariance matrix of

364     every parameter as priors, thereby shrinking the parameter search space. Motivated by recent

365     work showing that group-specific priors, compared to a single prior for the entire sample, can

366     better account for between-group differences in task performance, as well as improve

367     parameter robustness and recoverability (Valton, Wise, & Robinson, 2020), we used separate

368     mean and covariance matrices for the acute stress and no-stress control groups.

369        Highly similar parameter estimates were obtained after hierarchical fitting (for

370        parameter estimates after Bayesian hierarchical model fitting see Figure Supplement 5).

371        Similar to *post hoc* simulations using parameters from the non-hierarchically fit 2LR_$\gamma$

372        model, we observed moderate-to-high correlations between empirical and simulated data

373        using parameters obtained from the hierarchically fit model [$\rho_{RL\_control} = 0.65$, $p < 0.01$;

374        $\rho_{RL\_stress} = 0.84$, $p < 0.01$; $\rho_{EL\_control} = 0.37$, $p = 0.02$; $\rho_{EL\_stress} = 0.78$, $p < 0.01$; see Figure

375        Supplement 6]. All in all, these results confirm parameter stability within the 2LR_$\gamma$

376        architecture.

377        In light of model fitting results, post-hoc simulations, model and parameter

378        recoverability analyses, we used parameters and trial-by-trial predictions of the non-

379        hierarchically fit 2LR_ $\gamma$ model in all analyses reported below.

380     **Figure 4**



381

382     **Model selection, demonstrations, and *post hoc* simulations of the winning model.**

383     Panel **a**: Expectation of the posterior for all candidate models. Panel **b**: Model

384     demonstrations. To demonstrate how different parameter values within the 2LR_$\gamma$

385     architecture impact choice preferences for the optimal stimulus ("accuracy"), $\alpha_R$, $\alpha_E$, and $\gamma$

386     were set to 0.01/0.33/0.66/0.99, while $\beta$, a non-linear parameter, was set to 0.1/5/10/20.

387     Parameter effects were always demonstrated for a single parameter (columns), while all other

388     parameter values were kept constant ($\alpha_R$ and $\alpha_E$=0.25, $\gamma$=1, $\beta$=25). Greater values of $\alpha_R$ and

389     $\alpha_E$ selectively increase the speed with which the agent develops a preference for the optimal

390     RL and EL stimulus, respectively. Lower values of $\gamma$ produce an asymmetric decision-

391     making policy that emphasises reward value over action cost, leading to better performance

392     on RL versus EL trials, while greater values of $\gamma$ correct this asymmetric choice bias. Finally,

393     greater $\beta$ values lead to more deterministic sampling of optimal stimuli. Panel **c**: Post-hoc

394     simulations after fixing β and γ to group-level averages. Coloured lines represent mean ± SD.

395     Dashed lines denote chance level (0.5). *: $p < 0.05$, **: $p < 0.01$, ***: $p < 0.001$.

396

397     **Acute stress selectively reduces the difference between reward and action cost learning**
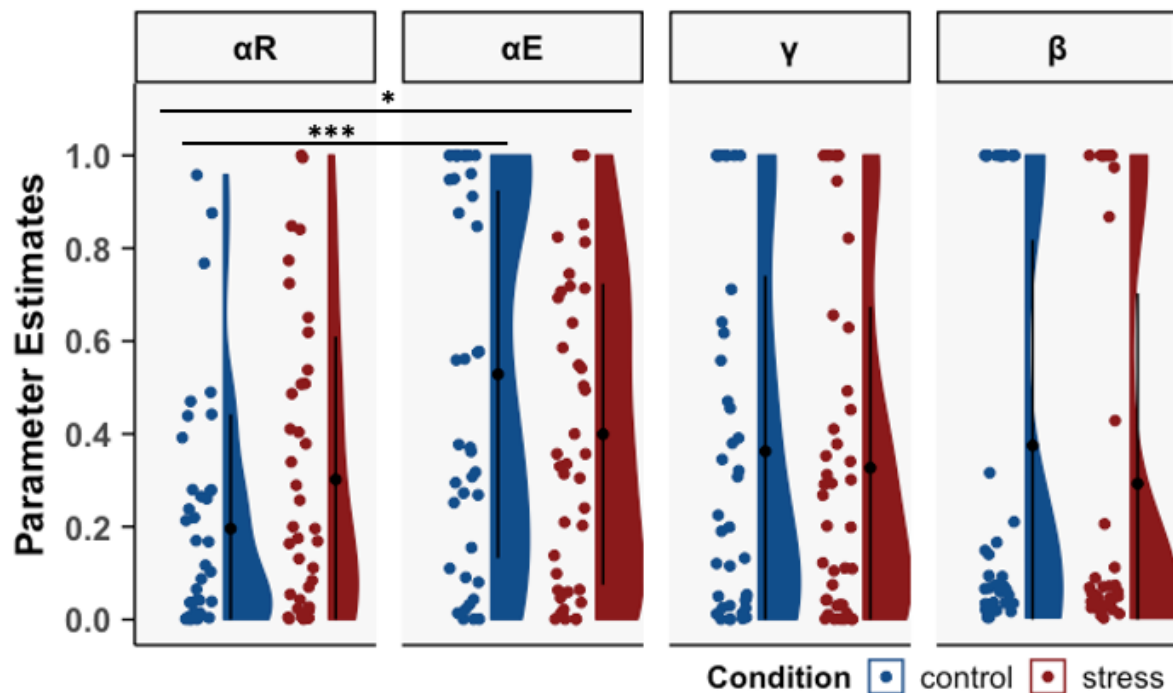
398     **rates**

399     Comparing 2LR_γ parameters between conditions, we observed a significant Condition-by-

400     Learning Rate ($\alpha_R$, $\alpha_E$) interaction [$F(1,78)$= 6.42, $p$=0.01, $n^2{}_G$= 0.03; 95% highest density

401     interval (HDI) for Bayesian mixed ANOVA = -0.405 to -0.023, mean= -0.219]. with greater

402     EPE relative to RPE learning rates in no-stress control participants [$t(39)$=-4.75, $p$<0.001],

403     while learning rates in the acute stress group did not significantly differ [$t(39)$=-1.61,

404     $p$=0.116]. No between-group differences in $\alpha_R$ and $\alpha_E$ or in the other parameters (γ, β) were

405     observed (all *p-values*>0.05) (Figure 5).

406        Paradoxically, symmetric reward value and action cost learning rates in the presence

407     of lower values of γ will lead to more efficient RL compared to EL. This is because lower

408     values of γ bias decisions towards reward value (via *greater* discounting of action cost) and

409     similar absolute values of $\alpha_R$/$\alpha_E$ will not counteract this bias. *Asymmetric* learning rates

410     ($\alpha_E$>$\alpha_R$) in combination with lower values of γ, however, will lead to more symmetric

411     performance on RL and EL trials via more efficient updating of action cost versus reward

412     expectations. This interpretation is supported by our demonstration of model parameters

413     (Figure 4b) and *post hoc* simulations (Figure 4c), as well as the observation that lower values

414     of γ (i.e., greater action cost discounting) were associated with greater learning rate

415     asymmetry ($\alpha_E$>$\alpha_R$; more efficient EL) in no-stress controls (ρ=-0.40, $p$=0.040), who

416     displayed similar RL and EL performance. These results demonstrate that, in a context where

417     all decisions involve a potential cost and benefit, acute stress selectively reduces the

418     difference between EPE and RPE learning rates, while leaving action cost discounting and

419    choice stochasticity unaffected. The direction of the change in learning rates (i.e., greater

420    similarity) implies a stress-induced failure to modulate learning rates in the service of

421    overcoming an asymmetric choice bias that emphasises reward value.

422        In analyses using posterior parameters obtained from the hierarchically fit model, we

423    recovered the key Condition-by-Learning Rate interaction (95% HDI for Bayesian mixed

424    ANOVA = -0.406 to -0.128, mean=-0.269) [and acute stress and no-stress control subjects

425    differed from each other on $\alpha_E$ (95% HDI = 0.0841 to 0.281, mean=0.183) but not $\alpha_R$ (95%

426    HDI = -0.186 to 0.0102, mean=-0.0872)] (Figure Supplement 7). Similar to the non-

427    hierarchically fit parameters, acute stress and control subjects did not differ on posterior

428    estimates of $\gamma$ (95% HDI = -0.0914 to 0.139, mean=0.0182) and $\beta$ (95% HDI = -0.0331 to

429    0.213, mean=0.0895) .

430

**Figure 5**



**Acute stress reduces the difference between reward and effort prediction error learning rates.**

Free parameters ($\alpha R$, $\alpha E$, $\gamma$, $\beta$) of the winning 2LR_$\gamma$ model for both groups. Black lines denote means $\pm$ SD, dots represent individual data points, and the violin-like shape denotes distribution and frequency of the data. *: $p < 0.05$, **: $p < 0.01$, ***: $p < 0.001$.

**Pupil size fluctuations track asymmetric cost-benefit reinforcement learning during acute stress**

We employed pupillometry to better understand whether task-relevant computational processes may be encoded by fluctuations in pupil dilation, which are thought to be controlled by ascending midbrain modulatory systems that play a role in value-based decision-making and the acute stress response (Arnsten, 2015; Hermans et al., 2011; Joshi, Li, Kalwani, & Gold, 2016).
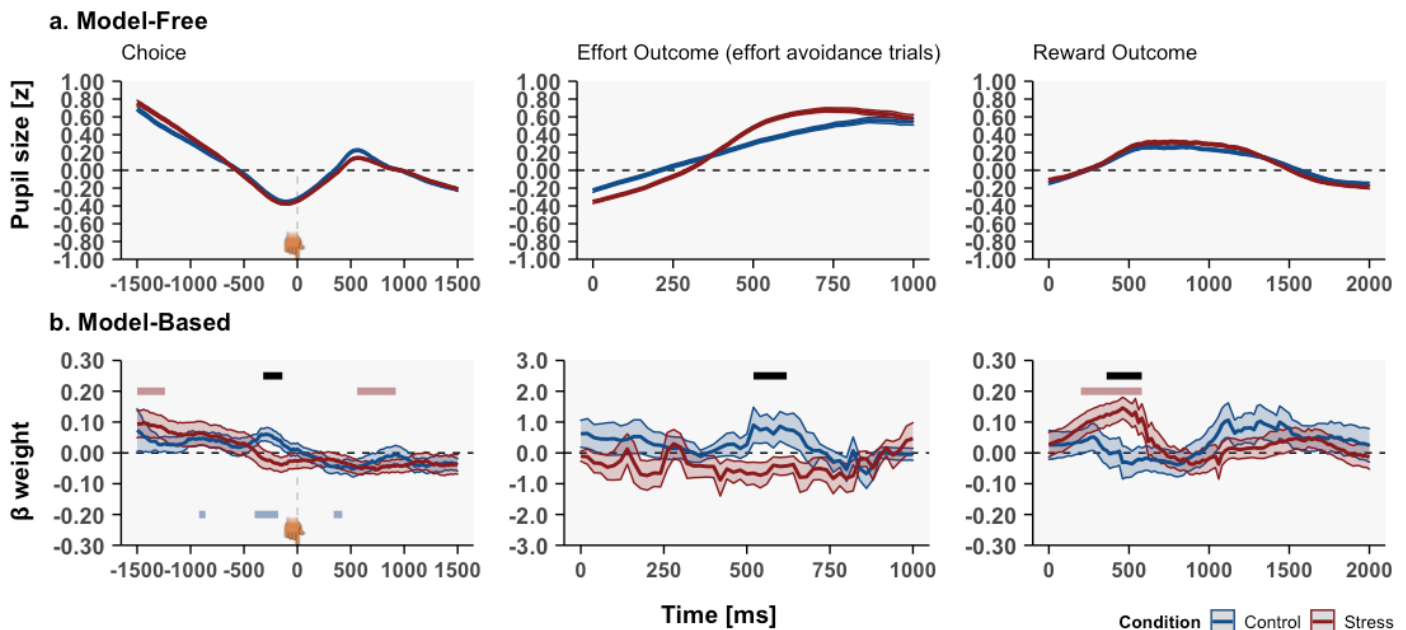
22

446        In model-free analyses – that is, comparing bins of pupillometry data between

447    conditions - we observed no main effect of Condition on pupil size fluctuations during

448    choice, effort outcome, and reward outcome epochs, suggesting that acute stress was not

449    associated with more general changes in pupil size (all bin-level $p>0.05$; Figure 6, a. Model-

450    free; Figure Supplement 8 for all effort trials).

451        Next, we conducted model-based pupillometry analyses (Lawson, Bisby, Nord,

452    Burgess, & Rees, 2020) to understand how trial-wise estimates of computational processes of

453    interest were encoded by fluctuations in pupil size. These analyses revealed effects of

454    Condition on pupil encoding of subjective decision value, EPEs and RPEs (Figure 6, b.

455    Model-based) in a manner commensurate with task performance results. First, immediately

456    prior to the stimulus choice, acute stress reduced pupil encoding of subjective decision value,

457    as evidenced by the absence of an association between pupil size and subjective decision

458    value (control>stress; stress n.s., control>0). Second, briefly after the presentation of effort

459    avoidance outcomes, both groups exhibited different pupil size-EPE associations, with no-

460    stress controls showing a non-significant numerically positive association between pupil size

461    and action cost prediction errors (control>stress, both groups n.s. different from 0). In model-

462    based analyses using all effort outcome trials, we were not able to uncover group differences

463    in pupil encoding of EPEs, which were likely eclipsed by prominent grip force-related effects

464    on pupil size (Figure Supplement 8). Third, during the reward outcome phase, acute stress

465    participants exhibited greater positive associations between pupil size and RPEs compared to

466    no-stress controls (stress>control, stress>0, control n.s.). The average pupil size-RPE slope

467    for bins in which no-stress control and acute stress participants differed (Figure 6b) correlated

468    significantly with stress-induced changes in SBP [$\rho_{stress}(38)= -0.41$, $p(permutation)=0.019$]

469    and PANAS negative affect changes [$\rho_{stress}= -0.46$, $p(permutation)=0.005$] in the stress group.

470        Crucially, group differences in pupil encoding of subjective decision value, EPEs, and

471    RPEs imply that the ascending neuromodulatory systems may have facilitated a stress-

472    induced shift in asymmetric cost versus benefit learning.

473

474    **Figure 6**



476    **Model-based analysis reveals altered pupil encoding of prediction errors and decision**

477    **value during acute stress .**

478    **a**. Model-free analyses of pupil size during choice, effort outcome, and reward outcome

479    phase revealed no main effect of Condition (no-stress control, acute stress). **b.** Model-based

480    analyses revealed a stress-induced shift in pupil encoding of subjective decision value (left),

481    action cost prediction errors (middle) and reward prediction errors (right). Black line

482    indicates significant main effect of Condition; blue and red line indicate significance against

483    zero for no-stress control and acute stress groups, respectively (cluster and bin level

484    $\alpha_{permute}$<0.05, 2000 permutations). Group differences in pupil encoding of action cost and

485    reward prediction errors were observed at similar times (note the x-axis differences for effort

486    outcome and reward outcome epochs). Source files of pupillometry data used for the analyses

487    are available in Figure 6 – Source Data 2.

488

489    **Discussion**

490    Stress-induced alterations in adaptive decision-making are commonly studied using

491    paradigms that isolate positive and negative reinforcement, such as the receipt of a reward or

492    avoidance of a loss. However, it remains poorly understood how acute stress affects the

493    complex process that entails learning about costs *and* benefits, a critical and pervasive feature

494    of everyday decisions. Participants completed a paradigm in which all actions (stimulus

495    choices) contained a potential cost (exerting physical effort) and a financial benefit (€0.20).

496    Crucially, acute stress induced a shift in reinforcement learning strategies that improved

497    maximization of monetary rewards relative to minimisation of energy expenditure. When

498    presented with novel stimulus arrangements and in the absence of feedback, individuals in

499    the acute stress condition, moreover, exhibited better discrimination of stimulus reward value

500    compared to action cost.

501         Relative improvements in reward versus action cost learning align well with previous

502    reports of enhanced reward learning during acute stress (Byrne, Cornwall, & Worthy, 2019;

503    Lighthall et al., 2013; Petzold et al., 2010), although such effects may depend on stressor

504    timing (Joëls, Pu, Wiegert, Oitzl, & Krugers, 2006), stressor type (Carvalheiro et al., 2020),

505    and/or sample characteristics (Evans & Hampson, 2015; Morris & Rottenberg, 2015). While

506    reports on action cost learning during acute stress are scarce, acute exposure to stress in

507    rodents impairs cost-benefit decisions via a selective change in sensitivity to physical effort, a

508    process mediated by corticotropin-releasing factor and dopamine (Bryce & Floresco, 2016).

509    Our analyses of win-stay/lose-shift rates indicate that asymmetric cost-benefit learning can be

510     driven by a relative increase in the sensitivity to monetary gains compared to the avoidance

511     of costly deterrents.

512        How might maximization of reward value take precedence over minimisation of

513     action cost? Acute stress leads to a redistribution of finite cognitive resources (Hermans et

514     al., 2014): this process limits the availability of computationally intensive strategies,

515     including working memory (Otto, Raio, Chiang, Phelps, & Daw, 2013; Qin, Hermans, van

516     Marle, Luo, & Fernández, 2009) and goal-directed instrumental actions (Lars Schwabe &

517     Wolf, 2011). Assuming that acute stress does not merely increase random responding - which

518     we verified via the choice stochasticity model parameter - a computationally cheap heuristic

519     in our task should present itself as better learning for one modality over the other. Increased

520     energy availability (Hermans et al., 2014), insensitivity to aversive stimuli (Timmers et al.,

521     2018), and impaired aversive value updating (Raio et al., 2017) under stress may have

522     reduced the ability – or urgency – to dedicate cognitive resources to strategies that minimize

523     action cost. Importantly, effort expenditure increases the perceived value of rewards

524     (Hernandez Lallement et al., 2014; Inzlicht, Shenhav, & Olivola, 2018). Thus, frequent

525     expenditure of physical effort, due to suboptimal action cost learning, may increase the

526     perceived value of rewards, and thus tilt learning towards the maximization of reward value.

527        Using a computational model of reinforcement learning (Skvortsova et al., 2017;

528     Skvortsova et al., 2014), we confirmed that biased cost-benefit learning can arise when

529     inappropriate (i.e. more similar) importance is afforded to teaching signals that convey

530     information about reward value (RPEs) and action cost (EPEs). Humans display presumably

531     instinctive biases, such as more efficient learning from better-than-expected outcomes

532     (Lefebvre, Lebreton, Meyniel, Bourgeois-Gironde, & Palminteri, 2017) (compared to worse-

533     than-expected outcomes) and asymmetric "Go"/approach learning (Guitart-Masip et al.,

534     2012) (compared to "No-Go"/avoidance learning), the latter being a bias that is also

535      modulated by acute stress (de Berker et al., 2016). From this perspective, no-stress controls,

536      who assigned greater importance to EPEs than RPEs, may have used a computationally

537      costly learning strategy that provides counterweight to a decision-making policy that is

538      biased towards the reward value of actions (captured by action cost discounting parameter $\gamma$).

539      Paradoxically, when decisions are by default tilted towards reward value, similar reward and

540      action cost learning rates will facilitate reward learning but hamper action cost learning.

541      Reduced learning rate asymmetry in the presence of action cost discounting may therefore

542      represent a computational reformulation of a heuristic that is employed when cognitively

543      demanding learning strategies are unavailable and the policy towards energy expenditure is

544      more liberal, such as during acute stress.

545          Importantly, stress-induced changes in task performance may crucially depend on the

546      release of catecholamines in neural circuits that support motivation and learning. Dopamine's

547      actions at D1 and D2 receptors in the basal ganglia mediate approach and avoidance learning

548      (Frank, Seeberger, & Reilly, 2004), and acute stress can improve associative learning by

549      augmenting reward-evoked DA bursts in selective striatal subdivisions (Stelly, Tritley,

550      Rafati, & Wanat, 2020). Dopamine's enhancement via L-DOPA administration, moreover ,

551      improves reward but not action cost learning (Skvortsova et al., 2017). To the degree that

552      pupillometry can be considered a proxy measure of activity of ascending neuromodulatory

553      systems, these findings are consilient with greater encoding of RPEs by pupil size

554      fluctuations during acute stress. Negative correlations between SBP and PANAS negative

555      affect and RPE-pupil size slopes suggest that primarily moderately stressed participants

556      displayed a preference for maximizing reward value, which might be consistent with an

557      inverted U-shape relationship between cognitive performance and DA transmission which is

558      modulated by stress (Arnsten & Goldman-Rakic, 1998; Baik, 2020). Noradrenaline, however,

559      mobilizes available energy to complete effortful actions and *locus coeruleus* neurons track

560    energy expenditure (Varazzani et al., 2015). Stress-induced sAA concentrations, increased

561    heart rate, and group differences in the association between pupil size fluctuations and EPEs

562    all point to the involvement of the noradrenaline system. Thus, our model-based pupillometry

563    and stress-induction results hint at stress-sensitive dopaminergic and noradrenergic

564    mechanisms that may regulate cost and benefit learning, which could be explored in future

565    work using targeted pharmacological approaches.

566        The results presented here may improve understanding of stress-related

567    psychopathology. While asymmetric cost-benefit learning during acute stress may be

568    beneficial to reach a desired goal state (e.g., safety) despite high action cost, such strategies

569    could also be maladaptive. For example, stress exposure can lead to drug or smoking relapse

570    (L. Schwabe, Dickinson, & Wolf, 2011), a context in which reward value and action cost

571    may be misaligned. Cost-benefit reinforcement learning may provide a useful framework to

572    test hypotheses regarding stress-related impairments in learning and decision-making.

573        Some study limitations need to be acknowledged. First, pupil dilation associated with

574    effort expenditure greatly reduced our power to detect robust associations between EPE

575    encoding and pupil size fluctuations. Future studies should, therefore, consider a temporal

576    delay between effort outcome and effort expenditure phases. Second, while our

577    computational model was able to recover overall task performance patterns in both groups,

578    such effects were subtle and dependent on the contribution of other (non-learning)

579    parameters, which may highlight the importance of interindividual differences in model

580    parameters.

581        To summarize, we present evidence of asymmetric effects of acute stress on cost

582    versus benefit reinforcement learning during acute stress, which computational analyses of

583    task behaviour explain as a failure to assign appropriate importance to RPEs versus EPEs,

584    and our model-based pupillometry tentatively link to activity of ascending midbrain

28

585     neuromodulatory systems. These results highlight for the first time how learning under acute

586     stress can be tilted in favour of acquiring good things and away from the avoidance of costly

587     things.

588

589     **Materials and Methods**

590     **Participants**

591     Adult participants were recruited via paper and online advertisements. All participants were

592     screened for a DSM-5 psychiatric and/or neurological disorder, substance use, endocrine

593     and/or vascular disorder, abnormal BMI (>40 or <18), smoking and drinking (>10

594     cigarettes/units per week), psychotropic medication use (lifetime) and hormonal

595     contraceptive use (current; female participants only). All participants completed the ~2-hour

596     experiment between 12:00h and 18:00h to minimize diurnal cortisol fluctuations (Bailey &

597     Heitkemper, 2001). Participants were instructed to refrain from alcohol (starting the evening

598     before the day of the experiment), smoking, food, caffeine intake, strenuous physical activity

599     and brushing their teeth (all >2 hr prior to experiment), which was verified verbally at the

600     start of the session. Four participants were excluded due to an equipment failure (n=4). Three

601     participants quit during stress-induction (n=2) or task procedures (n=1). Because chance-

602     level performance on reinforcement learning tasks might indicate a successful manipulation,

603     a lack of motivation, or a failure to comprehend the task instructions, participants that

604     performed at or below chance level (0.5) on both RL and/or both EL pairs near the end of the

605     experiment (final 10 presentations) were excluded (n=13; 6 acute stress, 7 no-stress control).

606     Including these participants did not alter our key finding that acute stress was associated with

607     asymmetric cost versus benefit learning (see Results). Pupillometry and neuroendocrine data

608     were not processed further for these participants. The study was approved by the ethics

609     committee of the Faculty of Psychology and Neuroscience, Maastricht University (ERCPN-

29

610     197_03_08_2018) and carried out in accordance with the Declaration of Helsinki.

611     Participants were remunerated in gift vouchers or research participation credits. Task

612     earnings were paid out in gift vouchers.

613

**Acute stress induction**

615     The MAST is a validated stress-induction paradigm combining both psychological and

616     physiological stressors, and robustly increases neuroendocrine, physiological, and subjective

617     indices of acute stress (Smeets et al., 2012). During a 5-min preparation phase, participants

618     are instructed about the upcoming task via oral and visually displayed instructions, followed

619     by a 10-min stress-induction phase consisting of alternating blocks of cold-water immersion

620     (non-dominant hand; 2°C) and backward counting in steps of 17 (while receiving negative

621     evaluative feedback from an experimenter), with a (non-recording) camera continuously

622     directed at the participant's face, which was displayed to the participant on a second display.

623     During the MAST no-stress control condition, participants immerse their hand in lukewarm

624     water (36°C) and perform simple mental arithmetic, e.g., counting from 1 to 25, without

625     receiving feedback or fake camera recordings.

626

**Neuroendocrine, physiological and subjective stress measurements**

628     sCORT and sAA were collected to measure stress-induced increases in hypothalamic-

629     pituitary-adrenal (HPA) axis and sympathetic-adrenal-medullary (SAM) axis activity,

630     respectively (Dickerson & Kemeny, 2004; Koh, Ng, & Naing, 2014). Saliva samples were

631     obtained using synthetic Salivette® devices (Sarstedt, Etten-Leur, the Netherlands) during 3-

632     min sampling periods at 6 time points. A baseline sample was collected 10 min prior to the

633     MAST (baseline: $t_1 = t_{-10}$) and five samples post-MAST ($t_2 = t_{+00}$, $t_3 = t_{+10}$, $t_4 = t_{+20}$, $t_5 = t_{+30}$, $t_6 = t_{+40}$). SAA assessments were obtained only for $t_1$-$t_4$, due to the rapid decay of sAA post-stress

635     induction (Dennis Hernaus, Quaedflieg, Offermann, Casales Santa, & van Amelsvoort, 2018;

636     Nater et al., 2005). For all participants, $t_2$ marked the starting point for the reward value

637     maximisation/action cost minimisation task. Samples were stored at -20ºC immediately after

638     completion of each session. SCORT and sAA levels were determined using a commercially

639     available luminescence immune assay kit (IBL, Hamburg, Germany) and kinetic reaction

640     assay (Salimetrics, Penn State, PA, USA), respectively.

641          Systolic blood pressure (SBP) and heart rate (HR) as an index of autonomic nervous

642     system (ANS) arousal (Schubert et al., 2009; Wright, O'Brien, Hazi, & Kent, 2014) were

643     assessed at $t_1$ and $t_2$ using an OMRON M4-I  blood pressure monitor (OMRON Healthcare

644     Europe B.V., Hoofddorp, The Netherlands). Subjective affect ratings were assessed at $t_1$  and

645     $t_2$ using the 20-item Positive and Negative Affect Scale (PANAS) (Watson, Clark, &

646     Tellegen, 1988).

647

648     **Reward maximization versus action cost minimization reinforcement learning task**

649     All participants completed a probabilistic stimulus selection paradigm during which they

650     learned to select stimuli with high reward value (20 Eurocents) and avoid stimuli with high

651     action cost (exerting force above a pre-calibrated individual threshold for a duration of

652     3000ms). This reinforcement learning task is conceptually similar to a previously-validated

653     probabilistic action selection task that has been employed to study the neural signatures of

654     reward and effort prediction errors and dopaminergic drug effects on reward-effort

655     computations (Skvortsova et al., 2017; Skvortsova et al., 2014). The paradigm was designed

656     in PsychoPy v3.0.0b11 (Peirce et al., 2019) and presented on a 24″ monitor (iiyama ProLite

657     b2483HSU). Physical effort (in mV/kgf) was registered using a hand-held dynamometer in

658     combination with a transducer amplifier (DA100C) and data acquisition system (MP160; all

659     manufactured by BIOPAC Systems, Inc). Individual effort thresholds used throughout the

660    task were obtained by calculating 50% of each participant's maximal voluntary contraction

661    (MVC) (Le Heron et al., 2018) reached over three calibration trials by squeezing the

662    dynamometer with the dominant hand.

663         On each of 120 trials, participants chose between two paired distinct black-and-white

664    images ("stimuli") that were probabilistically associated with both the receipt of a monetary

665    reward and exertion of physical effort (see Figure 1 for a graphical overview). At trial onset,

666    a fixation cross flanked by two images was presented; participants chose one image by

667    pressing the V/B button for the left/right option, respectively. A 440Hz/600Hz tone for

668    left/right choice (200ms) was presented to confirm the participant's choice. Next, a

669    thermometer with the command "SQUEEZE" or "DON'T SQUEEZE" was displayed. If

670    participants were required to exert effort, they were instructed to squeeze the dynamometer

671    until the mercury level reached the top. The mercury bar only moved if participants exerted

672    above-threshold levels of force and stopped moving if exerted force fell below. The

673    cumulative above-threshold time was 3000ms. If no effort production was required, an

674    animation of a rising mercury bar was displayed (3000ms). Finally, a screen was presented

675    showing either a €0.20 coin or a crossed-out coin, indicating no reward (3000ms).

676         Participants learned to choose the optimal (most reward or most effort avoiding)

677    stimulus for four distinct image pairs, 30 presentations each, with yoked reward and action

678    cost contingencies. For 2/4 pairs, participants could regularly acquire rewards by selecting

679    one (optimal) stimulus over the (suboptimal) other (henceforth, "reward learning"/RL pairs),

680    while the probability of having to exert effort was identical for both stimuli. For the other two

681    pairs, choices of one stimulus were more frequently followed by the avoidance of effort

682    ("effort learning"/EL pairs), while the probability of reward was kept constant between both.

683    For all pairs, the probability of the stimulus property that was *kept constant* (reward/effort)

684    was set to a 33.3% chance of positive outcome upon selection (reward/ effort avoidance) and

685    66.6% chance of negative outcome (no reward/effort).

686        To assess whether any acute stress effects on reward maximization (measured using

687    RL pairs) and effort cost minimization (measured using EL pairs) learning were potentially

688    mediated by task difficulty, we employed different difficulty levels for each RL and EL pair.

689    That is, for one RL and one EL ("easy") pair, a choice for the optimal stimulus was followed

690    by a positive outcome in 83% (vs. 17% negative outcome) of all trials (83% negative/17%

691    positive outcome for suboptimal stimulus); for the other RL and EL ("hard") pair a choice for

692    the optimal stimulus was followed by a positive outcome in 70% (vs. 30% negative outcome)

693    of all trials (and 70% negative/30% positive outcome) for the suboptimal stimulus. This

694    approach allowed us to disentangle whether acute stress primarily impacted domain-specific

695    (RL vs. EL) or general (easy vs. hard) reinforcement learning (the latter which also might

696    involve other cognitive skills that might be beneficial to performance and sensitive to change

697    under stress, such as working memory (Schoofs, Wolf, & Smeets, 2009). The task

698    contingencies described above were based on extensive pilot tests to identify a reinforcement

699    schedule that would enable us to detect stress-induced improvements *and* decreases in task

700    performance. We selected task contingencies based on pilot sessions involving a no-stress

701    control condition and chose a reinforcement schedule associated with non-ceiling/floor

702    performance on RL and EL trials.

703        Following the learning phase, participants completed a surprise test phase, similar to

704    previous work (D. Hernaus et al., 2019; D. Hernaus et al., 2018). This phase consisted of 64

705    trials in which participants were presented with the original four, as well as six novel,

706    stimulus combinations. Participants were asked to choose the stimulus with the highest

707    reward value or the lowest action cost - depending on a coin or thermometer image presented

708    in the middle of the screen - and received no choice feedback. This allowed us to assess

33

709      acquired choice tendencies, as well as generalizability of this information to novel situations.

710      The four original pairs were presented four times (total n=16) during which we only asked

711      participants to discriminate on the basis on the reward value (for RL) or action cost (for EL).

712      For novel stimulus combinations, we only presented stimuli that differed in reward

713      value/action cost if reward value discrimination/action cost discrimination was assessed (total

714      n=48: n=4 presentations for the 6 combinations).

715      For every participant, stimuli were randomly assigned to pairs, optimal/suboptimal

716      stimulus orientation was balanced (50% of all optimal stimulus presentations occurred on the

717      left-hand side) and misleading outcomes (e.g., negative outcomes for optimal stimuli) were

718      equally spaced out across the thirty presentations (and balanced for left/right side). Trial

719      presentation order was pseudo-randomized such that I) a given pair would never be presented

720      more than twice in a row and II) the gap between two presentations of a given pair was never

721      greater than four trials.

722      Prior to performing the actual task and prior before acute stress/no-stress control

723      procedures, participants received standard verbal instructions and completed a 16-trial

724      practice round of the learning phase. Participants were not informed about stimulus-outcome

725      contingencies; they were only advised to accrue as much money as possible and avoid

726      exerting unnecessary effort. A 60% accuracy performance threshold was used to confirm that

727      participants understood the general task procedure. The practice round was repeated if

728      participants failed to reach 60% accuracy. To prevent learning, we used deterministic

729      stimulus-outcome probabilities and different stimuli.

730

731      **Computational cost-benefit reinforcement learning model: model space**

732      In an attempt to uncover latent mechanisms by which acute stress affects reward

733      maximization and/or action cost minimization, we turned to cognitive computational

734      modelling. We employed a modified reinforcement learning framework based on Rescorla

735      and Wagner (Rescorla, 1972), and used in Skvortsova et al. (Skvortsova et al., 2017;

736      Skvortsova et al., 2014) to investigate whether acute stress impacted learning about

737      sensitivity to, and/or discounting of reward value and action cost. We first describe the model

738      space.

739           Various reinforcement learning models assume that choice preferences of an agent are

740      updated via the prediction error, i.e., the mismatch between outcome and expectation

741      (equation 1A, 1B) and the critical quantity that drives learning (Rescorla, 1972):

742

743      $$RPE_{(t)} = r_{(t)} - Q_{R(t)}(s, a) \quad (1A)$$

744      $$EPE_{(t)} = e_{(t)} - Q_{E(t)}(s, a) \quad (1B)$$

745

746           Here, $Q_{R(t)}(s, a)$ and $Q_{E(t)}(s, a)$ represent the expected reward value and action cost

747      (i.e., effort), where *s* reflects the given pair and *a* refers to the more abstract action of

748      selecting a stimulus (not to be confused with action selection), $r_{(t)}$ and $e_{(t)}$ represent the reward

749      and effort outcome for the chosen stimulus at trial *t*. $RPE_{(t)}$ and $EPE_{(t)}$, thus, represent the RPE

750      and EPE at trial t, respectively.

751           In order to allow for the possibility that humans do not calculate the prediction error

752      against the actual outcome but, rather, what the outcome "feels" like (Huys et al., 2013) , we

753      considered a scenario in which reward and effort outcomes are first multiplied by a free

754      parameter that captures the weight that reward and effort outcomes receive ("$W_R$" and $W_E$" in

755      equation 2A and 2B). As the value of these parameters approaches 1, rewards are

756      increasingly valued more positively, and effort more negatively. These parameters, therefore,

757      control the maximum size of the prediction error.

758

759 $$RPE_{(t)} = (r_{(t)} * W_R) - Q_{R(t)}(s, a) \quad (2A)$$

760 $$EPE_{(t)} = (e_{(t)} * W_E) - Q_{E(t)}(s, a) \quad (2B)$$

761

762       In various formulations of reinforcement learning, such as Q-learning (Watkins &

763 Dayan, 1992) and the actor-critic framework (Niv, 2009; Rescorla, 1972), the degree to

764 which prediction errors update choice preferences is represented by $\alpha$, the learning rate

765 (equation 3A), which determines how current prediction errors update choice preferences on

766 the subsequent trial. High values of $\alpha$ allow for rapid updating of choice preferences, while a

767 low $\alpha$ implies that choice preferences are updated at a slower pace and are thus co-

768 determined by outcomes further into the past.

769

770 $$Q_{R(t)}(s, a) = Q_{R(t-1)}(s, a) + \alpha_R * RPE_{(t-1)}(s, a) \quad (3A)$$

771 $$Q_{E(t)}(s, a) = Q_{E(t-1)}(s, a) + \alpha_E * RPE_{(t-1)}(s, a) \quad (3B)$$

772

773 Extensive evidence suggests that organisms use different learning systems for different types

774 of information, including reward value and action cost (Palminteri & Pessiglione, 2017;

775 Skvortsova et al., 2017; Skvortsova et al., 2014) (equation 3A/B). Thus, the use of separate

776 learning rates for RPEs and EPEs allows for asymmetrical learning about these types of

777 information.

778       While the learning rate controls the *speed* at which choice preferences are updated,

779 learning rate (nor reward/effort weight) alone does not explain how learned estimates of

780 reward value and action cost may compete at the decision stage (i.e., when participants

781 choose between two stimuli). Agents weight costs against benefits to calculate a subjective

782 decision value (Pessiglione et al., 2017; Skvortsova et al., 2017), which is used to guide

783 choices (equation 4).

36

784

785
$$Q_{(t)}(s, a) = Q_{R(t)}(s, a) - Q_{E(t)}(s, a) \quad (4)$$

786

787 In its simplest form, Q, the subjective decision value of a stimulus is represented by the

788 difference between the expected reward and action cost value at trial $t$ (equation 4)

789 (Skvortsova et al., 2014). However, this particular operationalization of subjective value does

790 not take into account the observation that humans tend to discount or prioritize certain types

791 of information in their decisions (Apps, Grima, Manohar, & Husain, 2015; Inzlicht et al.,

792 2018). We, therefore, allowed for variation in the calculation of subjective decision value via

793 action cost discounting (equation 5). While discounting rates can be linear or hyperbolic

794 (Hartmann, Hager, Tobler, & Kaiser, 2013), here we only considered linear discounting in

795 light of previous work using a similar task design (Skvortsova et al., 2017; Skvortsova et al.,

796 2014). As the value of γ approaches zero, action cost discounted increases leading the agent

797 to ignore action cost/only utilize reward value to make a decision.

798

799
$$Q_{(t)}(s, a) = Q_{R(t)}(s, a) - \gamma*Q_{E(t)}(s, a) \quad (5)$$

800

801 Once the subjective decision value has been computed, the degree to which

802 participants deterministically sample the optimal stimulus is captured by a softmax decision

803 function (equation 6).

804

805
$$pr(s, a) = \exp(Q_{(t)}(s, a)) / \text{sum}(\exp(\beta*Q_{(t)}(s))) \quad (6)$$

806

807 Here, pr is the probability of selecting an action, β is the inverse temperature parameter

808 that among others captures the balance between exploration and exploitation (Nassar &

37

809    Frank, 2016), $Q_{(t)}(s, a)$ is the net value of the chosen option and $Q_{(t)}(s)$ represents the net

810    values of both stimuli in the pair.

811        Within the above-described model space our predictions of acute stress effects on reward

812    maximization and action cost minimization could, thus, be explained by changes in

813    sensitivity to reward value and/or action cost ($W_R, W_E$), changes in how much weight RPEs

814    and EPEs are afforded (i.e., learning rates, $\alpha_R, \alpha_E$), and/or changes in the discounting of

815    reward value by action cost ($\gamma$). If acute stress leads to more random responses, such effects

816    should be captured by $\beta$.

817        Based on our predictions and the obtained pattern of results (most notably asymmetrical

818    RL/EL performance in the acute stress condition), we considered six candidate models that

819    could capture these various scenarios: I) a model with 2 distinct learning rates for reward and

820    effort ($\alpha R, \alpha E$) [2LR]; II) a model with 2 learning rates ($\alpha_R, \alpha_E$) and a discounting parameter

821    ($\gamma$) (2LR_ $\gamma$); III) a model with 2 learning rates ($\alpha_R, \alpha_E$), a reward weight ($W_R$) and an effort

822    weight parameter ($W_E$) (2LR_$W_R$_$W_E$), IV) a model with a *single* learning rate ($\alpha$), reward

823    weight ($W_R$), effort weight ($W_E$), and a discounting ($\gamma$) parameter (LR_ $W_R$__$W_E$_ $\gamma$); V) a

824    model with 2 learning rates ($\alpha_R, \alpha_E$), a reward weight ($W_R$), and a discounting ($\gamma$) parameter

825    (2LR_$W_R$_$\gamma$); VI) a model with 2 learning rates ($\alpha_R, \alpha_E$), a reward weight ($W_R$), effort weight

826    ($W_E$) and discounting ($\gamma$) parameter (2LR_ $W_R$_ $W_E$_ $\gamma$).

827        Lower/upper bounds for all parameters were set to [0,1] and all models contained a $\beta$

828    parameter. Consistent with previous work (Skvortsova et al., 2017; Skvortsova et al., 2014),

829    reward and action cost outcomes were set to [0,1 for no/yes reward] and [-1,0 for no/yes

830    effort avoidance], respectively.

**Pupillometry**

Fluctuations in pupil diameter were continuously measured using an SR-Research Eyelink 1000 Tower Mount infrared eye tracker while participants performed the reward maximization/action cost minimization reinforcement learning task (1000Hz sampling rate, except for three participants, whose data were obtained at 500Hz). Participants placed their head on an adjustable chin rest and against a forehead bar to minimize motion. Eye-tracker calibration was performed at the start of the paradigm, and subsequently every 10 min. Stimulus luminance was matched using the SHINE toolbox (Willenbockel et al., 2010) in MATLAB (v. 2014B; The MathWorks, Inc., Natick, Massachusetts, United States). Due to the COVID-19 pandemic, pupillometry data were not collected for the final eight participants. Three participants, moreover, failed the quality control for eye-tracking data (2 no-stress control/1 acute stress) leaving a final sample of 69 participants with eye-tracking data (34 no-stress control/35 acute stress).

Eye-tracking data were pre-processed using an open source pre-processing toolbox (Kret & Sjak-Shie, 2019) and in accordance with previous work (Jackson & Sirois, 2009). Blinks and other invalid samples, due to dilation speed, deviation from the trend line, and extreme values (Kret & Sjak-Shie, 2019) were removed, interpolated, smoothed (4Hz low-pass filter, fourth-order Butterworth filter) (Jackson & Sirois, 2009), z-scored and down-sampled to 50hz (i.e., 20ms). Bins with fewer than 80% valid samples were removed (Lawson et al., 2020). For analyses, we considered three epochs of interest: choice (-1500ms pre-choice - 15000ms post-choice), effort outcome (0-1000ms post-outcome), and reward outcome (0-2000ms post-outcome). We reduced the duration of the effort outcome epoch to 1000ms to minimize force exertion-related effects on pupil size (see below). Recent work has shown that expectation violations (prediction errors) are encoded by pupil size fluctuations within this timeframe (Lawson et al., 2020). Given that we observed large grip force-

856    associated effects on the pupillometry signal (see Figure Supplement 8 middle row, for a

857    comparison between effort and effort avoidance trials), we limited effort outcome analyses in

858    the main text to effort avoidance trials, although we also report analyses involving all effort

859    outcome trials in Figure Supplement 8.

860

861    **Statistical analyses**

862    Statistical analyses were conducted using R, version 3.6.2 (Team, 2020) and, where

863    applicable, results were visualised using Raincloud Plots (Allen, Poggiali, Whitaker,

864    Marshall, & Kievit, 2019). Acute stress measurements were analysed using mixed ANOVAs

865    involving Condition (between-factor condition: no-stress control, acute stress induction) and

866    Time (within-factor: 2 pre/post-MAST or 6 levels for sCORT).

867         For the reward maximisation/action cost minimisation reinforcement learning task, an

868    accuracy score was calculated by dividing the number of optimal stimulus choices by the

869    total trial amount ($n$=30 per pair). Mixed ANOVAs involving Condition, Trial Type (RL, EL)

870    and Difficulty (Easy, Hard pairs) were carried out. For analyses involving Time effects (i.e.,

871    repeated presentations of stimulus pairs), accuracy scores were averaged per bin of ten

872    presentations (presentation 1-10, 11-20, and 21-30). To better understand whether acute

873    stress effects on task performance were primarily driven by changes in sensitivity to positive

874    or negative outcomes, win-stay (repeating a choice following a positive outcome) and lose-

875    shift (choosing the other stimulus following a loss) rates were calculated for RL and EL trials

876    (Hanneke E. M. den Ouden et al., 2013). For RL trials, we calculated win-stay/lose-shift rates

877    using reward outcomes (yes/no reward); for EL trials we used effort outcomes (yes/no effort).

878    We refer to the 2-level factor representing win-stay/lose-shift rates as "Strategy". For surprise

879    test trials involving the original four pairs (n=4 presentations per pair), we investigated final

880    choice tendences using a one-sample t-test against chance level (0.5). Participants' ability to

881  discriminate stimuli based on reward value and action cost in novel stimulus arrangements

882  (n=48, 24 reward value and 24 action cost discrimination trials) were investigated using

883  mixed ANOVAs involving Condition and Trial Type.

884       Group differences in model parameters from the non-hierarchically fit model were

885  investigated using Condition-by learning rate ($\alpha_R$, $\alpha_E$) mixed ANOVAs and independent

886  samples t-tests. Given that we used separate priors for the two groups, we report the Bayesian

887  analogue of a t-test and mixed-ANOVA (Kruschke, 2014) - a more robust test of group

888  differences - for posterior parameters obtained from the hierarchically fit model (for

889  reference, we also report these analyses for the non-hierarchical data).

890       *Post hoc* (simple) main effect analyses for all ANOVAs were conducted using

891  independent sample (Condition), paired-samples (Time, Trial Type, Strategy), and one-

892  sample t-tests ($\neq 0$ or 0.5). Greenhouse–Geisser-corrected statistics were reported when

893  sphericity assumptions were violated. We report statistical significance as $p<0.05$ (two-

894  sided), but we note that most main and interaction effects involving Condition survived at a

895  more stringent threshold ($p<.01$), except for some strategy and surprise test phase effects,

896  which should be interpreted with caution. In case of statistically significant results,

897  generalized eta square (ges; $\eta^2_G$) was reported, with $\eta^2_G$ values of 0.02, 0.06, and 0.14

898  representing a small, medium, and large effect size, respectively (Lakens, 2013).

899       With respect to pupillometry, we conducted model-free and model-based analyses. In

900  model-free analyses, we investigated group differences in pupil size during the choice, effort

901  outcome, and reward outcome stage, for every bin of interest. To better understand how

902  putative activity of ascending neuromodulatory systems may drive stress-induced changes in

903  computational strategies that support reward maximisation/action cost minimisation learning,

904  we conducted model-based pupillometry analyses using computational parameters from the

905  winning (2LR_ $\gamma$) model (Lawson et al., 2020). First, we used linear regression to estimate

906  beta weights for the association between pupil size and computational estimates of task-

907  related behaviour for every participant, for every epoch, for every bin. For the choice phase,

908  we regressed trial-wise measures of pupil size against trial-wise estimates of the subjective

909  decision value (i.e., effort-discounted reward value) of the chosen stimulus. For the effort and

910  reward outcome phase, trial-wise EPEs and RPEs were the primary predictors of interest,

911  respectively. Trial number (1-120) and presented images/pair (RL_easy, RL_hard, EL_easy,

912  EL hard) served as additional predictors of interest for all models. Additional epoch-specific

913  variables of interest were included for the choice (optimal choice yes/no), effort (action cost

914  of chosen stimulus, effort avoidance yes/no), and reward (reward value of chosen stimulus,

915  reward yes/no, effort avoidance yes/no) outcome phase. Similar results were obtained when

916  repeating the analyses with more elaborate GLMs (e.g., the addition of yes/no most likely

917  outcome based on reward/effort outcome probabilities ["surprise"] and reward/action cost for

918  EL/RL trials). Secondly, in group-level GLMs, we compared the resulting beta weights I)

919  against zero (for the no-stress control/acute stress condition separately), to investigate when

920  the pupil encoded the computational process of interest, and II) between groups, to assess

921  stress-induced changes in  associations between pupil size and computational processes. To

922  control the false positive rate, we conducted permutation tests at the bin- and cluster-level

923  (2000 permutations, $\alpha_{permute}$=0.05). All correlations were performed using Spearman's $\rho$

924  correlations. Permutation tests were also conducted for correlation analyses involving acute

925  stress measures and pupil encoding of predictions errors.

926    **Figure 3 – Source Data 1**

927    **Source files for task performance data.**

928    This link contains all task performance data used for the analyses shown in Figure 3. Raw

929    data can be found under "task_performance".

930    https://osf.io/ydv2q/

931

932    **Figure 6 – Source Data 2**

933    **Source files for pupillometry data.**

934    This link contains pupillometry data used for the analyses shown in Figure 6. Raw data can

935    be found under "pupillometry".

936    https://osf.io/ydv2q/

937

938    **Acknowledgements**

939    We are indebted to Drs. Conny Quaedflieg, Edwin S. Dalmaijer. and Ross D. Markello for

940    advice on stress-induction procedures and paradigm development. We thank Dr. Michael

941    Frank for advice on hierarchical computational modelling procedures. We thank Katya Brat-

942    Matchett for involvement in participant recruitment and Truda Driessen for administrative

943    support.

944

945    **Competing interests**

946    D.H. has received financial compensation as a consultant for P1vital Products Ltd. These

947    activities were unrelated to the work presented in this manuscript. The authors declare no

948    competing interests.

949 **References**

950 Allen, M., Poggiali, D., Whitaker, K., Marshall, T. R., & Kievit, R. A. (2019). Raincloud

951       plots: a multi-platform tool for robust data visualization. *Wellcome Open Res, 4*, 63.

952       doi:10.12688/wellcomeopenres.15191.1

953 Apps, M. A. J., Grima, L. L., Manohar, S., & Husain, M. (2015). The role of cognitive effort

954       in subjective reward devaluation and risky decision-making. *Scientific reports, 5*(1),

955       16880. doi:10.1038/srep16880

956 Arnsten, A. F. (2015). Stress weakens prefrontal networks: molecular insults to higher

957       cognition. *Nat Neurosci, 18*(10), 1376-1385. doi:10.1038/nn.4087

958 Arnsten, A. F., & Goldman-Rakic, P. S. (1998). Noise stress impairs prefrontal cortical

959       cognitive function in monkeys: evidence for a hyperdopaminergic mechanism. *Arch*

960       *Gen Psychiatry, 55*(4), 362-368. doi:10.1001/archpsyc.55.4.362

961 Baik, J.-H. (2020). Stress and the dopaminergic reward system. *Experimental & Molecular*

962       *Medicine, 52*(12), 1879-1890. doi:10.1038/s12276-020-00532-4

963 Bailey, S. L., & Heitkemper, M. M. (2001). Circadian rhythmicity of cortisol and body

964       temperature: morningness-eveningness effects. *Chronobiology International, 18*(2),

965       249-261. doi:10.1081/CBI-100103189

966 Berghorst, L. H., Bogdan, R., Frank, M. J., & Pizzagalli, D. A. (2013). Acute stress

967       selectively reduces reward sensitivity. *Front Hum Neurosci, 7*, 133.

968       doi:10.3389/fnhum.2013.00133

969 Bryce, C. A., & Floresco, S. B. (2016). Perturbations in Effort-Related Decision-Making

970       Driven by Acute Stress and Corticotropin-Releasing Factor.

971       *Neuropsychopharmacology, 41*(8), 2147-2159. doi:10.1038/npp.2016.15

972     Byrne, K. A., Cornwall, A. C., & Worthy, D. A. (2019). Acute stress improves long-term

973          reward maximization in decision-making under uncertainty. *Brain and Cognition,*

974          *133*, 84-93. doi:https://doi.org/10.1016/j.bandc.2019.02.005

975     Cannon, W. B. (1915). *Bodily changes in pain, hunger, fear, and rage. CHAPTER XI: The*

976          *utility of the bodily changes in pain and great emotion*: D. Appleton and company.

977     Carvalheiro, J., Conceição, V. A., Mesquita, A., & Seara-Cardoso, A. (2020). Acute stress

978          impairs reward learning in men. *Brain Cogn, 147*, 105657.

979          doi:10.1016/j.bandc.2020.105657

980     Daw, N. D. (2011). Trial-by-trial data analysis using computational models. *Decision*

981          *making, affect, and learning: Attention and performance XXIII, 23*(1).

982     de Berker, A. O., Tirole, M., Rutledge, R. B., Cross, G. F., Dolan, R. J., & Bestmann, S.

983          (2016). Acute stress selectively impairs learning to act. *Scientific reports, 6*(1), 29816.

984          doi:10.1038/srep29816

985     de Kloet, E. R., Joëls, M., & Holsboer, F. (2005). Stress and the brain: from adaptation to

986          disease. *Nature Reviews Neuroscience, 6*(6), 463-475. doi:10.1038/nrn1683

987     Del Arco, A., Park, J., & Moghaddam, B. (2020). Unanticipated Stressful and Rewarding

988          Experiences Engage the Same Prefrontal Cortex and Ventral Tegmental Area

989          Neuronal Populations. *eneuro, 7*(3), ENEURO.0029-0020.2020.

990          doi:10.1523/ENEURO.0029-20.2020

991     den Ouden, H. E., Swart, J. C., Schmidt, K., Fekkes, D., Geurts, D. E., & Cools, R. (2015).

992          Acute serotonin depletion releases motivated inhibition of response vigour.

993          *Psychopharmacology (Berl), 232*(7), 1303-1312. doi:10.1007/s00213-014-3762-4

994     den Ouden, Hanneke E. M., Daw, Nathaniel D., Fernandez, G., Elshout, Joris A., Rijpkema,

995          M., Hoogman, M., . . . Cools, R. (2013). Dissociable Effects of Dopamine and

996   Serotonin on Reversal Learning. *Neuron, 80*(4), 1090-1100.

997   doi:https://doi.org/10.1016/j.neuron.2013.08.030

998 Dickerson, S. S., & Kemeny, M. E. (2004). Acute stressors and cortisol responses: a

999   theoretical integration and synthesis of laboratory research. *Psychol Bull, 130*(3), 355-

1000   391. doi:10.1037/0033-2909.130.3.355

1001 Evans, K. L., & Hampson, E. (2015). Sex-dependent effects on tasks assessing reinforcement

1002   learning and interference inhibition. *Frontiers in psychology, 6*, 1044-1044.

1003   doi:10.3389/fpsyg.2015.01044

1004 Frank, M. J., Seeberger, L. C., & Reilly, R. C. (2004). By Carrot or by Stick: Cognitive

1005   Reinforcement Learning in Parkinsonism. *science, 306*(5703), 1940.

1006   doi:10.1126/science.1102941

1007 Frey, A.-L., Frank, M. J., & McCabe, C. (2019). Social reinforcement learning as a predictor

1008   of real-life experiences in individuals with high and low depressive symptomatology.

1009   *Psychological Medicine*, 1-8. doi:10.1017/S0033291719003222

1010 Friedman, A., Homma, D., Bloem, B., Gibb, L. G., Amemori, K. I., Hu, D., . . . Graybiel, A.

1011   M. (2017). Chronic Stress Alters Striosome-Circuit Dynamics, Leading to Aberrant

1012   Decision-Making. *Cell, 171*(5), 1191-1205.e1128. doi:10.1016/j.cell.2017.10.017

1013 Guitart-Masip, M., Huys, Q. J. M., Fuentemilla, L., Dayan, P., Duzel, E., & Dolan, R. J.

1014   (2012). Go and no-go learning in reward and punishment: interactions between affect

1015   and effect. *Neuroimage, 62*(1), 154-166. doi:10.1016/j.neuroimage.2012.04.024

1016 Hartmann, M. N., Hager, O. M., Tobler, P. N., & Kaiser, S. (2013). Parabolic discounting of

1017   monetary rewards by physical effort. *Behavioural Processes, 100*, 192-196.

1018   doi:https://doi.org/10.1016/j.beproc.2013.09.014

1019      Hauser, T. U., Eldar, E., & Dolan, R. J. (2017). Separate mesocortical and mesolimbic

1020          pathways encode effort and reward learning signals. *Proceedings of the National*

1021          *Academy of Sciences, 114*(35), E7395. doi:10.1073/pnas.1705643114

1022      Hermans, E. J., Henckens, M. J. A. G., Joëls, M., & Fernández, G. (2014). Dynamic

1023          adaptation of large-scale brain networks in response to acute stressors. *Trends in*

1024          *neurosciences, 37*(6), 304-314. doi:https://doi.org/10.1016/j.tins.2014.03.006

1025      Hermans, E. J., van Marle, H. J. F., Ossewaarde, L., Henckens, M. J. A. G., Qin, S., van

1026          Kesteren, M. T. R., . . . Fernández, G. (2011). Stress-Related Noradrenergic Activity

1027          Prompts Large-Scale Neural Network Reconfiguration. *science, 334*(6059), 1151.

1028          doi:10.1126/science.1209603

1029      Hernandez Lallement, J., Kuss, K., Trautner, P., Weber, B., Falk, A., & Fliessbach, K.

1030          (2014). Effort increases sensitivity to reward and loss magnitude in the human brain.

1031          *Social cognitive and affective neuroscience, 9*(3), 342-349. doi:10.1093/scan/nss147

1032      Hernaus, D., Frank, M. J., Brown, E. C., Brown, J. K., Gold, J. M., & Waltz, J. A. (2019).

1033          Impaired Expected Value Computations in Schizophrenia Are Associated With a

1034          Reduced Ability to Integrate Reward Probability and Magnitude of Recent Outcomes.

1035          *Biol Psychiatry Cogn Neurosci Neuroimaging, 4*(3), 280-290.

1036          doi:10.1016/j.bpsc.2018.11.011

1037      Hernaus, D., Gold, J. M., Waltz, J. A., & Frank, M. J. (2018). Impaired Expected Value

1038          Computations Coupled With Overreliance on Stimulus-Response Learning in

1039          Schizophrenia. *Biol Psychiatry Cogn Neurosci Neuroimaging, 3*(11), 916-926.

1040          doi:10.1016/j.bpsc.2018.03.014

1041      Hernaus, D., Quaedflieg, C. W., Offermann, J. S., Casales Santa, M. M., & van Amelsvoort,

1042          T. (2018). Neuroendocrine stress responses predict catecholamine-dependent working

1043        memory-related dorsolateral prefrontal cortex activity. *Social cognitive and affective*

1044        *neuroscience, 13*(1), 114-123.

1045    Huys, Q. J., Pizzagalli, D. A., Bogdan, R., & Dayan, P. (2013). Mapping anhedonia onto

1046        reinforcement learning: a behavioural meta-analysis. *Biology of mood & anxiety*

1047        *disorders, 3*(1), 12-12. doi:10.1186/2045-5380-3-12

1048    Inzlicht, M., Shenhav, A., & Olivola, C. Y. (2018). The Effort Paradox: Effort Is Both Costly

1049        and Valued. *Trends in cognitive sciences, 22*(4), 337-349.

1050        doi:10.1016/j.tics.2018.01.007

1051    Jackson, I., & Sirois, S. (2009). Infant cognition: going full factorial with pupil dilation. *Dev*

1052        *Sci, 12*(4), 670-679. doi:10.1111/j.1467-7687.2008.00805.x

1053    Joëls, M., Pu, Z., Wiegert, O., Oitzl, M. S., & Krugers, H. J. (2006). Learning under stress:

1054        how does it work? *Trends in cognitive sciences, 10*(4), 152-158.

1055        doi:https://doi.org/10.1016/j.tics.2006.02.002

1056    Joshi, S., Li, Y., Kalwani, R. M., & Gold, J. I. (2016). Relationships between Pupil Diameter

1057        and Neuronal Activity in the Locus Coeruleus, Colliculi, and Cingulate Cortex.

1058        *Neuron, 89*(1), 221-234. doi:10.1016/j.neuron.2015.11.028

1059    Klein-Flügge, M. C., Kennerley, S. W., Friston, K., & Bestmann, S. (2016). Neural

1060        signatures of value comparison in human cingulate cortex during decisions requiring

1061        an effort-reward trade-off. *Journal of Neuroscience, 36*(39), 10002-10015.

1062    Koh, D., Ng, V., & Naing, L. (2014). Alpha Amylase as a Salivary Biomarker of Acute

1063        Stress of Venepuncture from Periodic Medical Examinations. *Frontiers in Public*

1064        *Health, 2*(121). doi:10.3389/fpubh.2014.00121

1065    Kret, M. E., & Sjak-Shie, E. E. (2019). Preprocessing pupil size data: Guidelines and code.

1066        *Behav Res Methods, 51*(3), 1336-1342. doi:10.3758/s13428-018-1075-y

1067    Kruschke, J. (2014). Doing Bayesian data analysis: A tutorial with R, JAGS, and Stan.

1068    Lakens, D. (2013). Calculating and reporting effect sizes to facilitate cumulative science: a

1069         practical primer for t-tests and ANOVAs. *Frontiers in psychology, 4*, 863-863.

1070         doi:10.3389/fpsyg.2013.00863

1071    Lawson, R. P., Bisby, J., Nord, C. L., Burgess, N., & Rees, G. (2020). The Computational,

1072         Pharmacological, and Physiological Determinants of Sensory Learning under

1073         Uncertainty. *Current Biology*. doi:https://doi.org/10.1016/j.cub.2020.10.043

1074    Le Heron, C., Plant, O., Manohar, S., Ang, Y. S., Jackson, M., Lennox, G., . . . Husain, M.

1075         (2018). Distinct effects of apathy and dopamine on effort-based decision-making in

1076         Parkinson's disease. *Brain, 141*(5), 1455-1469. doi:10.1093/brain/awy110

1077    Lefebvre, G., Lebreton, M., Meyniel, F., Bourgeois-Gironde, S., & Palminteri, S. (2017).

1078         Behavioural and neural characterization of optimistic reinforcement learning. *Nature

1079         Human Behaviour, 1*(4), 0067. doi:10.1038/s41562-017-0067

1080    Lighthall, N. R., Gorlick, M. A., Schoeke, A., Frank, M. J., & Mather, M. (2013). Stress

1081         modulates reinforcement learning in younger and older adults. *Psychology and aging,

1082         28*(1), 35-46. doi:10.1037/a0029823

1083    Meyniel, F., Goodwin, G. M., Deakin, J. W., Klinge, C., MacFadyen, C., Milligan, H., . . .

1084         Gaillard, R. (2016). A specific role for serotonin in overcoming effort cost. *Elife, 5*,

1085         e17282. doi:10.7554/eLife.17282

1086    Morris, B. H., & Rottenberg, J. (2015). Heightened reward learning under stress in

1087         generalized anxiety disorder: A predictor of depression resistance? *Journal of

1088         abnormal psychology, 124*(1), 115.

1089    Myung, J. I., Tang, Y., & Pitt, M. A. (2009). Evaluation and comparison of computational

1090         models. *Methods Enzymol, 454*, 287-304. doi:10.1016/s0076-6879(08)03811-1

1091    Nassar, M. R., & Frank, M. J. (2016). Taming the beast: extracting generalizable knowledge

1092         from computational models of cognition. *Current opinion in behavioral sciences, 11*,

1093         49-54. doi:10.1016/j.cobeha.2016.04.003

1094    Nater, U. M., Rohleder, N., Gaab, J., Berger, S., Jud, A., Kirschbaum, C., & Ehlert, U.

1095         (2005). Human salivary alpha-amylase reactivity in a psychosocial stress paradigm.

1096         *Int J Psychophysiol, 55*(3), 333-342. doi:10.1016/j.ijpsycho.2004.09.009

1097    Niv, Y. (2009). Reinforcement learning in the brain. *Journal of Mathematical Psychology,*

1098         *53*(3), 139-154. doi:https://doi.org/10.1016/j.jmp.2008.12.005

1099    Otto, A. R., Raio, C. M., Chiang, A., Phelps, E. A., & Daw, N. D. (2013). Working-memory

1100         capacity protects model-based learning from stress. *Proceedings of the National*

1101         *Academy of Sciences, 110*(52), 20941. doi:10.1073/pnas.1312011110

1102    Palminteri, S., & Pessiglione, M. (2017). Chapter 23 - Opponent Brain Systems for Reward

1103         and Punishment Learning: Causal Evidence From Drug and Lesion Studies in

1104         Humans. In J.-C. Dreher & L. Tremblay (Eds.), *Decision Neuroscience* (pp. 291-303).

1105         San Diego: Academic Press.

1106    Peirce, J., Gray, J. R., Simpson, S., MacAskill, M., Höchenberger, R., Sogo, H., . . . Lindeløv,

1107         J. K. (2019). PsychoPy2: Experiments in behavior made easy. *Behavior Research*

1108         *Methods, 51*(1), 195-203. doi:10.3758/s13428-018-01193-y

1109    Pessiglione, M., Vinckier, F., Bouret, S., Daunizeau, J., & Le Bouc, R. (2017). Why not try

1110         harder? Computational approach to motivation deficits in neuro-psychiatric diseases.

1111         *Brain, 141*(3), 629-650.

1112    Petzold, A., Plessow, F., Goschke, T., & Kirschbaum, C. (2010). Stress reduces use of

1113         negative feedback in a feedback-based learning task. *Behav Neurosci, 124*(2), 248-

1114         255. doi:10.1037/a0018930

1115    Pool, E., Brosch, T., Delplanque, S., & Sander, D. (2015). Stress increases cue-triggered

1116         "wanting" for sweet reward in humans. *Journal of experimental psychology. Animal*

1117         *learning and cognition, 41 2*, 128-136.

1118    Pruessner, J. C., Kirschbaum, C., Meinlschmid, G., & Hellhammer, D. H. (2003). Two

1119         formulas for computation of the area under the curve represent measures of total

1120         hormone concentration versus time-dependent change. *Psychoneuroendocrinology,*

1121         *28*(7), 916-931. doi:10.1016/s0306-4530(02)00108-7

1122    Qin, S., Hermans, E. J., van Marle, H. J. F., Luo, J., & Fernández, G. (2009). Acute

1123         Psychological Stress Reduces Working Memory-Related Activity in the Dorsolateral

1124         Prefrontal Cortex. *Biological psychiatry, 66*(1), 25-32.

1125         doi:https://doi.org/10.1016/j.biopsych.2009.03.006

1126    Raio, C. M., Hartley, C. A., Orederu, T. A., Li, J., & Phelps, E. A. (2017). Stress attenuates

1127         the flexible updating of aversive value. *Proceedings of the National Academy of*

1128         *Sciences of the United States of America, 114*(42), 11241-11246.

1129         doi:10.1073/pnas.1702565114

1130    Rescorla, R. A. (1972). A theory of Pavlovian conditioning: Variations in the effectiveness of

1131         reinforcement and nonreinforcement. *Current research and theory*, 64-99.

1132    Rigoux, L., Stephan, K. E., Friston, K. J., & Daunizeau, J. (2014). Bayesian model selection

1133         for group studies - revisited. *Neuroimage, 84*, 971-985.

1134         doi:10.1016/j.neuroimage.2013.08.065

1135    Russell, G., & Lightman, S. (2019). The human stress response. *Nat Rev Endocrinol, 15*(9),

1136         525-534. doi:10.1038/s41574-019-0228-0

1137    Schmidt, L., Lebreton, M., Cléry-Melin, M. L., Daunizeau, J., & Pessiglione, M. (2012).

1138         Neural mechanisms underlying motivation of mental versus physical effort. *PLoS*

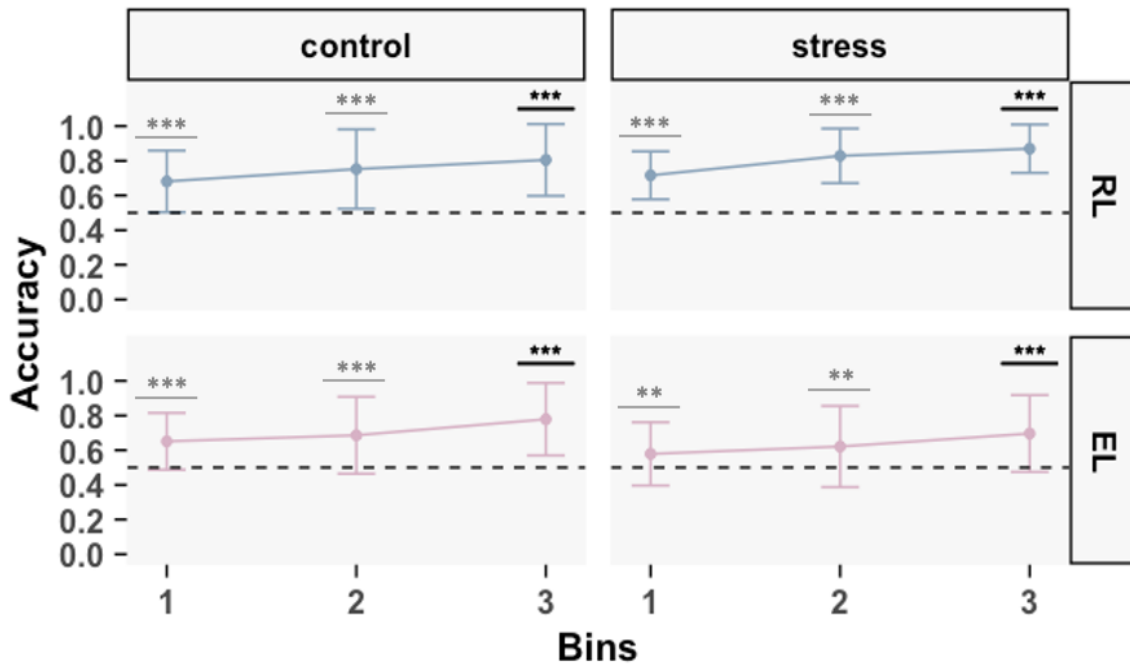1139         *Biol, 10*(2), e1001266. doi:10.1371/journal.pbio.1001266

1140    Schoofs, D., Wolf, O. T., & Smeets, T. (2009). Cold pressor stress impairs performance on

1141        working memory tasks requiring executive functions in healthy young men. *Behav*

1142        *Neurosci, 123*(5), 1066-1075. doi:10.1037/a0016980

1143    Schubert, C., Lambertz, M., Nelesen, R. A., Bardwell, W., Choi, J. B., & Dimsdale, J. E.

1144        (2009). Effects of stress on heart rate complexity--a comparison between short-term

1145        and chronic stress. *Biological psychology, 80*(3), 325-332.

1146        doi:10.1016/j.biopsycho.2008.11.005

1147    Schultz, W., Dayan, P., & Montague, P. R. (1997). A neural substrate of prediction and

1148        reward. *science, 275*(5306), 1593-1599. doi:10.1126/science.275.5306.1593

1149    Schwabe, L., Dickinson, A., & Wolf, O. T. (2011). Stress, habits, and drug addiction: a

1150        psychoneuroendocrinological perspective. *Exp Clin Psychopharmacol, 19*(1), 53-63.

1151        doi:10.1037/a0022212

1152    Schwabe, L., & Wolf, O. T. (2011). Stress-induced modulation of instrumental behavior:

1153        from goal-directed to habitual control of action. *Behavioural brain research, 219*(2),

1154        321-328.

1155    Shafiei, N., Gray, M., Viau, V., & Floresco, S. B. (2012). Acute stress induces selective

1156        alterations in cost/benefit decision-making. *Neuropsychopharmacology, 37*(10), 2194-

1157        2209. doi:10.1038/npp.2012.69

1158    Skvortsova, V., Degos, B., Welter, M.-L., Vidailhet, M., & Pessiglione, M. (2017). A

1159        selective role for dopamine in learning to maximize reward but not to minimize effort:

1160        evidence from patients with Parkinson's disease. *Journal of Neuroscience, 37*(25),

1161        6087-6097.

1162    Skvortsova, V., Palminteri, S., & Pessiglione, M. (2014). Learning to minimize efforts versus

1163        maximizing rewards: computational principles and neural correlates. *Journal of*

1164        *Neuroscience, 34*(47), 15621-15630.

1165    Smeets, T., Cornelisse, S., Quaedflieg, C. W., Meyer, T., Jelicic, M., & Merckelbach, H.

1166        (2012). Introducing the Maastricht Acute Stress Test (MAST): a quick and non-

1167        invasive approach to elicit robust autonomic and glucocorticoid stress responses.

1168        *Psychoneuroendocrinology, 37*(12), 1998-2008. doi:10.1016/j.psyneuen.2012.04.012

1169    Steinberg, E. E., Keiflin, R., Boivin, J. R., Witten, I. B., Deisseroth, K., & Janak, P. H.

1170        (2013). A causal link between prediction errors, dopamine neurons and learning.

1171        *Nature neuroscience, 16*(7), 966-973. doi:10.1038/nn.3413

1172    Stelly, C. E., Tritley, S. C., Rafati, Y., & Wanat, M. J. (2020). Acute Stress Enhances

1173        Associative Learning via Dopamine Signaling in the Ventral Lateral Striatum. *The

1174        Journal of Neuroscience, 40*(22), 4391. doi:10.1523/JNEUROSCI.3003-19.2020

1175    Team, R. C. (2020). R: A language and environment for statistical

1176        computing. *R Foundation for Statistical Computing*. Retrieved from https://www.R-

1177        project.org/

1178    Timmers, I., Kaas, A. L., Quaedflieg, C., Biggs, E. E., Smeets, T., & de Jong, J. R. (2018).

1179        Fear of pain and cortisol reactivity predict the strength of stress-induced hypoalgesia.

1180        *Eur J Pain, 22*(7), 1291-1303. doi:10.1002/ejp.1217

1181    Valton, V., Wise, T., & Robinson, O. J. (2020). Recommendations for Bayesian hierarchical

1182        model specifications for case-control studies in mental health. *arXiv preprint

1183        arXiv:2011.01725*.

1184    Varazzani, C., San-Galli, A., Gilardeau, S., & Bouret, S. (2015). Noradrenaline and

1185        dopamine neurons in the reward/effort trade-off: a direct electrophysiological

1186        comparison in behaving monkeys. *J Neurosci, 35*(20), 7866-7877.

1187        doi:10.1523/jneurosci.0454-15.2015

1188    Watkins, C. J. C. H., & Dayan, P. (1992). Q-learning. *Machine Learning, 8*(3), 279-292.

1189        doi:10.1007/BF00992698

1190    Watson, D., Clark, L. A., & Tellegen, A. (1988). Development and validation of brief

1191        measures of positive and negative affect: the PANAS scales. *J Pers Soc Psychol,*

1192        *54*(6), 1063-1070. doi:10.1037//0022-3514.54.6.1063

1193    Willenbockel, V., Sadr, J., Fiset, D., Horne, G. O., Gosselin, F., & Tanaka, J. W. (2010).

1194        Controlling low-level image properties: The SHINE toolbox. *Behavior Research*

1195        *Methods, 42*(3), 671-684. doi:10.3758/BRM.42.3.671

1196    Wilson, R. C., & Collins, A. G. (2019). Ten simple rules for the computational modeling of

1197        behavioral data. *Elife, 8*, e49547.

1198    Wright, B. J., O'Brien, S., Hazi, A., & Kent, S. (2014). Increased systolic blood pressure

1199        reactivity to acute stress is related with better self-reported health. *Scientific reports,*

1200        *4*, 6882-6882. doi:10.1038/srep06882

1201    Yohn, S. E., Errante, E. E., Rosenbloom-Snow, A., Somerville, M., Rowland, M., Tokarski,

1202        K., . . . Salamone, J. D. (2016). Blockade of uptake for dopamine, but not

1203        norepinephrine or 5-HT, increases selection of high effort instrumental activity:

1204        Implications for treatment of effort-related motivational symptoms in

1205        psychopathology. *Neuropharmacology, 109*, 270-280.

1206        doi:10.1016/j.neuropharm.2016.06.018

1207

1208 **Supplementary Figures**

1209    **Figure Supplement 1**



1210

1211 **Evidence of reward and action cost reinforcement learning.**

1212 Optimal stimulus choices ("accuracy") on reward learning (RL) and effort learning (EL)
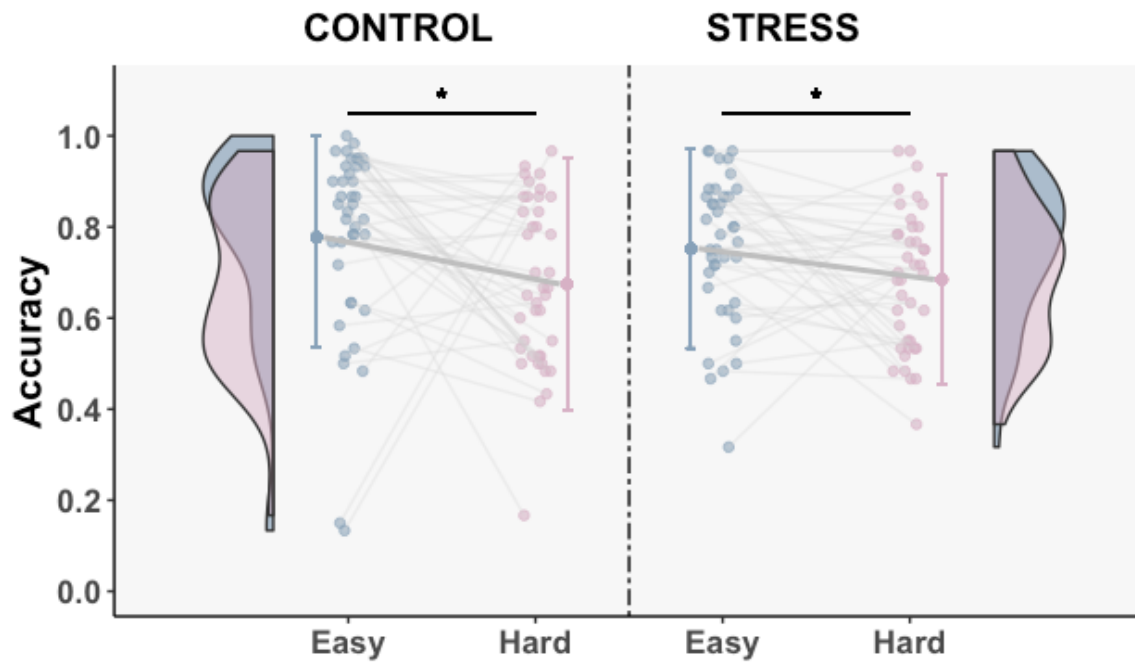
1213 (rows) trials for both conditions (columns). Trials were binned into groups of 10

1214 presentations. Participants performed significantly better than chance level in all bins. Means

1215 ± SD. Significant differences are denoted by asterisks (*: $p < 0.05$, **: $p < 0.01$, ***: $p <$
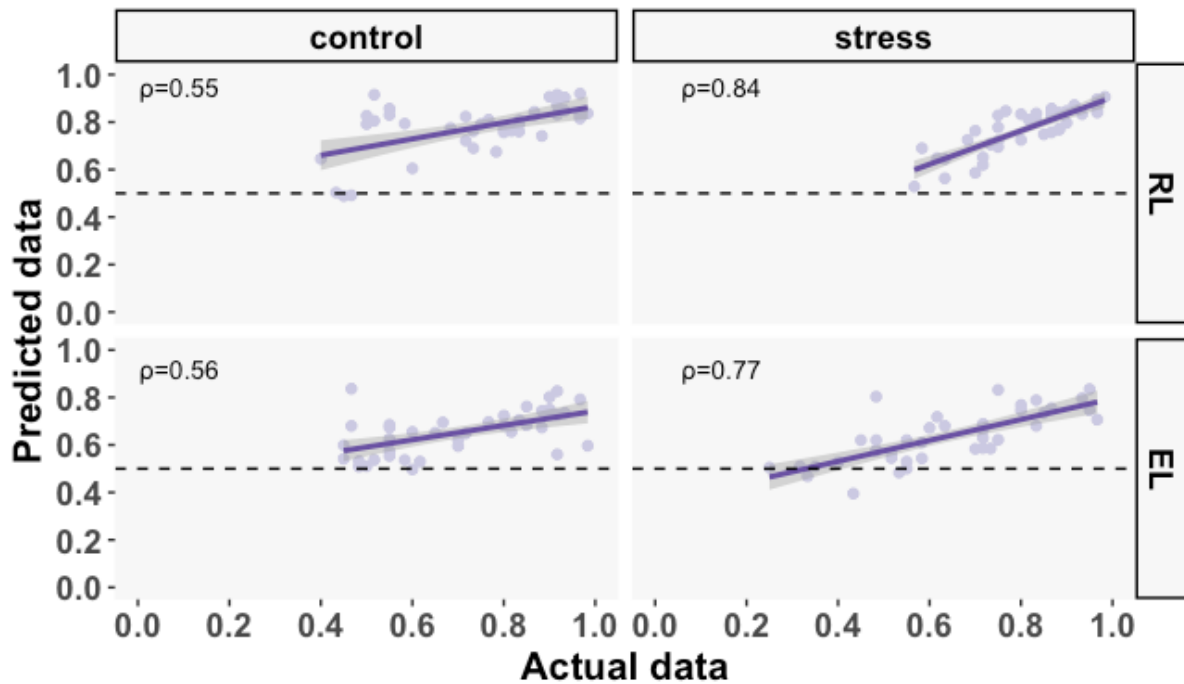
1216 0.001).

1217    **Figure Supplement 2**



1218

1219    **Acute stress does not affect difficulty learning.**

1220    Easy and hard pairs collapsed across RL/EL trials depicted for each condition separately.

1221    While all participants sampled the optimal choices more frequently for Easy vs. Hard pairs,

1222    no significant Condition-by-Difficulty interaction or between-group differences were

1223    observed. Means ± SD, individual data points, distribution and density of the data are

1224    displayed. Significant differences are denoted by asterisks (*: $p < 0.05$, **: $p < 0.01$, ***: $p <$

1225    0.001).

1226    **Figure Supplement 3**



1227

1228    **Surprise test phase performance.**

1229    The acute stress group performed better on reward than action cost discrimination trials.

1230    Means ± SD, individual data points, distribution and density of the data are displayed.

1231    Significant differences are depicted with asterisks (*: $p < 0.05$, **: $p < 0.01$, ***: $p < 0.001$).
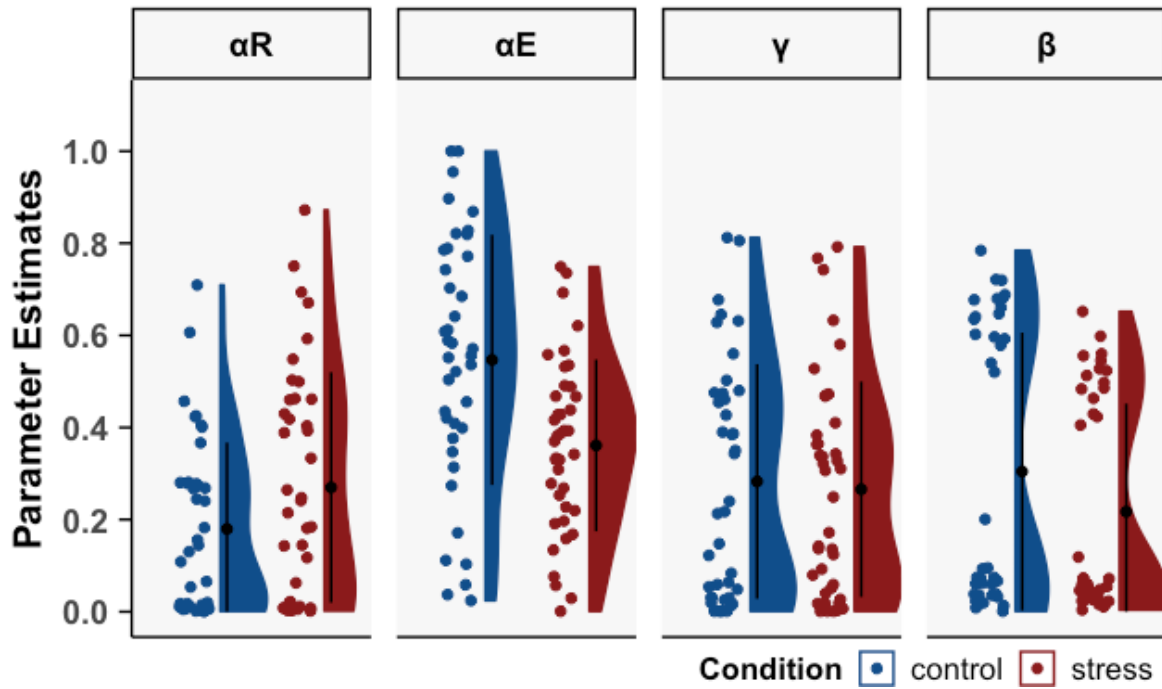
1232  **Figure Supplement 4**



1233

1234  **Correlations between empirical and simulated 2LR_γ choices.**

1235  Actual and *post hoc* simulated choices for RL and EL (rows) were highly correlated both for

1236  no-stress control and acute stress subjects (columns). Simulations were averaged across 10

1237  repetitions per subject. Solid and shaded lines represent mean $\pm$ CI$_{95\%}$. Dots represent

1238  individual data points. Horizontal dashed lines indicate chance level (0.5).

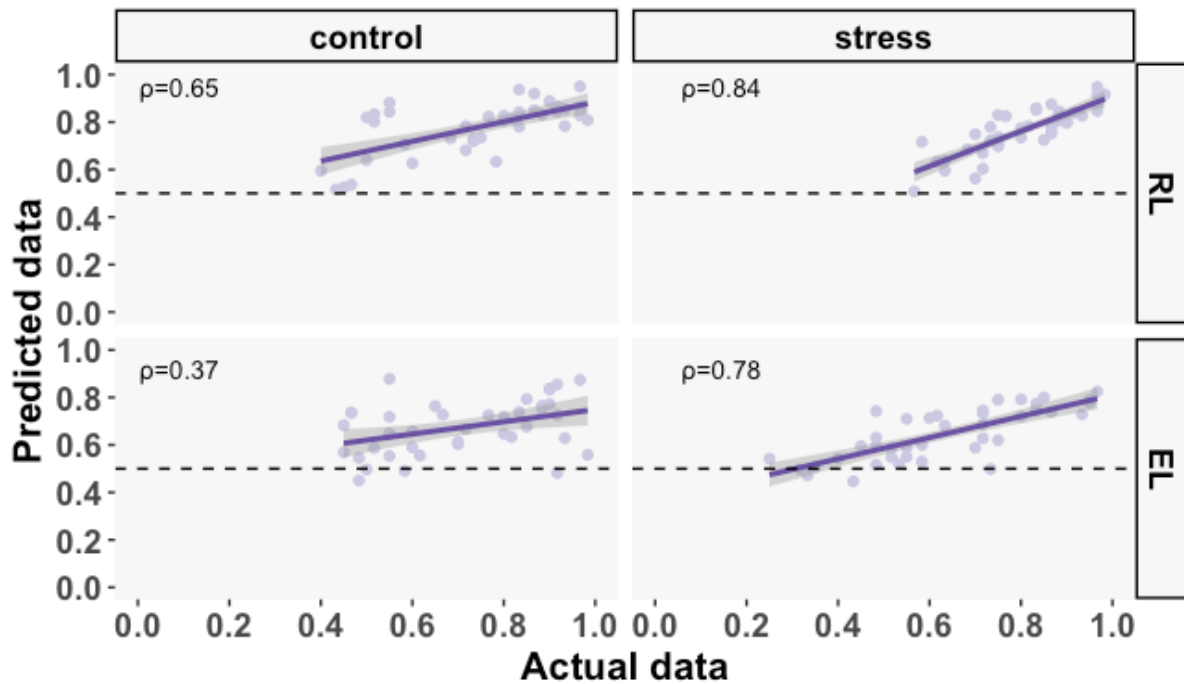1239    **Figure Supplement 5**



1240

1241    **Parameter estimates after Bayesian hierarchical model fitting.**

1242    Hierarchical model fitting reproduced the overall pattern of parameter estimates (Figure 5 for

1243    comparison) .

1244    **Figure Supplement *6***



1245

1246    **Correlations between empirical and simulated 2LR_γ choices after Bayesian**

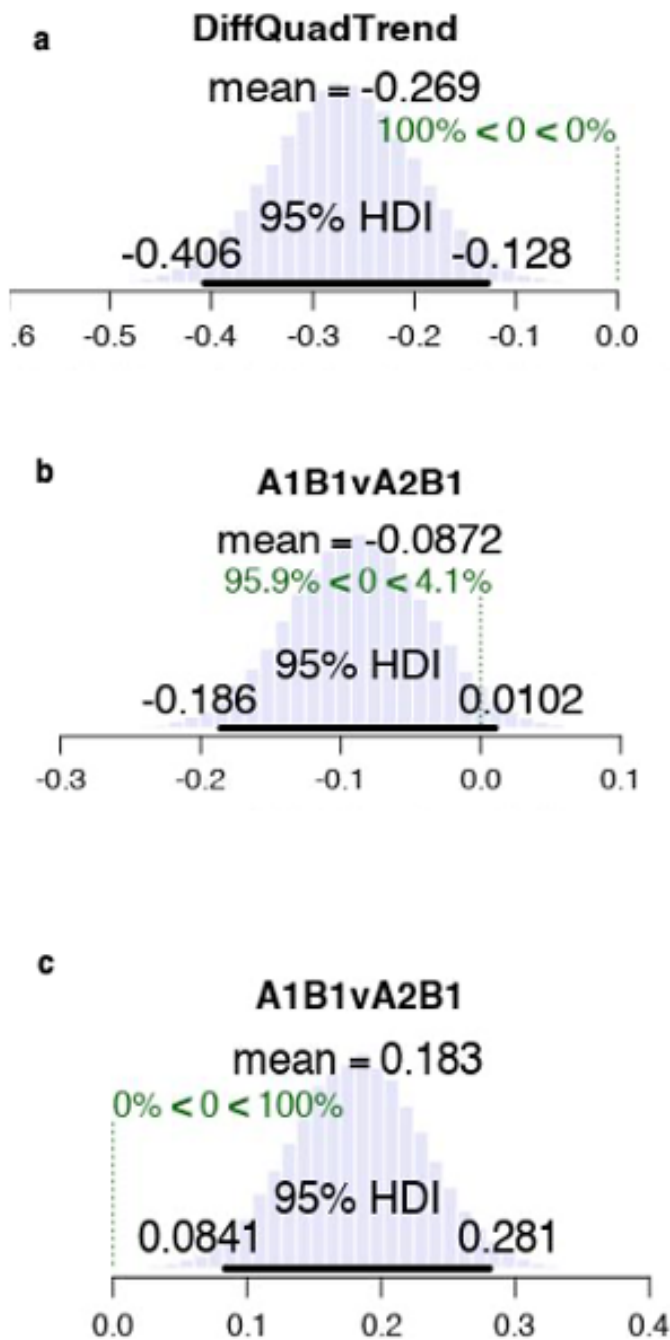1247    **hierarchical model fitting.**

1248    Correlations between actual and *post hoc* simulated choices for RL and EL (rows) for no-

1249    stress control and acute stress subjects (columns). Simulations were averaged across 10

1250    repetitions per subject. Solid and shaded lines represent mean $\pm$ CI$_{95\%}$. Dots represent

1251    individual data points. Horizontal dashed lines indicate chance level (0.5).

1252    **Figure Supplement 7**



1253

1254    **Bayesian estimation analysis to evaluate group differences in posterior parameter**
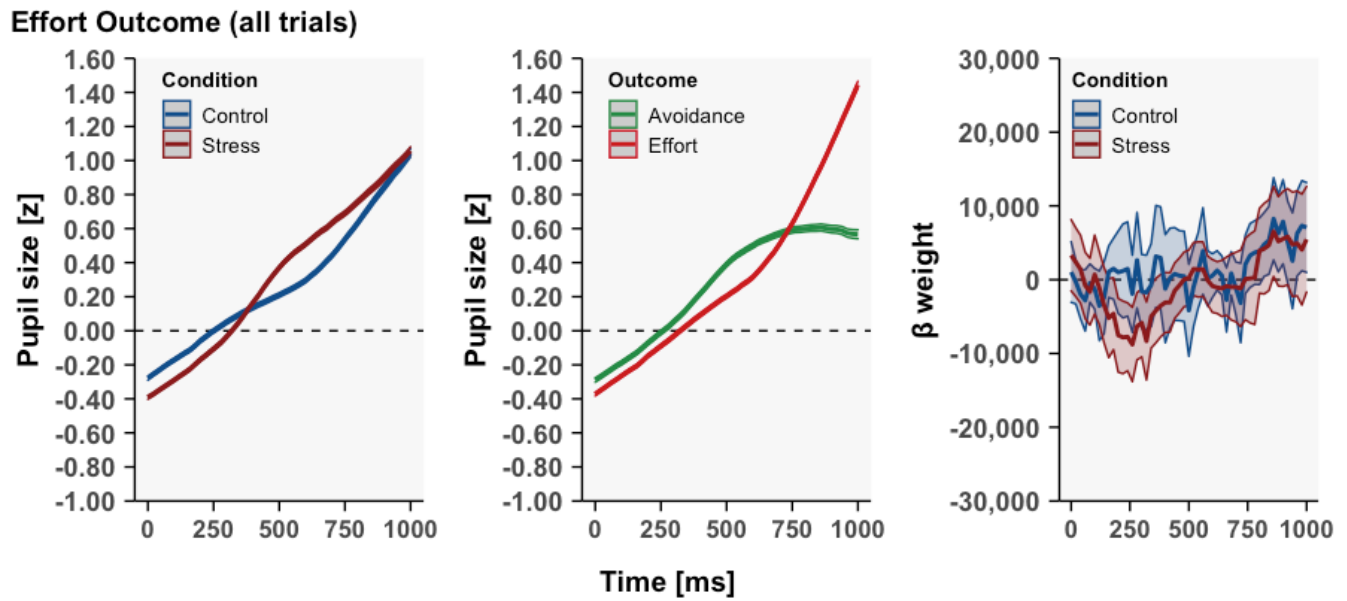
1255    **distributions**

1256    Panel A. Bayesian estimation (mixed-ANOVA) using posterior parameters (following

1257    hierarchical fitting) revealed evidence for a credible Condition-by-Learning Rate interaction.

1258    The observed mean difference from zero that falls outside the 95% HDI suggests that the

1259    difference between $\alpha_E$ and $\alpha_R$ was greater in no-stress controls compared to acute stress

1260    subjects. Panel B. Both groups did *not* differ in the magnitude of $\alpha_R$, as indicated by a 95%

1261    HDI that included 0. Panel C. Acute stress compared to no-stress control subjects exhibited a

1262    lower value of $\alpha_E$, as indicated by a 95% HDI that falls well above zero.

1263    **Figure Supplement 8**



1264

1265    **Pupillometry analyses using all effort outcome trials.**

1266    **Left:** Model-free analyses of pupil size using all effort outcome trials. **Middle**: Pupil size

1267    differences during effort/effort avoidance outcomes in the entire sample; force exertion was

1268    associated with large effects on pupil size and these trials were therefore excluded from

1269    analysis. **Right**: Model-based action cost prediction error analyses using all effort outcome

1270    trials.