

# 1 Reward prediction errors induce risk-seeking

2 Manuscript

3 Moritz Moeller\*1, Jan Grohn\*2, Sanjay Manohar\*\*12+, Rafal Bogacz\*\*1.

4 \*: equal contributions

5 \*\*: equal contributions

6 1: Nuffield Department of Clinical Neurosciences, University of Oxford

7 2: Department of Experimental Psychology, University of Oxford

8 +: corresponding author ([sanjay.manohar@psy.ox.ac.uk](mailto:sanjay.manohar@psy.ox.ac.uk))

## 9 Abstract

10 Reinforcement learning theories propose that humans choose based on the estimated values of  
11 available options, and that they learn from rewards by reducing the difference between the experienced  
12 and expected value. In the brain, such prediction errors are broadcasted by dopamine. However, choices  
13 are not only influenced by expected value, but also by risk. Like reinforcement learning, risk preferences  
14 are modulated by dopamine: enhanced dopamine levels induce risk-seeking. Learning and risk  
15 preferences have so far been studied independently, even though it is commonly assumed that they are  
16 (partly) regulated by the same neurotransmitter. Here, we use a novel learning task to look for  
17 prediction-error induced risk-seeking in human behavior and pupil responses. We find that prediction  
18 errors are positively correlated with risk-preferences in imminent choices. Physiologically, this effect is  
19 indexed by pupil dilation: only participants whose pupil response indicates that they experienced the  
20 prediction error also show the behavioral effect.

## 21 Introduction

22 Reward-guided learning in humans and animals can often be modelled simply as reducing the difference  
23 between the obtained and expected reward—a reward prediction error. This well-established  
24 behavioral phenomenon [Rescorla, 1972] has been linked to the neurotransmitter dopamine [Schultz,  
25 1997]. It has been shown that bursts of dopaminergic activity broadcast prediction errors to brain areas  
26 that are relevant for reward learning, such as the striatum, the amygdala, and the prefrontal cortex.

27 Another behavioral phenomenon that has been well studied is the effect of uncertainty and risk on  
28 decision making [Kahneman, 2013]. Here again, a different line of research has established an  
29 association between dopamine and risk-taking: dopamine-enhancing medication has been shown to  
30 increase risk-seeking in rats [St Onge, 2009], and drive excessive gambling when used to treat  
31 Parkinson's disease [Voon, 2006] [Gallagher, 2007] [Weintraub, 2010]. More recently, it has been  
32 demonstrated that phasic responses in dopaminergic brain areas modulate moment-by-moment risk-  
33 preference in humans: the tendency to take risks correlated positively with the magnitude of task-  
34 related dopamine responses [Chew, 2019]. A family of mechanistic theories of the basal ganglia network  
35 provides an explanation for these risk effects [Mikhael, 2016] [Moeller, 2019]. According to these  
36 models, positive and negative outcomes of actions are encoded separately in direct and indirect  
37 pathways of the basal ganglia. The balance between those pathways is controlled by the dopamine  
38 level. An increased dopamine level promotes the direct pathway and hence puts emphasis on potential  
39 gains, thus rendering risky options more attractive.

40 In summary, dopamine bursts are related to distinct behavioral phenomena—learning and risk-taking—  
41 by way of 1) acting as reward prediction errors, affecting synaptic weights during reinforcement  
42 learning, and 2) inducing risk-seeking behavior directly. There is no obvious a priori reason for those  
43 functions to be bundled together; in fact, one would perhaps expect them to work independently, and

44 their conflation might lead to interactions, unless some separation mechanism exists. There have been  
45 different suggestions for such separation mechanisms: it has been proposed that the tonic level of  
46 dopamine might modulate behavior directly, while phasic dopamine bursts provide the prediction errors  
47 necessary for reward learning [Niv, 2007]. Alternatively, cholinergic interneurons might flag dopamine  
48 activity that is to be interpreted as prediction errors by striatal neurons [Berke, 2018]. However, it has  
49 also been suggested that the architecture of the basal ganglia is well set-up to both learn from reward-  
50 prediction errors and use them to regulate risk-preferences [Mikhael, 2016] [Moeller, 2019].

51 Curiously, even though the multi-functionality of dopamine has been noted and separation mechanisms  
52 have been proposed, interference between the different functions has never been investigated  
53 experimentally. Here, we investigate this: if dopamine indeed provides both prediction errors for  
54 learning and modulates risk preferences, do these two processes interfere with each other, or are they  
55 cleanly separated by some mechanism? Can prediction errors induce risk-seeking?

56 A proven method to provoke prediction-error related dopamine bursts in humans is to present cues and  
57 outcomes in sequential decision-making tasks, hence causing prediction errors both when options are  
58 presented, and at the time of outcome [Seymour, 2004] [Pessiglione, 2006] [Niv, 2012]. To test whether  
59 such prediction errors induce risk seeking, we used a learning task in which prediction errors are  
60 followed by choices between options with different levels of risk. If there was a clear separation of roles,  
61 then risk preferences should be independent of prediction errors. Incomplete separation, in contrast,  
62 should result in a correlation between risk preferences and preceding prediction errors. In particular, we  
63 hypothesized that positive prediction errors (expectations exceeded) should induce risk seeking, while  
64 negative prediction errors should lead to risk aversion.

65 To consider the possibility that this effect might not appear equally strongly in all participants—which  
66 could be due to differences in behavioral strategy, neural information processing or risk- and learning-

67 related traits—we also tracked pupil dilation, which is comparatively robust, and known to reflect  
68 surprising events such as prediction errors events [Preuschoff, 2011] [Browning, 2015] [Lawson, 2017]  
69 [Cavanagh, 2014]. In particular, we hypothesized that participants who experience stronger prediction  
70 errors, as indexed by pupil dilation, also show a stronger behavioral effect of risk preferences.

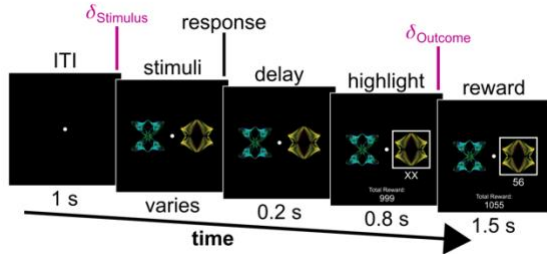
71 Our analysis proceeds in three steps: first, we conduct a model-free analysis of behavioral data to look  
72 for effects on the group level—do prediction errors make participants more risk seeking on average?  
73 Second, we move on to uncover individual differences. The effect we are interested in is likely not  
74 expressed homogenously; therefore, we use a trial-by-trial learning model to determine the effect size  
75 for individuals. Thirdly, we harness these individual differences by linking the strength of the behavioral  
76 effect to pupil dilation. This way, we validate our model on independent data, as well as explore a  
77 potential reason for the identified individual differences.

## 78 Results

### 79 The task

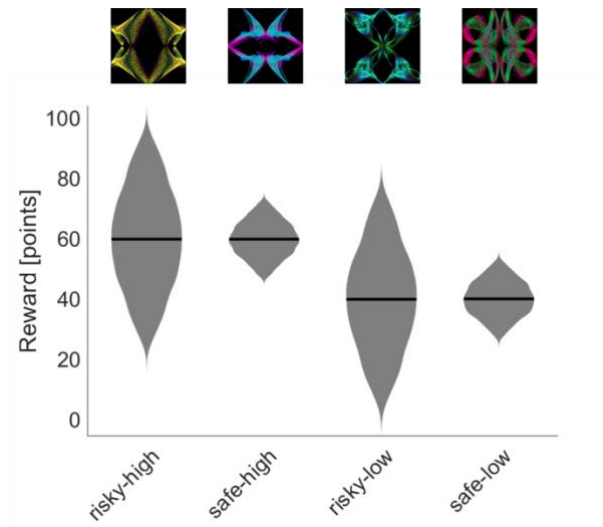
80 Our task consisted of sequences of two-alternative forced choice trials. On each trial, after an inter-trial  
81 interval (ITI) of 1 s, two stimuli (fractal images, Fig 1A) were drawn from a set of four stimuli and shown  
82 to the participant, who had to choose one. Following the choice, after a short delay of 0.8 s a numerical  
83 reward between 1 and 99 was displayed under the chosen stimulus for 1.5 s. Then, the next trial began.  
84 Participants were instructed to try to maximize the total number of reward points throughout the  
85 experiment. The reward on each trial depended on the participant's choice: each stimulus was  
86 associated with a specific reward distribution from which rewards were sampled. The four reward  
87 distributions associated with the four stimuli were approximately Gaussian and followed a two-by-two  
88 design: the mean of the Gaussian could be either high or low (60 or 40), and the standard deviation  
89 could be either large or small (20 or 5), resulting in four reward distributions in total (risky-high, risky-  
90 low, safe-high and safe-low, Fig 2B). Each participant (N=27, 3 excluded, see Methods and Fig S1)  
91 performed four blocks of 120 trials. During each block, all six possible stimuli pairings occurred equally  
92 often. Each block used a new set of four stimuli, mapped to the same four distributions.

93 A



94

B



95 Fig 1: A) Task structure. On each trial, participants were shown two out of four possible stimuli. They had  
96 to choose one of the two, which resulted in a reward. The reward was sampled from a distribution linked  
97 to the chosen stimulus. During each trial, prediction errors occurred at two times (indicated by purple  
98 lines). B) Reward structure. Each reward distribution is linked to one stimulus and is sampled from if that  
99 stimulus was chosen. The reward distributions are approximately normal; their parameters follow a two-  
100 by-two design: the mean could either be at 40 or at 60, the standard deviation could either be 5 or 20.

101

102 During each trial two distinct prediction errors occur. At stimulus onset it is revealed to participants  
103 whether the potential reward on this trial will likely be above or below average. This can be determined  
104 by considering the difference between the learned means of the two available options and the average  
105 reward associated with all four possible options. A positive prediction error occurs when the displayed  
106 options promise a higher than average reward, while a negative prediction error occurs when the  
107 expected reward is lower than average. This update of the reward prediction at stimulus onset will be  
108 called stimulus prediction error  $\delta_{Stimulus}$ . Stimulus prediction errors have previously been investigated:

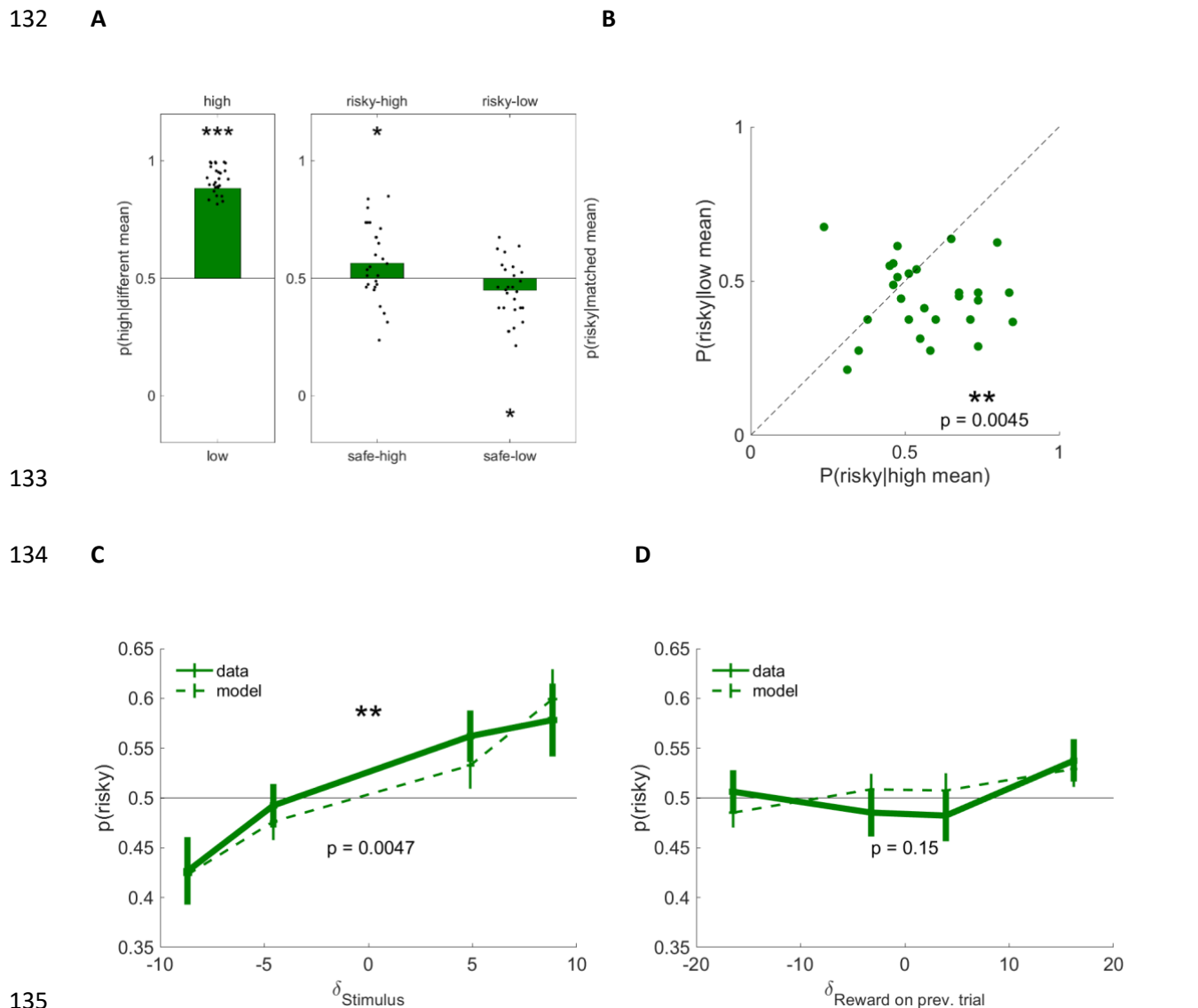
109 they are associated with phasic responses of dopamine neurons [Schultz, 1997], and have, for example,  
110 been used to assess the impact of dopamine on the formation of episodic memories [Jang, 2019]. After  
111 considering the options, the participant will make a choice, and be presented with a reward. The  
112 difference between the expectation and the actual outcome corresponds to a second reward prediction  
113 error, which we call the outcome prediction error  $\delta_{outcome}$ .

114 In our task, risk corresponds to the variability of the rewards associated with a stimulus. The reward  
115 distributions associated with some options are broad, while those related to other options are narrow  
116 (Fig 2B). It is "risky" to pick a stimulus associated with a broad reward distribution, since outcomes might  
117 deviate a lot from the expected outcome. Correspondingly, it is "safe" to pick a stimulus with a narrow  
118 distribution, since the outcomes will mostly be as expected. Note that some stimuli were matched to  
119 produce the same reward on average, while differing in variability. If participants have accurately  
120 learned the average reward of these options, then choices between those stimuli cannot be based on a  
121 value difference (since on average there is none); residual preferences must therefore be interpreted as  
122 risk preferences. Those choices between matched-mean stimuli were our way of reading out risk  
123 preferences, and we refer to such trials as matched-mean trials. In the other trials, one of the options  
124 provides substantially more reward than the other option (20 points difference on average). We refer to  
125 those trials as different-mean trials.

## 126 Behavior

127 To confirm that participants had understood the task and had learned the values associated with the  
128 four options, we first analyzed their choices in different-means trials. We observed a gradual shift from  
129 initial indifference to a strong preference for the high mean stimuli (proportion of correct choices after  
130 trial 40 > 0.5, t-test,  $t = 38.6$ ,  $p = 1.73 \times 10^{-24}$ ; see Fig 2A, first column). This suggests that participants  
131 understood the instructions and learned values accurate enough to maximize reward points.





143 *Prediction errors and value estimates were obtained by fitting a Rescorla-Wagner model to the choice*  
144 *data. Choices in matched-mean trials were binned by participant, value difference and stimulus*  
145 *prediction errors. The proportion of risky choices was then averaged across all but the prediction error*  
146 *bin. The solid green line shows the residual dependency of proportion of risky choices on stimulus*  
147 *prediction errors (error bars indicate the standard error of the mean). This binning method controls for*  
148 *confounding effects related to incidental differences in learned values and differences between*  
149 *individuals. The dashed line was obtained using the same binning method on predicted choice*  
150 *probabilities, obtained through a logistic regression fitted to predict choices from value differences,*  
151 *outcome prediction errors and participant ID (see Methods for details). D) Impact of outcome prediction*  
152 *errors on risk: identical to C), except this time using outcome prediction errors on the previous trial*  
153 *instead of stimulus prediction errors as predictor.*

154

155 In addition to this clear preference for high-mean options, we found a weak but significant preference  
156 towards the risky stimulus in risky-high versus safe-high choices (Fig 2A, second column; t-test:  $p =$   
157  $0.0343$ ,  $t = 2.23$ ), and a weak but significant preference against the risky stimulus in risky-low versus  
158 safe-low choices (Fig 2A, third column; t-test:  $p < 0.0317$ ,  $t = -2.27$ ). This suggests that on average,  
159 participants acted risk-seeking in high reward contexts, and risk-averse in low reward contexts. In  
160 addition to this group-level analysis, we investigated how preferences differed between the matched-  
161 mean conditions within each participant. We found that most of the participants were more risk seeking  
162 in the high-mean condition than in the low mean condition (Fig 2B; paired t-test:  $t = 3.11$ ,  $p = 0.0045$ ).  
163 These results are in line with previous findings [Wulff, 2018] [Madan, 2014], see Discussion for details.  
164 Next, we investigated whether these risk preferences could be due to prediction errors. Our hypothesis  
165 was that the dopamine release triggered by prediction errors might bias the participants' preferences

166 towards the risky option. To test this hypothesis, we fitted a basic Rescorla-Wagner (RW) model to each  
167 participant's behavior to obtain trial-by-trial estimates of subjective values and prediction errors (see  
168 Modelling and Methods for model specifications and fitting procedure). We then extracted both the  
169 stimulus prediction error  $\delta_{Stimulus}$  and the outcome prediction error  $\delta_{Outcome}$ , and checked whether  
170 these prediction errors were correlated with the risk preference displayed in the following choice. This  
171 was done by fitting logistic regressions to the choices recorded in matched-mean trials (see Methods for  
172 details of the procedure). We found that the probability of choosing the risky option was predicted by  
173 the stimulus prediction error, but not by the previous trial's outcome prediction error (Stimulus  
174 prediction error: Fig 2C; chi-squared test, chi-squared = 8.00, df = 1, p: 0.00468. Outcome prediction  
175 error: Fig 2D; chi-squared test, chi-squared = 2.11, df = 1, p = 0.146). This suggests that the stimulus  
176 prediction error immediately before the choice (0.97 s delay on average, with standard deviation 0.51)  
177 but not the outcome prediction error on the previous trial (3.47 s delay on average, with standard  
178 deviation 0.51) modulates risk preferences on a trial-by-trial basis.

## 179 Modelling

180 Having established that there was a correlation between prediction errors and risk-seeking, we tried to  
181 capture the effect in a reinforcement learning model. We designed a model that tracks the stimulus-  
182 specific mean rewards  $Q$ , as well as the stimulus-specific spreads  $S$ . More explicitly:  $Q_i$  represents an  
183 estimate of the average reward obtained after choosing stimulus  $i$ , while  $S_i$  represents an estimate of  
184 the mean absolute deviation or "spread" of that reward. Spread is one way to quantify risk, since it  
185 measures how unpredictable the stimulus is. Our model updates  $Q$  using a conventional Rescorla-  
186 Wagner rule,

$$187 \Delta Q_i = \alpha_Q \delta_{outcome},$$

188 where  $\delta_{outcome} = r - Q_i$  is the outcome prediction error, and  $\alpha_Q$  is the learning rate for value. The  
189 model updates the estimated spread  $S$  using a similar rule,

$$190 \quad \Delta S_i = \alpha_S (|\delta_{outcome}| - S_i),$$

191 where  $\alpha_S$  is the learning rate for risk. After sufficient burn-in, this rule produces  $S$  that fluctuate around  
192 the mean absolute deviation of the reward distributions, and hence provides an estimate of the risk  
193 associated with each stimulus. Our learning rules are analogous to plasticity rules that feature in a  
194 computational model of the basal ganglia (where the mean is encoded in the difference between  
195 synaptic weights of the direct and indirect pathway, while the spread is encoded in the sum of these  
196 weights) [Mikhael, 2016] [Moeller, 2019]. In those models, choices are based on subjective values  $V$   
197 which are assembled from the mean rewards  $Q$  and the dopamine-weighted spreads  $S$ . Following these  
198 models, we define the subjective value of reward in the following equation, where the spread is  
199 weighted by the stimulus prediction error—which is indicative of dopamine activity—on that trial:

$$200 \quad V_i = Q_i + \gamma \delta_{stimulus} S_i \quad (\text{Eq. 1})$$

201 where  $\delta_{stimulus} = \frac{1}{2} \sum_{i \in options} Q_i - \frac{1}{4} \sum_j Q_j$  (the stimulus prediction error represents the change in  
202 reward expectation before and after the presentation of the options). Note that we include the stimulus  
203 prediction error, but not the outcome prediction error on the previous trial, because only the former  
204 showed an effect on choices in our previous analysis (see Fig 2C and 2D). The parameter  $\gamma$  thus captures  
205 the extent to which recent dopaminergic prediction errors might modulate risk preference.

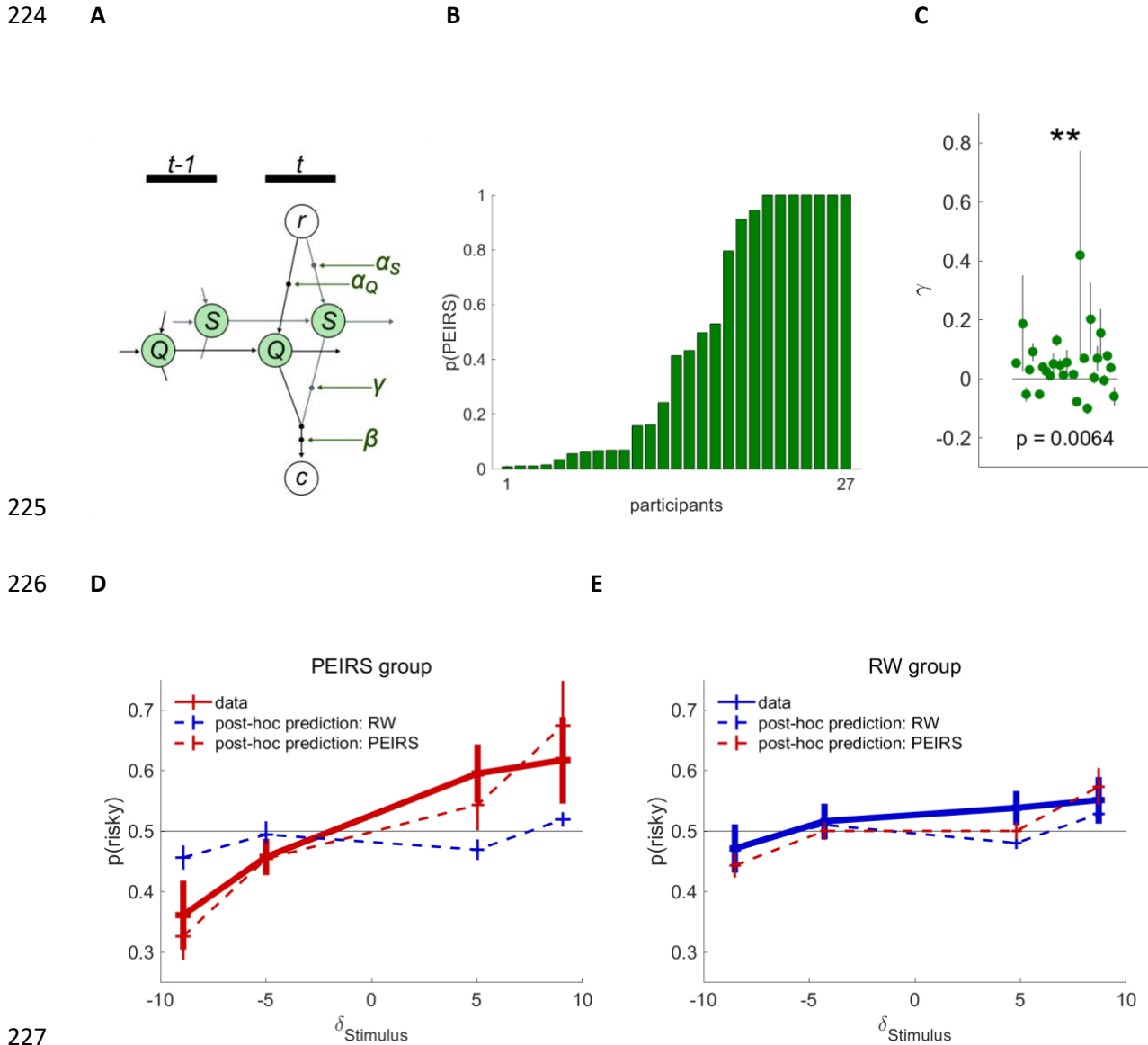
206 On every trial, subjective values  $V$  are computed from the learned  $Q$  and  $S$  for both available options.  
207 Those values are then softmax-transformed into a probability distribution, from which choices are  
208 sampled. This model, which we call Prediction Error Induced Risk Seeking (PEIRS), can be understood as

209 a generalization of the RW model, which is contained in it as a special case ( $\gamma = 0$  decouples risk from  
210 choices and recovers the conventional value based RW model).

211 We fitted our PEIRS model (as well as a conventional RW model) to the choice data of each participant  
212 individually, to obtain estimates on the strength and direction of prediction error induced risk seeking  
213 for each participant (see Methods for details). A model comparison for each participant individually  
214 showed that 11 out of 27 participants were better described by PEIRS than by RW (Fig 3B). This means  
215 that for 11 out of 27 participants in our cohort (about 40 %), prediction error induced risk seeking is  
216 strong enough to merit extra parameters.

217 We next investigated the posterior parameter distributions for  $\gamma$  that we obtained from the fit. We  
218 found that they are grouped around a positive mean significantly different from zero (Fig 3C; one-tailed  
219 t-test:  $t = 2.67$ ,  $p = 0.0064$ ). That  $\gamma$  tends to be positive across the cohort is in line with the dopaminergic  
220 interaction we propose. Note that the model did not feature any bias for  $\gamma$  to be positive; the positive  
221 tendency that we observe in the fitted values is due entirely to biases in our participant's behavior.

222 Overall, our basic analysis of behavior as well as the model comparison both suggest that there is a  
223 significant positive interaction between prediction errors and ensuing risk-seeking.



235 *Estimates of the parameter  $\gamma$  from all participants are shown (green dots). The error bars represent the*  
236 *standard deviation of the corresponding posterior distribution. D) Prediction-error induced risk seeking in*  
237 *participants for which PEIRS wins the model comparison. The red solid line represents choice data of*  
238 *these participants, the dashed lines represent choice predictions extracted from model fits. Choice data*  
239 *and choice predictions are plotted in the same way as in Fig 2C and 2D, merely displaying posterior*  
240 *predictions from our generative models instead of predictions of simple logistic models. D) Prediction-*  
241 *error induced risk seeking in participants for which RW wins the model comparison. Similar to C), but*  
242 *based on a complementary subgroup of participants.*

243

244 To check whether the model indeed captured the effect that it was intended to capture, we performed  
245 post-hoc simulations [Palminteri, 2017]. Using the posterior predictive density over choices that we  
246 obtain as an output of the fit, we generated post-hoc predictions for all choices (i.e. we used the fitted  
247 models to predict probabilities for all choices). Such predictions were generated for all participants, both  
248 from the PEIRS model and the RW model, leaving us with three data sets: a data set simulated from the  
249 fitted RW model, a data set simulated from the fitted PEIRS model and the data set obtained from  
250 humans in our experiment. We then split all these data sets according to the model comparison results  
251 (i.e. whether a participant's choices are best described by the PEIRS or by the RW model), and used the  
252 same binning scheme as for the recorded choices to check whether the two models predicted any  
253 dependency of risk preferences on reward prediction errors (Fig 3D and 3E, dashed lines).

254 The behavior simulated from the RW model did not show any substantial dependency between  
255 prediction errors and risk-taking, even when fitted to participants whose choices were better explained  
256 by the PEIRS model (blue dashed lines in Fig 3D and 3E). The PEIRS model, on the other hand, produced  
257 an approximately linear dependency between risk-taking and prediction errors for the participants best

258 described by the PEIRS model, but did not produce any dependency for the participants best described  
259 by the RW model (red dashed lines in Fig 3D and 3E). The tendencies simulated from the PEIRS model  
260 coincide with the tendencies that were observed (experimentally observed tendencies correspond to  
261 the solid lines in Fig 3D and 3E; compare the thick lines to the blue dashed lines). We concluded that the  
262 PEIRS model successfully captured our participant's risk preferences both qualitatively (linear upwards  
263 trend) and quantitatively (both intercept and slope coincide). The RW model, on the other hand, was  
264 not able to capture the risk preferences, even with fitted parameters.

265 We ran three additional tests to check the robustness of our results and the validity of our conclusions.  
266 First, we assessed the reliability of our parameter estimates by performing a standard parameter  
267 recovery analysis for both models. We found that all parameters could be recovered with little  
268 distortion, for both models and realistic parameter settings (see Fig S2). Second, we tested whether  
269 reward predictions (rather than prediction errors) might be the cause of risk-seeking. A model where  
270 reward predictions induced risk-seeking did not fit the data as well as the PEIRS model. Additional  
271 analyses based on linear models confirmed that prediction errors are more likely than reward  
272 predictions to cause the observed risk preferences (see Fig S3). Third, we tested whether our results  
273 depended on of the linearity of the utility function. We found that the interaction between risk-seeking  
274 and reward prediction errors was present even when we accounted for a nonlinear utility function (see  
275 Fig S4). These tests suggest that our results are robust if assumptions are modified, and provide  
276 additional support for our conclusions.

## 277 Pupillometry

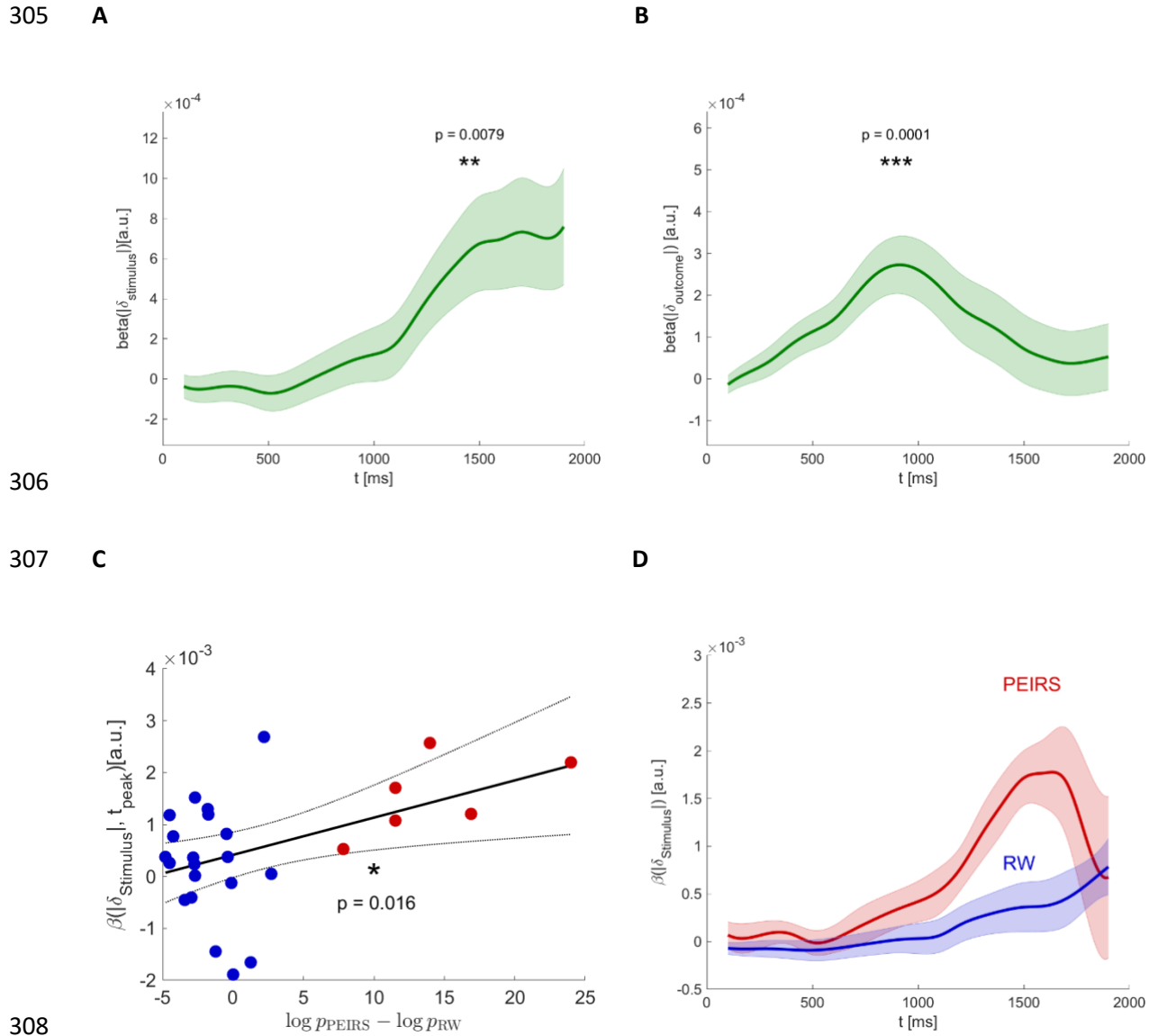
278 A range of studies have demonstrated a dilation of the pupil in response to surprising events  
279 [Preuschoff, 2011] [Browning, 2015] [Lawson, 2017] [Cavanagh, 2014]. Those phenomena have recently  
280 been synthesized into a coherent theory: pupil dilation is triggered by belief updates and scales with the



281 mismatch between prior and posterior beliefs [Zénon, 2019]. We sought to capitalize on this effect, to 1)  
282 establish the occurrence of the two prediction errors during our task through a physiological marker,  
283 and 2) to understand the individual differences suggested by our behavioral modelling better.

284 As a first step, we investigated whether pupil dilation reflected updates in reward expectation (i.e.  
285 prediction errors). We used the absolute value of the two task-related prediction errors,  $|\delta_{Stimulus}|$  and  
286  $|\delta_{Outcome}|$ , as a measure of mismatch between prior and posterior reward expectation. Trial-by-trial  
287 estimates of those prediction errors were extracted from the PEIRS model fits. Regression analyses were  
288 used to determine whether pupil dilation after stimulus onset encoded  $|\delta_{Stimulus}|$ , and whether dilation  
289 after reward presentation encoded  $|\delta_{Outcome}|$ . To avoid confounding factors such as reward or outcome  
290 prediction errors, we censored all data points collected after reward presentation in the analysis of the  
291 stimulus prediction error. For both prediction errors, we found delayed phasic responses which peaked  
292 1.6 s after stimulus onset and 0.9 s after reward presentation, respectively (Stimulus prediction error:  
293 Fig 4A; t-test:  $t = 2.89$ ,  $p = 0.0079$ . Outcome prediction error: Fig 4B; t-test:  $t = 4.61$ ,  $p = 0.00010$ .  
294 Statistical significance was established through leave-one-out unbiased peak detection, see Methods).  
295 The responses were similar for both prediction errors, except for a longer delay between stimulus  
296 prediction error onset and the peak of the pupil response. There might be many reasons for this  
297 difference in delay. Among those, differences in information processing might play a role: generating a  
298 stimulus prediction error involves two stimuli, hence attention mechanisms, in addition to retrieval of  
299 value estimates from memory. Generating the outcome prediction error, on the other hand, just  
300 requires the processing of a number.

301 We concluded that both prediction errors occur as assumed in our modelling, not only as cognitive  
302 variables, but as measurable physiological events with appropriate timing. This means that our model  
303 must at least partially represent the neural processes that occur during decision making, since it  
304 provided us with latent variables that are correlated with physiological variables in a meaningful way.



309 *Fig 4: A) Pupil dilation encodes the magnitude of the stimulus prediction errors. The green line represents*  
 310 *the regression coefficient across participants as a function of time, extracted from a mixed-effects model.*  
 311 *Responses were aligned at stimulus onset. The shaded area indicates the standard error of the estimate,*  
 312 *as provided by the regression model. For display, the trace was smoothed using spline interpolation. B)*  
 313 *Pupil dilation encodes the magnitude of the outcome prediction errors. Similar to A), with responses*  
 314 *aligned at reward presentation. C) Effect strength in behavior predicts effect strength in pupil dilation.*  
 315 *The plot shows the correlation between the log-odds that a participant generated choices according to*

316 *the PEIRS model (x-axis) and the regression coefficient for the stimulus prediction error as a predictor of*  
317 *pupil dilation, at time of maximum effect strength (y-axis). The coloring of the dots corresponds to the*  
318 *split of the cohort used in D). D) Group split to illustrate how behavior predicts pupil response. The red*  
319 *curve represents the mean pupil response across those participants that show strong prediction error*  
320 *induced risk-seeking ( $p(\text{PEIRS}) > 0.95$ ); the blue curve represents the mean pupil response of all other*  
321 *participants.*

322

323 Next, we aimed to correlate individual differences in behavior with individual differences in pupil  
324 responses. Our behavioral modelling allowed us to stratify our cohort with respect to the strength of  
325 individual prediction-error induced risk seeking (see section Modelling). We quantified the effect  
326 strength through the logarithmic odds ratios of the probability that a participant is better described by  
327 the PEIRS model than by the RW model (see Fig 3A for a plot of those probabilities). The strength of the  
328 pupil response was quantified through the respective correlation coefficient at the time of strongest  
329 effect (determined through leave-one-out peak detection to avoid bias). Using a linear model to relate  
330 pupil effect strength and behavioral effect strength, we found that participants responded stronger to  
331 stimulus prediction errors if they showed more pronounced prediction-error induced risk seeking in  
332 their behavior (Fig 4C, adjusted  $R^2$ : 0.1864,  $p < 0.05$ ). To better illustrate this, we split our cohort into  
333 two groups: those that showed significant prediction error induced risk seeking (i.e.  $p(\text{RW}) < 0.05$ , PEIRS  
334 group) and the rest ( $p(\text{RW}) > 0.05$ , RW group). While there is no noticeable response in the RW group,  
335 the PEIRS group shows a very pronounced effect (Fig 4D). Overall, this suggests that the degree to which  
336 some individual responds to the stimulus prediction error (as measured through the strength of the  
337 pupil response) is correlated to the strength of prediction error induced risk-seeking in the behavior of  
338 that individual.

339 To rule out a possible circularity that might confound these results (both the log-odds-ratio and the  
340 predictor variable for the response come from the same model fit), we conducted additional analyses  
341 based on estimates of the stimulus prediction error that were independent from the model  
342 (Supplemental Materials, Fig S5). We found that the effect appears unchanged for model-free estimates  
343 of the stimulus prediction error, which rules out the possibility of a model artefact.

## 344 Discussion

345 Different behavioral phenomena—learning from prediction errors and biased risk-preferences --are  
346 attributed to the same neuromodulator, dopamine. Using a task where reward prediction errors are  
347 immediately followed by decisions that involve risk, we showed that reward prediction errors and the  
348 probability of risk-taking are positively correlated: positive reward prediction errors induce risk seeking,  
349 negative ones inhibit it. In particular, our results show that the strength of the reward prediction error  
350 (as indexed by pupil dilation) determines the effect on risk-preferences. This result suggests that the two  
351 roles of dopamine (teaching signal and risk modulator) interfere with each other. It provides evidence  
352 against the conjecture that the roles are well separated.

353 Decision making under uncertainty has been extensively studied in behavioral economics. One main  
354 finding in this field, codified in prospect theory, is that humans tend to be risk-averse if decisions  
355 concern gains, and risk-seeking if decisions concern losses [Kahneman, 2013]. However, those classic  
356 findings rely on explicit knowledge about the probabilities involved in the decisions. Several more recent  
357 studies indicate that those tendencies reverse when risks and probabilities are learned from experience  
358 (i.e. by trial and error): if learning is incremental and based on feedback, humans tend to make risky  
359 decisions about gains and risk-averse decisions about losses [Wulff, 2018]. This phenomenon has been  
360 termed the description-experience gap. In cognitive neuroscience and psychology, some studies have  
361 reproduced this phenomenon [Madan, 2014], while others report risk aversion in the gain domain [Niv,  
362 2012]. This diversity might be associated with the degree of implicitness of the knowledge that is gained  
363 during the task: [Niv, 2012] used classical bimodal reward distributions (e.g. 40 points with probability  
364 50 %, 0 points otherwise) which participants might be able to resolve after a few trials. Here, we used a  
365 strongly random reward schedule (normal distributions, see Fig. 1B), which made anything but implicit

366 learning intractable. Our main behavioral findings (Fig. 2A and 2B) are in line with the description-  
367 experience-gap, and differ to those of [Niv, 2012].

368 The dopamine-related interpretation of our results is based on previously reported causes and  
369 consequences of phasic dopamine signaling. On the causes side, it is well established that changes in  
370 reward anticipation, brought about by reward-predicting cues, provoke dopamine bursts that originate  
371 in brain areas such as VTA, and are broadcast to structures such as the striatum and the medial  
372 prefrontal cortex [Schultz, 1997][Seymour, 2004] [Pessiglione, 2006] [Niv, 2012]. In our task, such  
373 prediction errors occurred at the presentation of the available options. On the consequences side, it has  
374 been shown that phasic dopamine activity can affect risk-preferences in decision making [Chew, 2019].  
375 Our task featured decisions between options that provided the same average reward, but different  
376 spreads of the individual rewards. Choices on those trials were not biased by value differences, and  
377 hence well suited to read out risk preferences. The simultaneous occurrence of these two dopamine-  
378 related phenomena explains our result: risk-seeking followed positive prediction errors and risk aversion  
379 followed negative prediction errors.

380 For our behavioral results, interpretations other than our dopaminergic explanation may be evoked: the  
381 behavior in a similar task [Madan, 2014] was interpreted as the result of memory replay: experiences  
382 ("Obtained reward X after choosing option Y") might not only be used for immediate value updates but  
383 might also be stored in a memory buffer. This buffer can then be used for offline learning from past  
384 experiences in times of inactivity, such as during the inter-trial interval. It was proposed by [Madan,  
385 2014] that experiences are more likely to enter the buffer if they are extreme. If entering the buffer is  
386 biased in this way, then so are the values learned from replaying those experiences. In our task,  
387 extreme might mean that the reward was extremely high or low. The corresponding bias would drive  
388 choice towards the stimuli that produce the highest rewards, and away from those that produce the  
389 lowest, and thereby lead to a pattern similar to the one we observed.

390 Which theory is closer to the truth? It is difficult to compare the memory theory directly to prediction-  
391 error induced risk-seeking; it is unclear how to obtain trial-by-trial choice predictions from the memory  
392 model, which rules out a formal model comparison. Indeed, the memory model has so far only been  
393 fitted to and assessed based on summary statistics of a large collection of trials. Further, the memory  
394 model has so far not been equipped with a mechanistic underpinning and was therefore not validated  
395 on physiological variables such as pupil dilation. In contrast, prediction-error induced risk-seeking can be  
396 fitted trial-by-trial, allowing it to make predictions not only about summary statistics but about the  
397 evolution of preferences during the task as well as about the immediate impact of extreme events such  
398 as large prediction errors. The corresponding latent variables can be correlated with physiological  
399 variables, proving that they can explain aspects of pupil dilation in addition to behavior (Fig. 3C and 3D).

400 If one interprets our results as resulting from dopaminergic interaction, one is forced to give up on the  
401 idea that direct and indirect dopaminergic effects are strictly separated. This conclusion is consistent  
402 with other recent findings: it has been shown that phasic dopamine correlates with motivational  
403 variables [Hamid, 2016] and movement vigor [da Silva, 2018] just as well as with reward prediction  
404 errors. These findings cast doubt on the separation into tonic and phasic and on separations in general.

405 In summary, our findings show that there is an interaction between prediction errors and risk seeking  
406 that matches what one would expect from dopaminergic interactions. We further show that this effect  
407 is detectable even on the individual level in a sizable part of our group, and that between-participant  
408 variability in behavior can be linked to differences in pupil responses—the stronger the pupil response  
409 to stimulus prediction errors, the stronger the prediction error induced risk seeking.

## 410 Methods

### 411 Participants

412 We tested 30 participants (15 female, median age: 26, range: 18-42). Our participants did not suffer  
413 from visual, motor or cognitive impairments. They were recruited and tested voluntarily, all  
414 experimental procedures were approved by the local ethics committee. Our results are based on 27 of  
415 the 30 participants. Three participants were excluded from the analysis due to their failure to  
416 understand the task. We evaluated the participants' understanding of the task by scoring their  
417 preferences in mixed-mean choices during the second half of the blocks. Participants were included in  
418 the analysis if they chose the high-valued option in more than 70 % of those trials (Fig S1).

### 419 Logistic regressions

420 Logistic regressions were conducted using mixed-effects modelling. The target variable  $y$  was defined as  
421  $y = 1$  if the risky option was chosen and  $y = 0$  else. The predictors of interest were the prediction  
422 errors  $\delta_{Stimulus}$  and  $\delta_{Outcome}$  that preceded the choice. We further included  $Q_{risky\ option} -$   
423  $Q_{save\ option}$  as a predictor to control for residual value differences. Individual differences were  
424 accounted for by a random intercept and random slopes for each predictor. The p-values we report for  
425 single predictors were obtained from chi-squared tests on likelihood ratio statistics. Those were  
426 computed through comparisons between the fit with all predictors included and the fit without the  
427 predictor of interest (but with the respective random slope).

### 428 Models

429 The RW model as well as the PEIRS model feature a softmax choice rule:

$$430 \quad P(\text{choice} = i) = \frac{e^{\beta V_i}}{\sum_j e^{\beta V_j}}$$



431 The models differ in how those subjective values  $V$  are constructed. In the RW model, the subjective  
432 value of an option is simply the learned value of this option:  $V_{RW,i} = Q_i$ . In the PEIRS model, the  
433 subjective value is determined according to Eq. 1. For both the PEIRS and the RW model we set the  
434 initial value  $Q_0$  to the empirical mean of 50. For the PEIRS model the initial value of the spread  $S_0$  is left  
435 as a free parameter. All in all, the RW model features two free parameters  $(\alpha_Q, \beta)$ , while the PEIRS  
436 model features five  $(\alpha_Q, \alpha_S, \beta, \gamma, S_0)$ .

### 437 [Model fit, comparison and regularization](#)

438 Fits and model comparisons were performed using the VBA toolbox [Daunizeau, 2014]. This toolbox  
439 implements a Variational Bayes scheme. It takes a set of measurements, a generative probabilistic  
440 model that describes how the measurements arise (which usually contains some latent, i.e. unobserved,  
441 variables) and prior distributions over the model parameters as input, and outputs among other things  
442 an approximate posterior distribution over model parameters, an approximate posterior distribution  
443 over the latent variables, and an upper bound for the model evidence. We fitted both models to each  
444 participant.

445 Our model comparison is based on the approximate model evidences  $L(model|data)$  that the toolbox  
446 provides. Assuming that a participant generated data according one of the models, the probability of  
447 that participant using model  $m$  is given by  $p(m|data) = \frac{L(m|data)}{\sum_{m'} L(m'|data)}$  (see the documentation of the  
448 VBA toolbox or [Stephan, 2009] for reference). We use these probabilities as an index of effect strength.

449 We estimate parameters using the posterior distributions over parameters that the toolbox outputs.  
450 Point estimates of parameters are obtained by computing the expected value of the posterior  
451 distributions. The same procedure is applied for latent variables, such as values and prediction errors.

452 The questions we pursue in this study involve physiological factors, such as dopamine and pupil dilation.  
453 To be useful to answer our questions and make valid predictions, it is important that our models  
454 operate in a physiologically plausible regime. One important requirement was that strong, systematic  
455 prediction errors should only occur during the learning phase at the beginning of each block, and not  
456 persist after choice behavior has stabilized. We found that our models did not fulfil this requirement by  
457 default, and hence introduced a regularization: from trial 61 onwards, we penalized prediction errors by  
458 introducing a prior centered around zero. This was implemented by adding an additional observed  
459 variable which was normally distributed around the outcome prediction error:  $\delta_{outcome}^{observed} \sim$   
460  $N(\delta_{outcome}, \sigma^2)$ . We then provided the model with “measurements”  $\delta_{outcome}^{observed} = 0$  for trials 61 to 100.  
461 During model inversion, those “observations” penalized  $\delta_{outcome}$  that differed strongly from zero. This  
462 regularization applies to both the RW and the PEIRS model.

## 463 Pupillometry

464 We recorded time series of pupil diameters for every trial, using an EyeLink 1000 system. The raw  
465 measurements were preprocessed (smoothing, blink correction) using standard methods [Manohar,  
466 2019]. Then, the traces were aligned to the relevant temporal markers (stimulus onset, or reward  
467 onset). We used the mean over 500 ms prior to the alignment point to define a trial-wise baseline. All  
468 traces were divided and shifted by that baseline, resulting in traces reflecting the relative change of  
469 pupil diameter after the alignment point. Finally, traces were downsampled to 10 Hz.

470 To uncover the pupil response to the stimulus prediction error, we aligned the pupil time courses at  
471 stimulus onset. After stimulus onset, participants would eventually make a choice (with variable delay,  
472 the median reaction time was 0.86 s) and receive a reward (with a 1 s delay) after their choice. Since the  
473 reward or the resulting outcome prediction error might confound our regression analysis, we censored  
474 out all data after reward presentation. This means that the number of observations on which

475 regressions can be based rapidly declines after the median reward presentation time, which is at 1.86 s  
476 after stimulus onset. Estimates obtained later are increasingly unreliable, since they are based on  
477 insufficient data. We hence conducted our analyses for the interval 0 s to 1.9 s after stimulus onset. This  
478 allows us to obtain reliable estimates of the statistics, while still avoiding confounding effects related to  
479 reward presentation.

480 To test whether the pupil responses to the prediction errors are statistically significant, we needed to  
481 perform a test in a single time-point corresponding to the largest effect. To avoid circularity, the time of  
482 peak effect was identified using a leave-one-out method: We first calculated time-series of regression  
483 weights for each participant individually. Then, for each participant, we determined the peak effect  
484 strength of the response. To achieve this without introducing bias, we temporarily excluded the  
485 participant in question and determined the time bin in which the responses of the other participants  
486 were most significant. This was achieved by executing t-tests on the response strengths in each time  
487 bin, and selecting the bin with the smallest p-value. We then took the left-out participant's response  
488 strength from that time bin, considering it their response strength at the peak of the group response. In  
489 a final step, we pooled all those individual response strengths at peak effect and used a t-test to check  
490 whether they deviated significantly from zero.

## 491 Acknowledgements

492 This work has been supported by MRC grants MC\_UU\_12024/5, MC\_UU\_00003/1 and BBSRC grant  
493 BB/S006338/1 held by RB, an MRC clinician scientist fellowship MR/P00878X to SGM and  
494 studentships from the MRC (MR/K501256/1 and MR/N013468/1) and St John's College to JG.

## 495 References

496 [Schultz, 1997] Schultz, W., Dayan, P., & Montague, P. R. (1997). A neural substrate of prediction and  
497 reward. *Science*, 275(5306), 1593-1599.

498 [Rescorla, 1972] Rescorla, R. A., & Wagner, A. R. (1972). A theory of Pavlovian conditioning: Variations in  
499 the effectiveness of reinforcement and nonreinforcement. *Classical conditioning II: Current research and*  
500 *theory*, 2, 64-99.

501 [Berridge, 2012] Berridge, K. C. (2012). From prediction error to incentive salience: mesolimbic  
502 computation of reward motivation. *European Journal of Neuroscience*, 35(7), 1124-1143.

503 [Niv, 2007] Niv, Y. (2007). Cost, benefit, tonic, phasic: what do response rates tell us about dopamine  
504 and motivation?. *ANNALS-NEW YORK ACADEMY OF SCIENCES*, 1104, 357.

505 [Berke, 2018] Berke, J. D. (2018). What does dopamine mean? *Nature neuroscience*, 21(6), 787.

506 [Pessiglione, 2006] Pessiglione, M., Seymour, B., Flandin, G., Dolan, R. J., & Frith, C.D. (2006). Dopamine-  
507 dependent prediction errors underpin reward-seeking behaviour in humans. *Nature*, 442(7106), 1042.

508 [Chew, 2019] Chew, B., Hauser, T. U., Papoutsis, M., Magerkurth, J., Dolan, R. J., & Rutledge, R. B. (2019).  
509 Endogenous fluctuations in the dopaminergic midbrain drive behavioral choice variability. *Proceedings*  
510 *of the National Academy of Sciences*, 116(37), 18732-18737.

511 [St Onge, 2009] St Onge, J. R., & Floresco, S. B. (2009). Dopaminergic modulation of risk-based decision  
512 making. *Neuropsychopharmacology*, 34(3), 681.

513 [Preuschoff, 2011] Preuschoff, K., t Hart, B. M., & Einhauser, W. (2011). Pupil dilation signals surprise:  
514 Evidence for noradrenaline's role in decision making. *Frontiers in neuroscience*, 5, 115.

515 [Palminteri, 2017] Palminteri, S., Wyart, V., & Koehlin, E. (2017). The importance of falsification in  
516 computational cognitive modeling. *Trends in cognitive sciences*, 21(6), 425-433.

517 [Madan, 2014] Madan, C. R., Ludvig, E. A., & Spetch, M. L. (2014). Remembering the best and worst of  
518 times: Memories for extreme outcomes bias risky decisions. *Psychonomic bulletin & review*, 21(3), 629-  
519 636.

520 [Hamid, 2016] Hamid, A. A., Pettibone, J. R., Mabrouk, O. S., Hetrick, V. L., Schmidt, R., Vander Weele, C.  
521 M., ... & Berke, J. D. (2016). Mesolimbic dopamine signals the value of work. *Nature neuroscience*, 19(1),  
522 117.

523 [Engelhard, 2019] Engelhard, B., Finkelstein, J., Cox, J., Fleming, W., Jang, H. J., Ornelas, S., ... & Witten, I.  
524 B. (2019). Specialized coding of sensory, motor and cognitive variables in VTA dopamine neurons.  
525 *Nature*, 1.

526 [Daunizeau, 2014] Daunizeau, J., Adam, V., & Rigoux, L. (2014). VBA: a probabilistic treatment of  
527 nonlinear models for neurobiological and behavioural data. *PLoS Computational Biology*, 10(1),  
528 e1003441.

529 [Zénon, 2019] Zénon, A. (2019). Eye pupil signals information gain. *Proceedings of the Royal Society B*,  
530 286(1911), 20191593.

- 531 [Browning, 2015] Browning, M., Behrens, T. E., Jocham, G., O'reilly, J. X., & Bishop, S. J. (2015). Anxious  
532 individuals have difficulty learning the causal statistics of aversive environments. *Nature neuroscience*,  
533 *18*(4), 590.
- 534 [Lawson, 2017] Lawson, R. P., Mathys, C., & Rees, G. (2017). Adults with autism overestimate the  
535 volatility of the sensory environment. *Nature neuroscience*, *20*(9), 1293.
- 536 [Cavanagh, 2014] Cavanagh, J. F., Wiecki, T. V., Kochar, A., & Frank, M. J. (2014). Eye tracking and  
537 pupillometry are indicators of dissociable latent decision processes. *Journal of Experimental Psychology:*  
538 *General*, *143*(4), 1476.
- 539 [Mikhael, 2016] Mikhael, J. G., & Bogacz, R. (2016). Learning reward uncertainty in the basal ganglia.  
540 *PLoS computational biology*, *12*(9), e1005062.
- 541 [Moeller, 2019] Möller, M., & Bogacz, R. (2019). Learning the payoffs and costs of actions. *PLoS*  
542 *computational biology*, *15* (2), e1006285.
- 543 [Voon, 2006] Voon, V., Hassan, K., Zurowski, M., Duff-Canning, S., De Souza, M., Fox, S., ... & Miyasaki, J.  
544 (2006). Prospective prevalence of pathologic gambling and medication association in Parkinson disease.  
545 *Neurology*, *66*(11), 1750-1752.
- 546 [Weintraub, 2010] Weintraub, D., Koester, J., Potenza, M. N., Siderowf, A. D., Stacy, M., Voon, V., ... &  
547 Lang, A. E. (2010). Impulse control disorders in Parkinson disease: a cross-sectional study of 3090  
548 patients. *Archives of neurology*, *67*(5), 589-595.
- 549 [Gallagher, 2007] Gallagher, D. A., O'Sullivan, S. S., Evans, A. H., Lees, A. J., & Schrag, A. (2007).  
550 Pathological gambling in Parkinson's disease: risk factors and differences from dopamine dysregulation.  
551 An analysis of published case series. *Movement disorders: official journal of the Movement Disorder*  
552 *Society*, *22*(12), 1757-1763.

553 [Niv, 2012] Niv, Y., Edlund, J. A., Dayan, P., & O'Doherty, J. P. (2012). Neural prediction errors reveal a  
554 risk-sensitive reinforcement-learning process in the human brain. *Journal of Neuroscience*, 32(2), 551-  
555 562.

556 [Kahneman, 2013] Kahneman, D., & Tversky, A. (2013). Prospect theory: An analysis of decision under  
557 risk. In *Handbook of the fundamentals of financial decision making: Part I* (pp. 99-127).

558 [Wulff, 2018] Wulff, D. U., Mergenthaler-Canseco, M., & Hertwig, R. (2018). A meta-analytic review of  
559 two modes of learning and the description-experience gap. *Psychological bulletin*, 144(2), 140.

560 [Jang, 2019] Jang, A. I., Nassar, M. R., Dillon, D. G., & Frank, M. J. (2019). Positive reward prediction  
561 errors during decision-making strengthen memory encoding. *Nature human behaviour*, 1.

562 [Stephan, 2009] Stephan, K. E., Penny, W. D., Daunizeau, J., Moran, R. J., & Friston, K. J. (2009). Bayesian  
563 model selection for group studies. *Neuroimage*, 46(4), 1004-1017.

564 [da Silva, 2018] da Silva, J. A., Tecuapetla, F., Paixão, V., & Costa, R. M. (2018). Dopamine neuron activity  
565 before action initiation gates and invigorates future movements. *Nature*, 554(7691), 244.

566 [Seymour, 2004] Seymour, B., O'Doherty, J. P., Dayan, P., Koltzenburg, M., Jones, A. K., Dolan, R. J., ... &  
567 Frackowiak, R. S. (2004). Temporal difference models describe higher-order learning in humans. *Nature*,  
568 429(6992), 664.

569 [Manohar, 2019] Manohar, S. G. (2019, October 19). Matlib: MATLAB tools for plotting, data analysis,  
570 eye tracking and experiment design (Public). <https://doi.org/10.17605/OSF.IO/VMABG>

# 1 Reward prediction errors induce risk-seeking

2 Supplementary Materials

3 Moritz Moeller\*1, Jan Grohn\*2, Sanjay Manohar\*\*12+, Rafal Bogacz\*\*1.

4 \*: equal contributions

5 \*\*: equal contributions

6 1: Nuffield Department of Clinical Neurosciences, University of Oxford

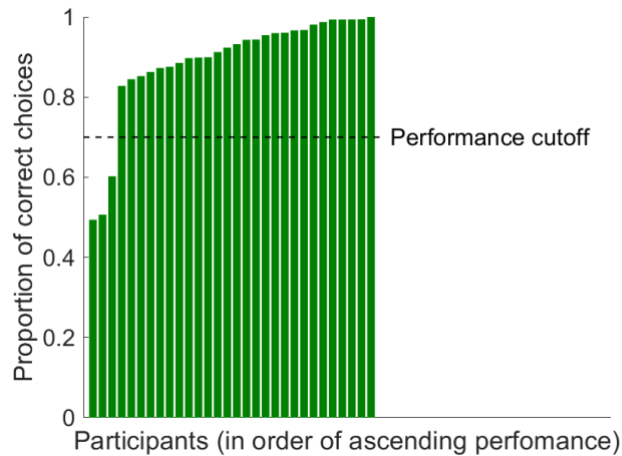
7 2: Department of Experimental Psychology, University of Oxford

8 +: corresponding author ([sanjay.manohar@psy.ox.ac.uk](mailto:sanjay.manohar@psy.ox.ac.uk))



## 9 Performance evaluation

10 We assessed individual performances post-hoc, using the proportion of correct choices after trial 60 as  
11 the criterion. The data are provided in Fig S1.



12

13 *Fig S1: Individual performance evaluation. Every bar represents one participant. The height of the bar*  
14 *indicates the proportion of correct choices (i.e. choosing the high-valued option over the low-valued*  
15 *option) in the last 60 trials of each block. Three participants (corresponding to the first three bars in this*  
16 *graph) fell below our cutoff of 70 % and were hence excluded from the study.*

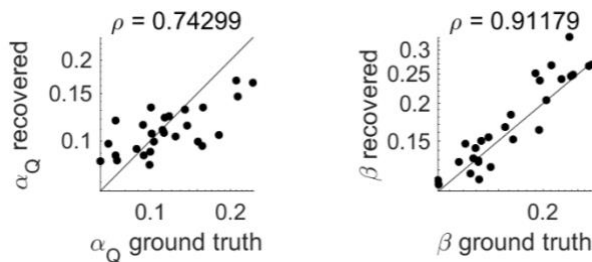
17

## 18 Parameter recovery

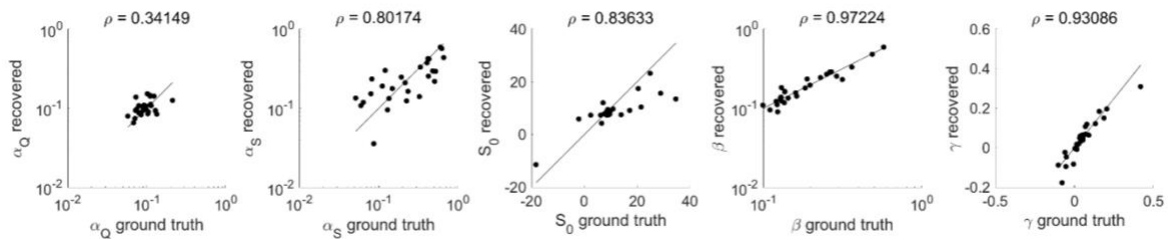
19 To assess the reliability of the parameter estimates produced by our fitting procedure, and hence build  
20 confidence in the conclusions based on the fits, we conducted a parameter recovery analysis for both  
21 models. For each model and participant, we used the posterior distributions over parameter space  
22 obtained from the fit to get point estimates of the parameters that best describe the recorded behavior.  
23 The resulting parameter set was then used to run a simulation, aiming to produce simulated data with  
24 characteristics like those of the empirical data. As a next step, we fitted the same model that was used

25 for simulation to the simulated data, and again obtained estimates of the parameters, which could now  
26 be compared with the ground truth (the parameters that were used to simulate, and that were  
27 supposed to be recovered). For both models, we could robustly recover all parameters with minimal  
28 distortion (Fig S2).

29 **A**



31 **B**



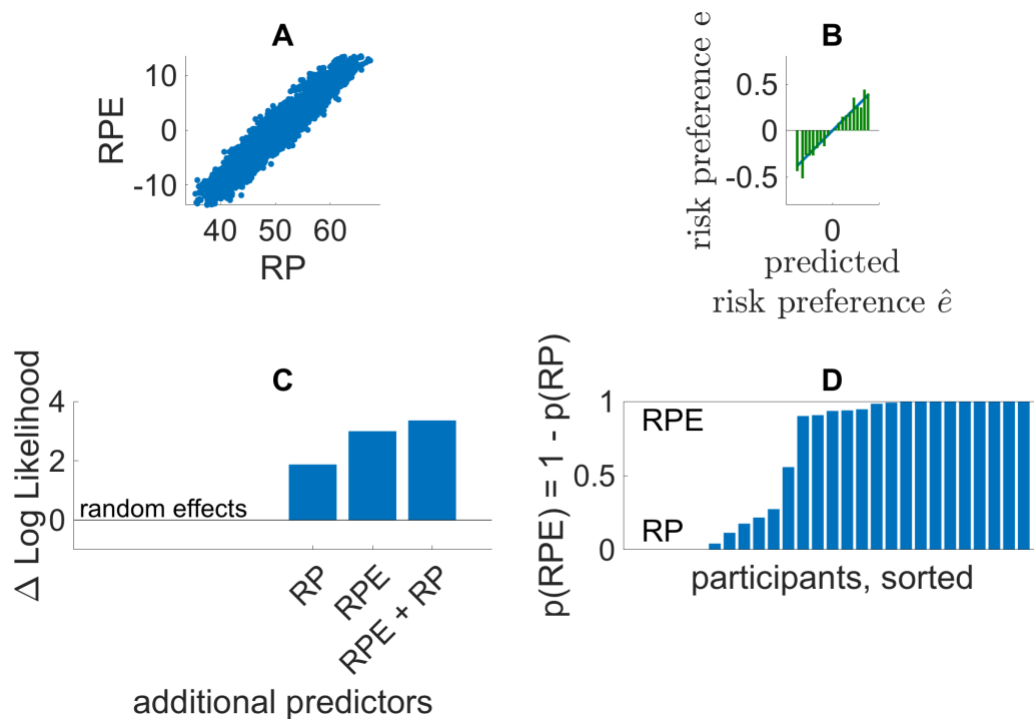
33 *Fig S2: Parameter recovery for reinforcement learning models. A) RW model. Each panel corresponds to*  
34 *one parameter of the RW model. The x-axes correspond to the ground truth values (those used for the*  
35 *simulations), the y-axes correspond to the recovered parameters extracted from fits to the simulated*  
36 *data. Every dot corresponds to one participant. Above the plot, we report the correlation coefficient*  
37 *between the actual and the recovered parameters across the population. The black line indicates*  
38 *equality, i.e.  $x = y$ . B) Parameter recovery for the PEIRS model. Same conventions as in A).*

39

40 We conclude that our models do not suffer from ambiguity or under-determination/over-  
41 parametrization. This means that both estimates of parameters and latent variables are to be taken as  
42 meaningful, unambiguous quantities.

### 43 Might reward predictions cause risk preferences?

44 We showed that seemingly irrational risk preferences can be explained by reward prediction errors  
45 (RPEs, changes in reward expectation) that occur immediately before the choice. A confounding variable  
46 in this analysis is the reward prediction (RP) itself: it could be that the anticipation of high rewards  
47 causes risk-seeking, while anticipating low rewards causes risk-aversion. If so, it would still seem as if  
48 RPEs cause risk preferences, since RPEs and RPs are highly correlated in our experiment (Fig S3 A).



49  
50 *Fig S3: Reward predictions as a confounding variable. A) Correlation between trial-wise reward*  
51 *prediction errors (RPE) and the reward predictions (RP) that followed stimulus presentation. These*  
52 *prediction errors were estimated by fitting a standard RW model to each participant individually. B)*

53 *Model fitted to residual preferences. Risk preferences were predicted with a linear mixed effects model.*  
54 *Preference was modelled as a linear function of RP and RPE. Individual differences in regression*  
55 *coefficients were modelled as random effects of subject ID ( $e \sim 1 + RP + RPE + (1 + RP + RPE|ID)$ ),*  
56 *the last term corresponds to individual differences in slopes and intercept). Residual preferences were*  
57 *binned according to predicted preferences, averaged per bin and are displayed as green bars. The*  
58 *predicted preferences are represented by the blue line. C) Relative log likelihoods of models with different*  
59 *sets of predictor variables. The bars indicate the increase in likelihood relative to a baseline model that*  
60 *used only random effects. D) Model comparison results on participant level. The bars represent the*  
61 *likelihood that the data recorded from a given participant was generated from the PEIRS model rather*  
62 *than the PIRS model.*

63

64 To test whether risk preferences are due to RPEs rather than RPs, we conducted two additional  
65 analyses. First, we used linear models to test which signal—RPEs or RPs—is a better predictor for risk  
66 preferences. We started by extracting preferences  $e$  that could not be explained by standard learning  
67 effects. This was done by predicting choices  $\hat{c}$  on matched-mean trials with a standard RW model, and  
68 subtracting them from the measured choices  $c$  ( $c = 1$  when the risky option was chosen, and  $c = 0$   
69 otherwise). The residual preferences  $e = c - \hat{c}$  contain the risk preferences we seek to explain. Next,  
70 we used linear models to predict the residual preferences  $e$ . As predictors, we considered RPE, RP and  
71 the corresponding random effects. Taken together, those signals partially explain the residual  
72 preferences (adjusted  $R^2$ : 0.0603, Fig S3 B). Finally, we checked how much explanatory power each  
73 signal contributes by comparing log likelihoods (LL) corresponding to different predictors. If risk  
74 preferences were due to RPEs, we should expect that 1) adding only RPEs should increase LL more than  
75 adding only RPs, and that 2) adding RPs on top of RPEs should not increase LL substantially. Point 2)  
76 specifically holds if RP does not contain additional relevant information about risk-preferences over and

77 above those that it shares with RPE. We found that 1) and 2) hold (Fig S3 C). This suggests that risk  
78 preferences are best explained by RPEs. In our experiment, they can also be predicted by RPs, but only  
79 because RPs are correlated with RPEs and thereby gain some of the RPEs' predictive power.

80 We run another analysis to corroborate this result: to test whether RPs could explain our data better  
81 than RPEs, we defined another trial-by-trial model (PIRS, "Predictions Induce Risk Seeking") similar to  
82 the PEIRS model. PIRS is identical to PEIRS with one exception: it is the prediction and not the prediction  
83 error that interacts with risk in the decision rule (Eq. 1 in the main text). We then performed a model  
84 comparison between PIRS and PEIRS. We found that our data is more likely to be generated by PEIRS  
85 than by PIRS (odds ratio about 3:1 for PEIRS), and that most participants are better fitted by PEIRS (Fig  
86 S3 D). This result aligns with the result we obtained using linear models to predict residual preferences,  
87 and suggests that it is the RPE, and not the RP, that might cause risk preferences.

## 88 Nonlinear utility

89 To set up models such as ours, one must choose a way to relate the point score that participants are  
90 shown to the abstract reward signal that features in RL models (i.e. one must choose a utility function  
91 that maps points to reward). For our analysis, we chose a simple linear mapping. Thus, we implicitly  
92 assume that points would directly translate into reward. However, it has been shown that often, utility  
93 functions are not as simple—for example, in behavioral economics the utility of money is frequently  
94 modeled using concave functions. Crucially, nonlinear utility functions can lead to apparent risk  
95 preferences. Do our results and conclusions still hold if we drop the assumption of a linear utility  
96 function, and allow for non-linear utility curves?

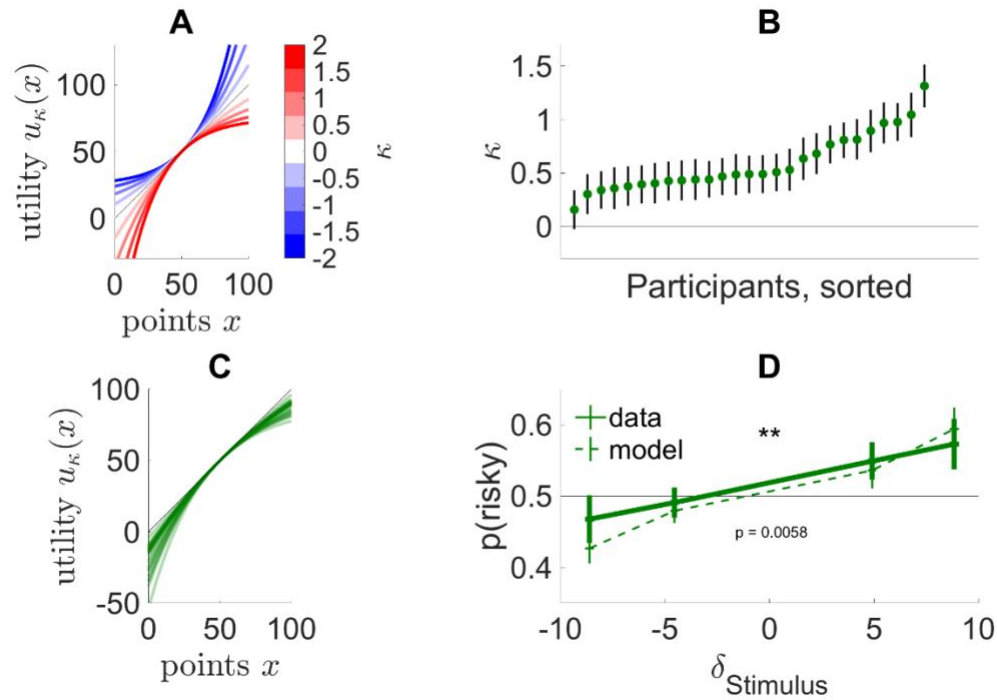
97 To test this, we started by choosing a parametric family of utility functions. We then determined the  
98 most likely utility function for each participant by fitting a RW-model with parametric utility to their

99 choices. Finally, we checked whether there was still a correlation between risk preferences and  
100 preceding prediction errors after the nonlinear utility curves were considered.

101 To model non-linear utility curves, we chose an exponential family centered at 50 points, defined by

102 
$$u_{\kappa}(x) = 50 + \frac{50}{\kappa} \left( 1 - e^{-\kappa \left( \frac{x}{50} - 1 \right)} \right).$$

103 The functions are shifted such that  $u(50) = 50$  for all values of  $\kappa$ , to keep initial values independent of  
104  $\kappa$  (see Fig S4 A for some exemplars). Next, we fitted standard RW models to the choices of each  
105 participants, using  $r_{\kappa} = u_{\kappa}(x)$ . From this, we obtained estimates of  $\kappa$  for each participant. We found  
106 that almost all  $\kappa$  were positive (Fig S4 B), suggesting concave mappings from points to subjective value  
107 for almost all participants (Fig S4 C). Finally, we performed the same analysis as in Fig 2C in the main  
108 manuscript, checking whether there was a correlation between the likelihood of risk-seeking and the  
109 magnitude of the prediction error immediately before the choice. We found a significant correlation  
110 similar to the corresponding curve based on linear utility (Fig S4 D). This suggests that our findings are  
111 robust, and hold even when the assumption of linear utility is relaxed.



112

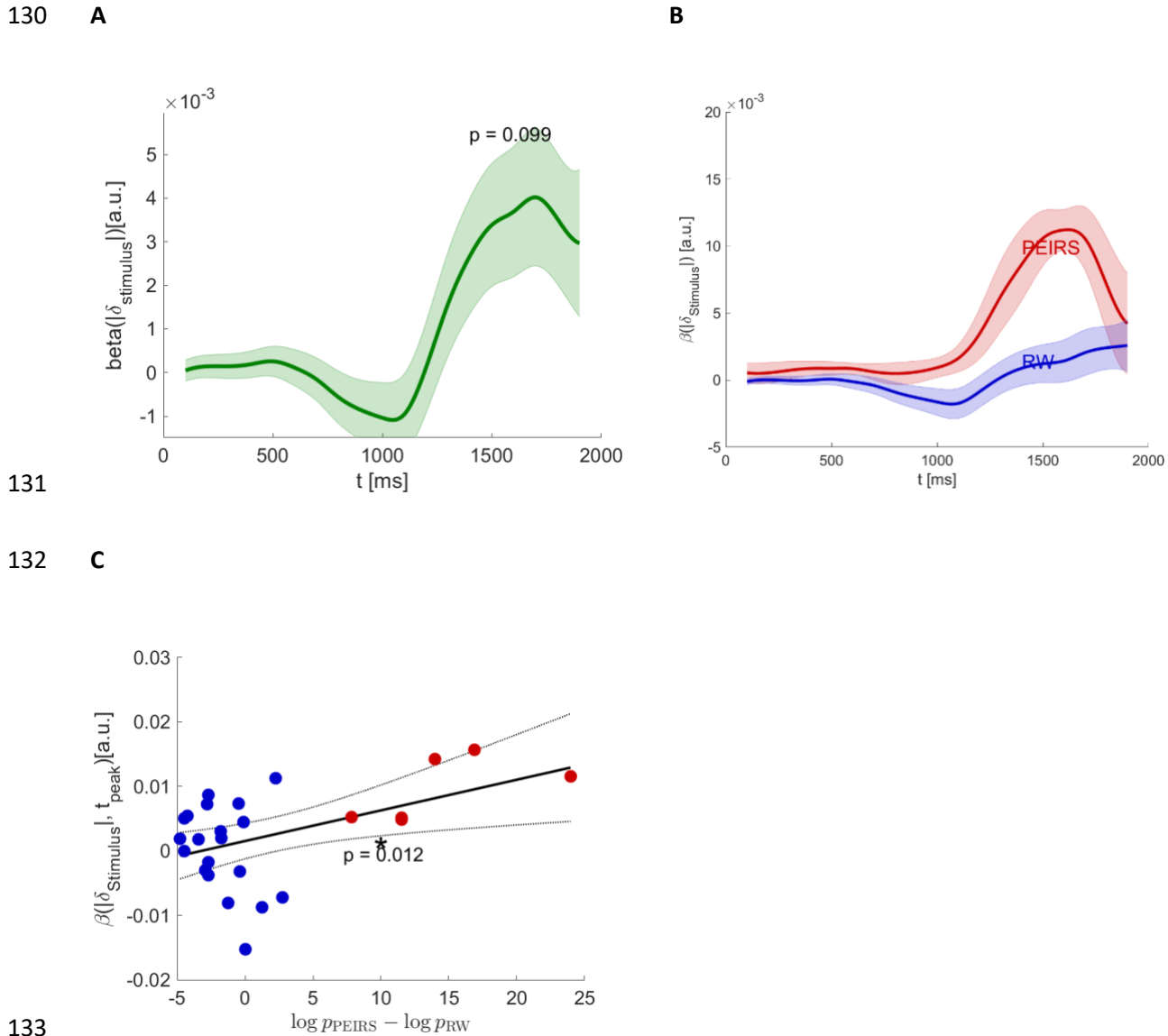
113 *Fig S4: Robustness under nonlinear utility mappings. A) Family of utility functions. The parameter  $\kappa$*   
 114 *which controls the curvature is represented by color. B) Estimates for  $\kappa$ . Mean and standard deviation of*  
 115 *the posterior distribution of  $\kappa$  are indicated by a green dot and black lines. The statistics for the posterior*  
 116 *distribution are provided for each participant individually, ordered by the mean of the posterior. C)*  
 117 *Estimated utility curves. Posterior estimates in C) were converted in utility functions and superimposed.*  
 118 *Each green line corresponds to the most likely utility function of one participant. The lines are*  
 119 *transparent to aid visibility. D) Similar to Fig 2C in the main text, but with stimulus prediction errors and*  
 120 *values taken from a RW model with the non-linear utility functions depicted in C).*

## 121 Ground truth pupillometry

122 The predictor variables of our pupil-related regression analyses (the absolute stimulus prediction error  
 123 and the absolute outcome prediction error) are model-dependent variables. One might thus suspect  
 124 that the correlations displayed in Fig 4C and Fig 4D might be spurious: pupil responses are defined with

125 respect to a model variable (the stimulus prediction error) and are predicted by another model-  
126 depended quantity (the logarithmic odds ratio). To rule out potential confounding effects, and to make  
127 sure that the pupil responses do in fact provide an external validation of our behavioral modelling, we  
128 conducted the same analyses based on the ground truth (model-free) prediction error instead of the  
129 model-based stimulus prediction error (Fig S5).





134 *Fig S5: Ground truth pupillometry results. A, B and C) same as Fig 4 A, C and D) but using model-*  
135 *independent ground truth prediction error instead of the prediction error extracted from the model fit.*

136

137 We found that all results described in the main text hold similarly if the analysis is conducted on  
138 the ground truth instead of the model-based variable. Our reasoning is thus not circular, and the result  
139 are not due to modelling artifacts.

140 The ground truth stimulus prediction error is defined as

141 
$$\delta_{Stimulus,GT} = E_{option\ shown}(R) - E_{all\ option}(R) = E_{option\ shown}(R) - 50$$

142 Since  $E_{option\ shown}(R)$  could only take the values 40, 50 or 60 by experimental design,  $\delta_{Stimulus,GT}$

143 could only take the values -10, 0 or 10, and the predictor variable  $|\delta_{Stimulus,GT}|$  could only take the

144 values 0 or 10. Therefore, the ground-truth prediction error used for the control analyses is equivalent

145 to a contrast between matched-mean and different-mean conditions.