# 1  Uncertainty-guided learning with scaled

# 2  prediction errors in the basal ganglia

3  Moritz Moeller 1, Sanjay Manohar 12, Rafal Bogacz 1+.

4  1: Nuffield Department of Clinical Neurosciences, University of Oxford

5  2: Department of Experimental Psychology, University of Oxford

6  +: corresponding author (rafal.bogacz@ndcn.ox.ac.uk)

7

8  Short title: Learning with scaled prediction errors

# Abstract

To accurately predict rewards associated with states or actions, the variability of observations has to be taken into account. In particular, when the observations are noisy, the individual rewards should have less influence on tracking of average reward, and the estimate of the mean reward should be updated to a smaller extent after each observation. However, it is not known how the magnitude of the observation noise might be tracked and used to control prediction updates in the brain reward system. Here, we introduce a new model that uses simple, tractable learning rules that track the mean and standard deviation of reward, and leverages prediction errors scaled by uncertainty as the central feedback signal. We provide a normative analysis, comparing the performance of the new model with that of conventional models in a value tracking task. We find that the new model has an advantage over conventional models when tested across various levels of observation noise. Further, we propose a possible biological implementation of the model in the basal ganglia circuit. The scaled prediction error feedback signal is consistent with experimental findings concerning dopamine prediction error scaling relative to reward magnitude, and the update rules are found to be consistent with many features of striatal plasticity. Our results span across the levels of implementation, algorithm, and computation, and might have important implications for understanding the dopaminergic prediction error signal and its relation to adaptive and effective learning.

## Author Summary

The basal ganglia system is a collection of subcortical nuclei in the mammalian brain. This system and its dopaminergic inputs are associated with learning from rewards. Here, dopamine is thought to signal errors in reward prediction. The structure and function of the basal ganglia system are not fully understood yet—for example, the basal ganglia are split into two antagonistic pathways, but the reason for this split and the role of the two pathways are unknown. Further, it has been found that under some circumstances, rewards of different sizes lead to dopamine responses of similar size, which cannot be explained with the reward prediction error theory. Here, we propose a new model of learning in the basal ganglia—the scaled prediction error model. According to our model, both reward average and reward uncertainty are tracked and represented in the two basal ganglia pathways. The learned reward uncertainty is then used to scale dopaminergic reward prediction errors, which effectively renders learning adaptive to reward noise. We show that such learning is more robust than learning from unscaled prediction errors and that it explains several physiological features of the basal ganglia system.

## Introduction

39

For any organism, better decisions result in better chances of survival. Reward prediction is an important

40

aspect of this—for example, if an organism can predict the size of a food reward associated with some

41

behavior, it can decide whether it is worth to engage in that behavior or not. Reward predictions are

42

typically based on values learned from previous reward observations. An extensive literature describes the

43

role of reward prediction in behavior, as well as the related neural mechanisms (1).

44

Piray and Daw (2) argue that when trying to predict rewards, the organism faces two challenges. The first

45

challenge is the dynamic nature of the environment: reward sizes and contingencies might change over

46

time, in ways that cannot be predicted. Such genuine changes in the environment can be quantified by the

47

typical rate of change, which is called *process noise*. The second challenge is *observation noise*: even if

48

the environment is stable, rewards will vary from experience to experience. This could be due to the

49

random nature of the environment, but also to variability in the organism's own behavior, or to noise in

50

the organism's perception and evaluation systems.

51

The stock market serves as a nice example of the two types of noise: consider the day-to-day change of a

52

stock price as a reward signal (if the stock price rises from 20 GBP to 21 GBP overnight, then the

53

shareholders win 1 GBP in that transition). Most of the variability of that signal will be due to random

54

fluctuations—this can be classified as observation noise. However, a part of the signal's variability will

55

reflect genuine lasting changes in the stock prize, for example a rise in price when a new product is

56

released. This part should be classified as process noise.

57

What is the best reward prediction method an organism could use when facing process noise and

58

observation noise? Similar problems occur in engineering, for example in the context of navigation.

59

There, a very versatile solution has been found. That solution, called the *Kalman filter* (3), is very widely

60

used—it even played a role in the moon landing (4). The Kalman filter describes how estimates of a

61

variable must be updated when new noisy observations of that variable become available. For certain

62

63    types of signals, it can be shown that the Kalman filter is indeed the ***optimal*** method for prediction in the

64    presence of noise. The method has proven useful not only in engineering, but also as a model of neural

65    and behavioral processes (5-8).

66    However, if one wants to use a Kalman filter to predict rewards, one runs into a problem: the Kalman

67    filter requires estimates of the magnitudes of both process noise and observation noise as parameters.

68    Where to take these values from? An organism might either use fixed (perhaps genetically determined)

69    values or estimate the values somehow. The former option bears a risk: if the world changes, the quality

70    of the organism's predictions might decline strongly. The latter option raises the next question: how is this

71    estimation done?

72    Solutions for this have been proposed. For example, Piray and Daw (9) present a model that tracks both

73    process noise and observation noise alongside reward, allowing for ***adaptive*** Kalman filtering. However,

74    their model (a variant of the particle filter) is targeted at the computational level, i.e., it is set up to

75    investigate how the simultaneous adaptation to two noise types of noise affects learning. Questions

76    concerning the underlying biological mechanisms remain largely unaddressed. The model of Piray and

77    Daw (9) is hence not suitable to describe biological learning on the mechanistic or the algorithmic level.

78    This leads us to the central question of this paper: ***how might organisms track observation noise in a***

79    ***biologically plausible, computationally simple way, and use it for adaptive reward prediction?*** We

80    propose that observation noise is tracked in the basal ganglia, and that it is used to improve learning

81    performance by normalizing reward prediction errors.  This proposal is based on the observation of

82    dopamine activity patterns consistent with normalized prediction errors (10), as well as on previous

83    suggestions that reward uncertainty might be represented in the basal ganglia (11, 12). Such scaling of

84    prediction errors is conceptually related to scaling mechanisms that occur in free energy models, as well

85    as to techniques such as adaptive momentum that are used to improve fitting algorithms (see Discussion

86    for details).

5

87    Below, we give a detailed analysis of how the basal ganglia circuit might carry out the computations

88    necessary to track and utilize reward observation noise. To provide some context for this analysis, we

89    now move on to a brief review of the main features of that part of the brain.

90    Fig 1 shows a highly simplified version of the cortico-basal-ganglia-thalamic circuit, with three important

91    brain regions—the cortex, the striatum, and the thalamus—arranged along the vertical axis.
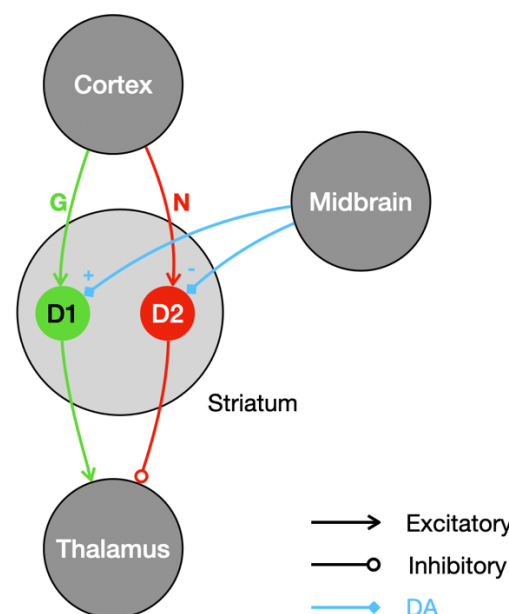
92



93

94    *Fig 1. The simplified basal ganglia circuit. Selected nuclei and connections are shown as circles and*

95    *arrows. Green connections correspond to the direct pathway; red connections correspond to the indirect*

96    *pathway. Dopamine projections are shown in blue.*

97

98    The striatum is the largest nucleus within the basal ganglia system. It includes two populations of medium

99    spiny projection neurons (SPNs): the D1 and the D2 population (D1 and D2 are the types of dopamine

100   receptors that the corresponding neurons express). This division of the striatum gives rise to two parallel

101   descending pathways called the direct/Go and the indirect/No-go pathway, shown in green and red

6

102 respectively in Fig 1 (13-15). The cortical inputs to these two pathways are modulated by the strengths of

103 the synapses between cortex and the striatal populations, collectively labeled $G$ and $N$ in Fig 1.

104 The thalamus receives the output of the basal ganglia circuit. The effects of the direct and the indirect

105 pathway on the thalamus are differential: the direct pathway effectively excites the thalamus; the indirect

106 pathway effectively inhibits it. Note that the projections from the striatum to the thalamus in Fig 1 are

107 abstractions—in the brain, there are several intermediate nuclei between the striatum and the thalamus.

108 The final key element of the basal ganglia system are the dopamine projections from midbrain regions

109 that target the striatal D1 and D2 populations. The effects of dopamine release on the striatal populations

110 are twofold: dopamine modulates activity, but also triggers plasticity. The direction of those effects

111 depends on the receptor type of the target neuron: dopamine increases excitability and potentiates

112 synapses in the D1 population while it decreases excitability and depresses synapses in the D2 population

113 (13-15).

114 Overall, we may view the basal ganglia circuit as two parallel descending pathways that converge on the

115 level of the thalamus, where they have opposite effects. Those pathways are differentially modulated by

116 dopamine, which also controls synaptic plasticity between the cortex and the striatum.

117 Concerning the function of the elements of this model of the basal ganglia, we follow a popular view

118 often used in modelling (11, 16, 17): the cortex supplies contextual information, i.e., cues, stimuli,

119 sensory data or information on the state of the environment; the other populations (D1, D2 and Thalamus)

120 encode actions. Each action is represented by a distinct subpopulation of each nucleus, and the

121 connectivity between the nuclei is action specific. For example, assume there is a subpopulation in D1

122 associated with pressing a lever. A corresponding subpopulation could be found in D2 as well as in the

123 thalamus, which is known to relay motor commands to the relevant cortical areas (18). The two striatal

124 subpopulations associated with the lever press would then project exclusively to the lever-press

125 subpopulation in the thalamus, together forming what is often called an action channel (19). Learning is

126   assumed to take place at the interface between the cortex and the striatum (which, in this model, can be

127   considered a state-action mapping). Learning is implemented through dopamine-mediated plasticity of

128   cortico-striatal synapses. These synapses within an action channel are assumed to store information on the

129   action value (the mean reward associated with the action). Action values determine the relative

130   activations of action channels (i.e., the difference in activation between the Go- and the No-go pathways),

131   and hence contribute to action selection at the level of the thalamus. It has been proposed (11) that

132   cortico-striatal synapses additionally encode reward uncertainty (in the sum of the weights in the Go- and

133   the No-go pathways), as we explain in detail below.

134   In the following sections, we present and analyze a model—the scaled prediction error (SPE) model—that

135   tracks observation noise and uses it for adaptive reward prediction. In the first part of the paper, we

136   introduce the model and test its performance. There, we show that it outperforms the classic Rescorla-

137   Wagner (RW) model of associative learning (20) which does not adapt to observation noise, using

138   simulations of a reward prediction task. In the second part of the paper, we discuss neural mechanisms

139   that might implement the SPE model in the basal ganglia circuit. We first focus on dopamine signals, and

140   then move on to the mechanisms behind tracking observation noise and scaling prediction error.

# Results

141

## *The model*

142

143 The SPE model is a model of reward prediction—it predicts the magnitude of the next reward based on

144 previous reward observations. It can be understood as approximate Bayesian inference with respect to a

145 particular model of the reward generation process. This model is given by

146
$$r_t \sim N(\mu, \sigma)$$

147                                                                                               Eq. 1

148 with $r_t$ the reward in trial $t$, $\mu$ the mean reward and $\sigma$ the reward observation noise. The SPE model can

149 be derived by approximating Bayesian inference of the parameters $\mu$ and $\sigma$ from the observed rewards

150 (we show this derivation in Appendix S1). It does this by maintaining estimates $m$ and $s$ of those

151 parameters, which it updates whenever a new reward is observed. Though the model is designed to infer

152 the mean and standard deviation of stationary reward processes, it can also be applied to reward processes

153 with a drifting mean (i.e., to processes with non-zero process noise). We show this below in our

154 simulations.

155 The SPE update rules are

156
$$\delta = \frac{r_t - m_{t-1}}{s_{t-1}}$$

157                                                                                               Eq. 2

158
$$m_t = m_{t-1} + \alpha_m \delta$$

159                                                                                               Eq. 3

160
$$s_t = s_{t-1} + \alpha_s(\delta^2 - 1).$$

161                                                                                               Eq. 4

9

162    In these equations $\delta$ is a scaled reward prediction error, while $\alpha_m$ and $\alpha_s$ denote the learning rates for the

163    mean and standard deviation, respectively. We will next show that the equations can indeed recover mean

164    reward and its standard deviation.

## *Fixed point analysis*

166    Do the SPE rules do what they are meant to do? Here, we use a stochastic fixed-point analysis to show

167    that in theory, $m$ and $s$ should converge to the mean and the standard deviation of the reward signal. Let

168    us assume that rewards are indeed generated by sampling from a distribution with mean $\mu$ and standard

169    deviation $\sigma$ (this could be a normal distribution or any other distribution with well defined mean and

170    standard deviation).

171    We consider a situation in which the learner has already found the correct values of the variables it

172    maintains, i.e., $m = \mu$ and $\sigma = s$. From there, what are the **expected updates**? A straightforward

173    calculation yields

174
$$E(\Delta m) = \alpha_m E\left(\frac{r - m}{s}\right) = \alpha_m\left(\frac{Er - \mu}{\sigma}\right) = \alpha_m\left(\frac{\mu - \mu}{\sigma}\right) = 0$$

175                                                                                          Eq. 5

176
$$E(\Delta s) = \alpha_s\left(\frac{E(r - m)^2}{s^2} - 1\right) = \alpha_s\left(\frac{E(r - \mu)^2}{\sigma^2} - 1\right) = \alpha_s\left(\frac{\sigma^2}{\sigma^2} - 1\right) = 0$$

177                                                                                          Eq. 6

178    with $Er = \mu$ and $E(r - \mu)^2 = \sigma^2$ by definition. We find that the expected change away from $(m, s)$

179     $= (\mu, \sigma)$ is zero, which makes $(\mu, \sigma)$ a stochastic fixed-point. We may conclude that in equilibrium, the

180    rules given in Eq. 2 – 4 should give us unbiased estimates of the reward mean and standard deviation.

181    It can be shown that the SPE model is related to the Kalman filter—it can be viewed as an approximation

182    to a steady-state Kalman filter, which becomes more accurate if observation noise dominates process

10

183    noise (see Appendix S2 for details). Finally, note that the RW model is a special case of the SPE model

184    (i.e., if $\alpha_S = 0$).

## *Performance*

186    How do the SPE rules compare to established rules such as the Rescorla-Wagner model with respect to

187    accurate reward predictions? To compare the performances of SPE learning and RW learning, we apply

188    both to a reward prediction task: sequences of rewards are generated according to

$$r_t \sim N(\mu_t, \sigma)$$

$$\mu_{t+1} \sim N(\mu_t, \nu)$$

191    In the above equation $\nu$ is the standard deviation of the process noise. Both learners observe the reward

192    signal and provide reward predictions at every trial. The learners' performance is judged by measuring the

193    average precision of their predictions.

194    The task is designed to challenge the learners with rewards that change over time, forcing them to

195    continuously learn. Note that the reward-generating process here is more complex than the generative

196    model from which SPE learning is derived. This is not a problem—the SPE model is robust with respect

197    to violation of its assumptions, as we shall see below. Of course, one could derive learning rules tailored

198    to this reward process. These would involve a representation of the process noise, or volatility. This is not

199    our goal here. Instead, we are interested in the SPE learning rules as a model of basal ganglia learning and

200    want to test their performance.

201    We compare the models for different levels of observation noise $\sigma$, while keeping the process noise $\nu$

202    constant at $\nu = 1$. The results of those comparisons are presented in Fig 2 (see Methods for details of the
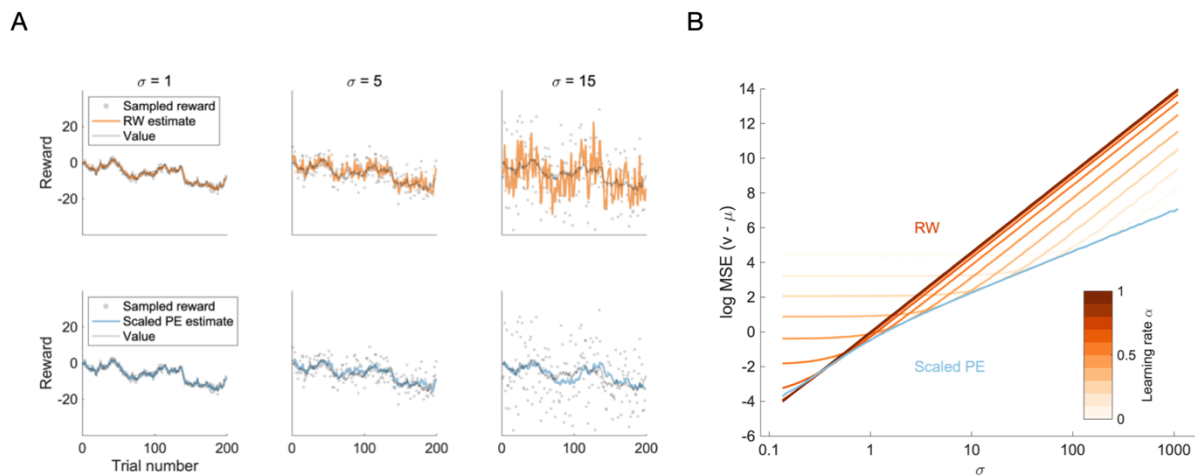
203    implementation).

**Fig 2. Reward prediction performance of the RW and SPE models.** *A The first 200 trials of reward prediction for the RW learner (upper row, orange color) and the SPE learner (lower row, blue color). The true value (grey line), the observed rewards (grey dots) and the learner's estimate (colored line) are shown as a function of trial number. Columns correspond to selected levels of observation noise ($\sigma = 1, 5, 15$). B Learning performance averaged over trials. We show the logarithm of the mean squared difference between the mean of the reward distribution and the learner's prediction thereof, as a function of the observation noise $\sigma$. Orange lines correspond to RW learners, the blue line corresponds to a SPE learner parametrized with $\alpha_m = 1$ and $\alpha_S = 0.01$. The different shades of orange correspond to different learning rates, as indicated by the color bar.*

Looking first at the time series in Fig 2A, we find that there is a qualitative difference between the RW learner in the top row and the SPE learner in the bottom row: as the noise level $\sigma$ increases, the RW learner's predictions increasingly fluctuate, since the reward prediction errors (and hence the updates) scale proportionally to the amplitude of the observation noise. This is not so for the SPE learner, whose predictions fluctuate as much for low noise levels as they do for high noise levels.

221     This effect is also visible in the aggregated performance measure, shown in Fig 2B: the mean squared

222     errors of the learners' predictions grow with observation noise for all learners, but they grow stronger for

223     the RW learners. We find a very stereotyped effect for the average performance of RW learners: as the

224     level of noise increases, prediction accuracy does not change much up to a certain point (call this the

225     plateau) and grows steadily after that point (call this the slope). This is the case irrespective of the

226     learning rate. Smaller learning rates have a plateau that extends to higher noise levels but also provides a

227     lower accuracy. The steepness of the slope is invariant across learning rates.

228     To gain an intuition for the shape of these curves, let us compare two different situations. First, consider

229     very low levels of observation noise. In this regime, reward observations are almost identical to

230     observation the underlying mean reward. Hence, if the observed reward changes, this mostly reflects a

231     genuine change of the underlying mean reward. To keep reward predictions precise, such changes should

232     be followed. However, for learning rates smaller than one, the RW model does not fully follow the

233     changes of the reward signal—we may call this underfitting, as the model ignores meaningful variation in

234     the signal. The resulting error dominates the performance. The magnitude of this error depends on the

235     volatility of the signal. Since the volatility is kept constant in the simulations in Fig 2B, we see a

236     performance plateau at low levels of observation noise.

237     Now, consider very high levels of observation noise. In this case, reward observations are very

238     inaccurate—an observation tells us very little about the underlying mean reward. Changes in observed

239     rewards mostly reflect the noisiness of the observations and should be ignored. However, for learning

240     rates larger than zero, the RW model does not fully ignore those fluctuations, but follows them. This can

241     be called overfitting, as the model tries to adapt to random fluctuations.

242     Overall, the behavior of the RW learners is such that for each given level of observation noise there is an

243     optimal learning rate: if one selects any one position on the x-axis of Fig 2A, there is always a single

244     orange curve with the lowest y-coordinate (and hence the smallest average error) at that position. In

245     general, we find: the higher the observation noise, the lower the optimal learning rate. This appears

13

246    consistent with intuition—if observation noise is high, there is less useful information in any single

247    observation and an organism should therefore update its estimate more carefully.

248    The SPE learner shows different behavior. There is also a slope (prediction accuracy steadily decreases

249    with increasing observation noise), but no plateau. The steepness of the slope changes at $\sigma = 1$. For

250    higher levels of observation noise, the slope of the SPE learner is shallower than those of the RW

251    learners. We find that for any given level of noise $\sigma$ larger than one, the performance of the SPE model is

252    about as good as the performance of the best RW model. This suggests that in the regime of high

253    observation noise we might view the SPE model as an RW learner that reaches optimal performance by

254    fine-tuning itself to the estimated level of observation noise.

255    Can one do better than this? In fact, one can show that the SPE model (parametrized with $\alpha_m = 1$) is

256    approximately optimal in the situation investigated here: for high levels of observation noise, SPE

257    learning approximates the steady-state Kalman filter (we show this in Appendix S2), which is

258    approximately optimal for the types of signals we use here.

259    As mentioned above, to use a Kalman filter, one needs to provide it with the correct values of $\sigma$ and $\nu$.

260    This is also true for the steady-state version of the Kalman filter, but it is not the case for the SPE model:

261    here one only needs to provide $\alpha_m$—which corresponds to $\nu$ (see Appendix S2), —but not $\sigma$, which the

262    model can track by itself. We can thus think of SPE learning as *adaptive* steady-state Kalman filtering.

263    However, to work optimally, the SPE model still needs to be provided with the correct value for $\alpha_m$. To

264    make the model more autonomous, one might extend it with a mechanism to track $\nu$ alongside $\sigma$, for

265    example the mechanism proposed by Piray and Daw (2). This is an interesting direction for further

266    research but goes beyond the scope of this work.

267    In summary, we find that the SPE model is approximately optimal for signals with $\nu < \sigma$. In particular, it

268    will be at least as good as any RW learner, and about as good as a steady-state Kalman filter. SPE learners

269    thus appear particularly well suited to track signals with unknown or changing levels of observation noise,

270     as they can adapt themselves to whatever level of noise they experience. In contrast, an RW learner would

271     either have to be fine-tuned based on prior knowledge, or it would perform suboptimally due to under- or

272     overfitting.

## *The neural implementation of SPE learning*

273

274     Could the SPE learning rules be implemented in the dopamine system and the basal ganglia pathways? In

275     this section, we propose a possible mechanism. We suggest that striatal dopamine release broadcasts

276     scaled prediction errors, $\delta = \frac{r-m}{s}$, and that the update rules given in equations Eq. 3 and Eq. 4 are

277     implemented by dopamine-dependent plasticity in the striatum. In the next subsections, we will analyze

278     the plausibility of these suggestions. First, we discuss the relationship between dopamine responses and

279     scaled prediction errors. Then we discuss how the SPE learning rules can be mapped on striatal plasticity

280     rules. Finally, we propose a mechanism that might implement the scaling.

### *Scaled prediction errors are consistent with dopamine activity*

281

282     In a seminal study, Tobler et al. (10) investigated how the responses of dopamine neurons to

283     unpredictable rewards depended on reward magnitude, using electrophysiology in monkeys. Three

284     different visual stimuli were paired with three different reward magnitudes (0.05 ml, 0.15 ml and 0.5 ml

285     of juice). After being shown one of the stimuli, the monkeys received the corresponding reward with a

286     probability of 50%. Seeing the stimulus allowed the monkey to predict the magnitude of the reward that

287     could occur, but not whether it would occur on a given trial. Reward delivery thus came as a surprise and

288     evoked a dopamine response. Interestingly, these responses did not scale with the magnitude of the

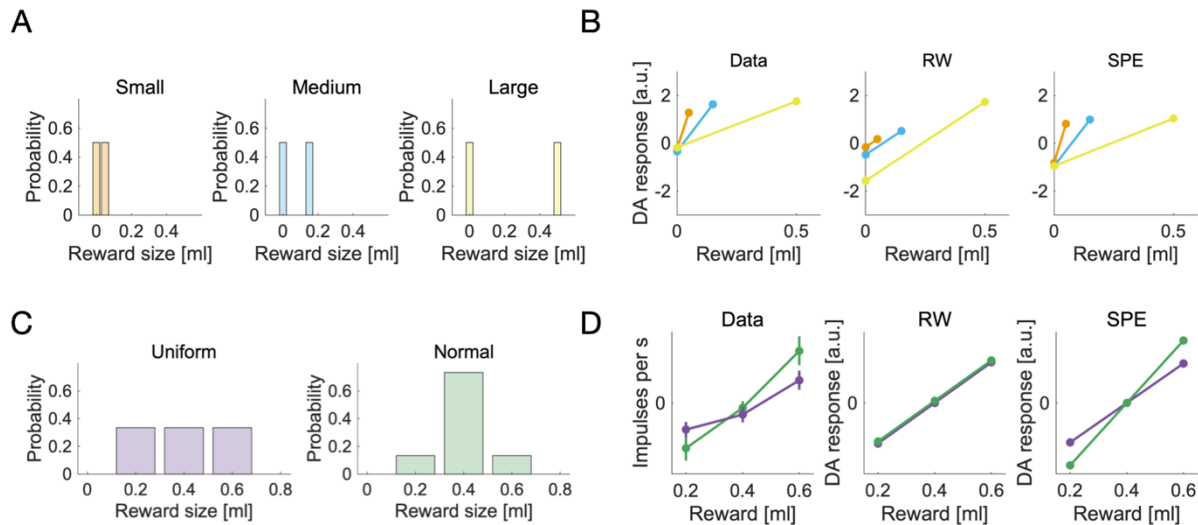289     received rewards. The measured dopamine responses are shown in Fig 3B.

15

**Fig 3. Dopamine responses to unpredictable rewards—experimental data and simulations. *A** The*

*reward distributions used by Tobler et al. (10). Each distribution corresponds to an experimental*

*condition. **B** Dopamine responses to rewards sampled from the distributions in A are shown as a function*

*of reward magnitude, for the three different conditions. The representation of data is similar to that in*

*figure 4C of Tobler et al. (10). We show experimental data, extracted from figure 4C (animal A) of Tobler*

*et al. (10) and simulated data, using a standard RW model and the SPE model. The colors relate the*

*dopamine responses in B to the reward distributions in A. **C** The reward distributions used by*

*Rothenhoefer et al. (21). The panel is reproduced from Rothenhoefer et al. (21), figure 1A. **D** Dopamine*

*responses to rewards sampled from the distributions in C. We show the empirical values, reproduced*

*from Rothenhoefer et al. (21), figure 2E, and the responses according to the RW model computed*

*analytically as $\delta = r - \mu$, and the SPE model computed as $\delta = \frac{r - \mu}{\sigma}$, where $\mu$ and $\sigma$ are the mean and*

*standard deviation of corresponding reward distributions in C. Purple lines correspond to the uniform*

*reward distribution, green lines correspond to the normal reward distribution.*

This result was unexpected—standard RW learning would predict that the residual prediction errors in

rewarded trials should grow linearly with reward magnitude. Our new SPE rules, on the other hand,

predict exactly what has been observed. See Fig 3B for simulated and experimental DA responses.

16

307     One may object that the results of Tobler et al. (10) might also be explained by scaling with respect to the

308     reward range—reward range and reward standard deviation cannot be dissociated in that experiment.

309     While that is true, another recent experiment can dissociate them: Rothenhoefer et al. (21) used two

310     reward distributions with the same reward range but different reward standard deviations in a Pavlovian

311     conditioning task (see Fig 3C).

312     After exhaustive training, single unit recordings were performed to measure dopamine responses to

313     rewards that deviated from the expected value. It was found that the same deviation from the expected

314     value caused stronger dopamine responses for the distribution with the smaller standard deviation (Fig

315     3D, first panel). This is consistent with scaling by reward standard deviation, but not with scaling by

316     reward range---both distributions had the same range, so scaling by range should yield similar responses

317     for both conditions. These experimental data cannot be accounted for by the RW model (Fig 3D, second

318     panel), but can be reproduced by the SPE model (Fig 3D, third panel).

319     *The SPE learning rules are consistent with striatal plasticity*

320     After establishing that dopaminergic scaled prediction errors are plausible, we now move on to discuss

321     how the update rules given in Eq. 3 and Eq. 4 could be implemented in the basal ganglia circuit.

322     Mikhael and Bogacz (11) proposed a distributed encoding of the two reward statistics ($m$ and $s$) in the

323     two main basal ganglia pathways: in their model, the mean of the reward signal is encoded in the

324     difference between synaptic inputs to striatal neurons in direct (Go) and indirect (NoGo) pathways,

325     whereas the standard deviation of the signal is encoded in the sum of these inputs. Formally, we write

$$m = \frac{1}{2}(G - N)$$

327                                                                                                             Eq. 7

$$\lambda s = \frac{1}{2}(G + N)$$

329        Eq. 8

330     In Eq. 7 and 8, $G$ and $N$ denote the synaptic inputs in the direct and indirect pathway respectively (22),

331     and $\lambda$ is a coefficient determining the accuracy with which the standard deviation can be encoded (as

332     explained below). These assumptions can be used to rewrite the learning rules given in Eq. 3 and 4 in

333     terms of $G$ and $N$. In particular, note that by combining Eq. 7 and 8, we see that $G = m + \lambda s$ and

334     $N = \lambda s - m$. Therefore, we can derive the update rules for $G$ (or $N$) by adding (or subtracting) Eq. 3 and

335     4, and obtain

336 $$\Delta G = \alpha_m f_\beta(\delta) - \lambda \alpha_s$$

337        Eq. 9

338 $$\Delta N = \alpha_m f_\beta(-\delta) - \lambda \alpha_s$$

339        Eq. 10

340     with $f_\beta(\delta) = \beta \lambda \delta^2 + \delta$ and $\beta = \alpha_s/\alpha_m$. It is worth emphasizing that Eq. 9 and 10 are equivalent to Eq. 3

341     and 4 (because they are just rewritten in terms of different variables). Therefore, the model described by

342     Eq. 9 and 10 estimates exactly the same mean and variances as a model described by Eq. 3 and 4, and

343     hence it produces identical performance in Fig 2 and dopaminergic responses in Fig 3.

344     One important issue while considering biological plausibility of the model is the fact that the synaptic

345     weights on the indirect pathway $N$ cannot be negative, while the model assumes that these weights encode

346     $N = \lambda s - m$. Imposing a constraint of $N$ being non-negative will limit the ability of the network to

347     accurately estimate standard deviation of rewards to cases when it is sufficiently high (i.e. $\sigma \geq \mu/\lambda$).

348     Hence the parameter $\lambda$ controls the accuracy with which the standard deviation can be estimated.

349     However, according to Eq. 8 there is a cost of high accuracy, because a high value of $\lambda$ will result in

350     overall larger values of the synaptic weights (analogously as in the model of Mikhael and Bogacz (11)),

351     and hence higher metabolic cost of the computations.

352  Eq. 9 and 10 show three main features: 1) different overall effects of dopamine on plasticity in each

353  pathway, 2) nonlinear effects of dopaminergic prediction errors represented by the transformations $f_\beta$ and

354  3) synaptic unlearning represented by decay terms. We will discuss the experimental data supporting the

355  presence of these features in turn.

356  First, the efficacy of direct pathway synapses is assumed to increase as a result of positive reward

357  prediction errors (i.e., $\delta > 0$), and decrease as a result of negative reward prediction errors (i.e., $\delta < 0$).

358  The opposite is assumed to hold for indirect pathway synapses: their efficacy should decrease with

359  positive prediction errors and increase with negative prediction errors. This premise corresponds to the

360  sign of the prediction error in Eq. 9 and 10, and it is consistent with data obtained in experiments (23).

361  Second, it is assumed that for striatal neurons in the direct pathway, positive prediction errors have a

362  stronger effect on plasticity than negative prediction errors. This assumption is expressed in the shape of

363  the function $f_\beta(\delta)$, which is plotted in Figure 4Ai. Note that the slope for positive $\delta$ is steeper than for

364  negative $\delta$, implying that positive prediction errors should lead to bigger changes in $G$ than negative

365  prediction errors. The computational role of this nonlinearity is to filter the reward prediction errors: it

366  amplifies the positive components while dampening the negative components of the signal.

367  For striatal neurons in the indirect pathway, the SPE model assumes the opposite: negative prediction

368  errors should have a stronger plasticity effect than positive prediction errors, because the weight

369  modification is proportional to $f_\beta(-\delta)$, which is plotted in Figure 4Aii. Mikhael and Bogacz (11) argue

370  that this premise is realistic, based on the different affinities of the D1 and D2 receptors that are present in

371  striatal neurons in the direct and indirect pathways respectively: while D1 receptors are mostly

372  unoccupied at baseline dopamine levels, D2 receptors are almost saturated—this is visualized in Fig 4B.

373  Due to this baseline setting, additional dopamine should lead to a large difference in the occupation of D1

374  receptors and hence affect the neurons on the direct pathway, but only a small change in the occupancy of

19

375    D2 receptors thus little influence the neurons on the indirect pathway. A decrease in dopamine, on the

376    other hand, is strongly felt in D2 receptor occupancy but does not change D1 receptor occupancy much.
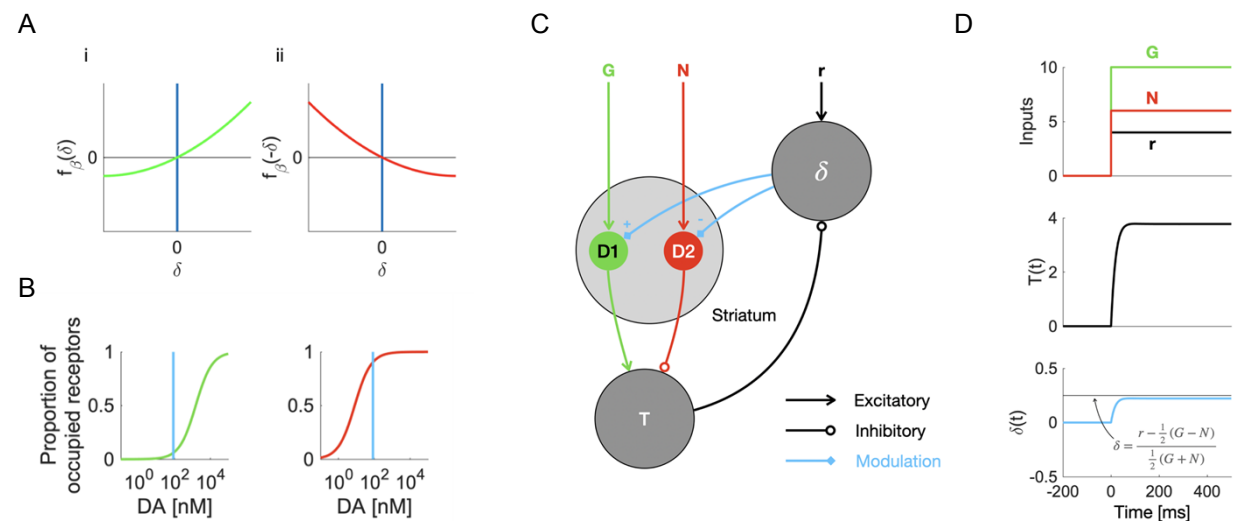
377

378



379    *Fig 4. Plasticity and computations in the basal ganglia circuit. A The nonlinear transformation of*

380    *dopaminergic prediction errors in the SPE model. The transformation in the direct pathway (i) and the*

381    *transformation in the undirect pathway (ii) are mirror images of each other. **B** We plot the proportion of*

382    *occupied receptors in the striatum as a function of dopamine concentration. The curves are based on the*

383    *results of Dreyer et al. (25). The blue vertical lines indicate the baseline dopamine concentration in the*

384    *ventral striatum, based on the results of Dodson et al. (26). The green curve corresponds to the*

385    *occupancy of D1 receptors, the red curve corresponds to the occupancy of D2 receptors. Panel B is a*

386    *partial reproduction of figure 3D of Möller and Bogacz (27). **C** The connectivity underlying a dynamical*

387    *model of the simplified basal ganglia circuit. Circles correspond to neural populations; arrows between*

388    *them indicate connections. **D** The computation of a scaled prediction error in continuous time, according*

389    *to a dynamical model of the basal ganglia. We show how the relevant variables, T and δ, evolve as a*

390    *function of time, assuming a step-function activation for the input nodes G, N and r. The black line in the*

391    *lowest panel indicates the level of dopamine required for exact SPE learning.*

392    Third, an activity-dependent decay (or 'unlearning') is assumed to occur in the synaptic weights

393    whenever they are activated in the absence of prediction errors. This is reflected in terms $-\lambda\alpha_s$ in Eq. 9

394    and 10. On the neural level, that premise translates into mild long-term depression after co-activation of

395    the pre- and postsynaptic cells at baseline dopamine levels. Recently, this effect has been observed at

396    cortico-striatal synapses in vivo (24): in anaesthetized rats, presynaptic activity followed by postsynaptic

397    activity caused LTD at baseline dopamine levels (i.e. in the absence of dopamine-evoking stimuli).

398    In summary, we discussed the three premises of the learning rules—the different overall effects of

399    dopamine on plasticity in each pathway, the nonlinear effects of dopaminergic prediction errors and

400    synaptic unlearning. We saw that all three premises are supported by the physiological properties of

401    striatal neurons on the direct and indirect pathway.

402    *Scaled prediction errors and the basal ganglia circuit*

403    Next, we discuss how the scaled prediction error $\delta = \frac{r-m}{s}$ might be computed in the basal ganglia

404    system. Expressed in terms of $G$ and $N$, the scaled prediction error is given as

405
$$\delta = \frac{r - \frac{1}{2}(G-N)}{\frac{1}{2\lambda}(G+N)}$$

406                                                                                                              Eq. 11

407    This seems to be a complicated combination of terms, and it is difficult to see how a simple network

408    might compute it. Surprisingly, there is a simple approximate implementation based on a feedback loop.

409    Here, we will describe that mechanism, using a minimal dynamical model of the basal ganglia network.

410    First, where is the feedback loop? Let us assume that the prediction error $\delta$ is encoded in the activity of a

411    population of dopaminergic neurons. This population receives inhibitory input $T$ from the thalamus, and

412    excitatory input $r$ that encodes a reward signal. Formally, we assume $\delta = r - T$. We follow Möller and

21

413     Bogacz (27) in assuming that the thalamic activity reflects the total output from the basal ganglia—the

414     difference between the activity in the direct and indirect pathways—which is captured by

415

$$T = \frac{1 + \delta/\lambda}{2} G - \frac{1 - \delta/\lambda}{2} N$$

416           Eq. 12

417     The first term $\frac{1 + \delta/\lambda}{2} G$ corresponds to the activity in the direct pathway, which is proportional to synaptic

418     input $G$, and is increased by the dopaminergic modulation, because the gain of striatal neurons in the

419     direct pathway is enhanced by dopamine. The second term $\frac{1 - \delta/\lambda}{2} N$ corresponds to the activity in the

420     indirect pathway, which is attenuated by dopamine, because the gain of striatal neurons in the indirect

421     pathway is reduced by dopamine (13, 14). The proposed model contains a feedback loop: dopamine

422     release modulates the thalamic activity, which itself inhibits dopamine release.

423     To examine the computation of scaled prediction errors, we model the relevant populations' activities as

424     leaky integrators with effective connectivity as sketched in Fig 4C, using differential equations in

425     continuous time.

426     The dynamical system sketched in Fig 4C corresponds to a set of differential equations,

427

$$\tau_\delta \dot{\delta} = -\delta + (r - T)$$

428           Eq. 13

429

$$\tau_T \dot{T} = -T + \frac{1 + \delta/\lambda}{2} G - \frac{1 - \delta/\lambda}{2} N$$

430           Eq. 14

431     Here, $\tau_\delta$ and $\tau_T$ are the characteristic timescales of the striatal dopamine release and thalamic activation.

432     The system is set up such that its equilibrium point is consistent with our trial-wise description ($\delta = r - T$

433     and $T = \frac{1 + \delta/\lambda}{2} G - \frac{1 - \delta/\lambda}{2} N$ at $\dot{\delta} = \dot{T} = 0$). This asserts that the two levels of description are consistent

22

434    with each other. Using these equilibrium equations, we can determine the equilibrium value of $\delta$ (by

435    inserting one equation into the other and solving for $\delta$). We find

436

$$\delta = \frac{r - \frac{1}{2}(G - N)}{1 + \frac{1}{2\lambda}(G + N)}$$

437                                                                         Eq. 15

438    For $\frac{1}{2\lambda}(G + N) \gg 1$, this approximates the scaled prediction error in Eq. 11. This suggests that the circuit

439    can compute an approximation to the scaled prediction error. The approximation will be accurate as long

440    as $s$ is sufficiently large (recall from Eq. 8 that $\frac{1}{2\lambda}(G + N) = s$). Although the additional term 1 in the

441    denominator prevents perfect scaling, it might in fact be beneficial: it could prevent catastrophically large

442    prediction errors that might cause the instabilities when the denominator becomes very small.

443    So far, it looks as though the circuit has an equilibrium point at approximately the right value. However, it

444    is not yet clear whether and how this equilibrium is reached. To learn more about these aspects, we need

445    to simulate the system. To simulate the computation of the prediction error, we assume $G$, $N$ and $r$ to be

446    provided externally, for example through cortical inputs. $G$ and $N$ then represent precisely timed reward

447    predictions, while $r$ represents the reward signal itself. We model $G$, $N$ and $r$ as step-functions that jump

448    from zero to their respective values at the same time, as illustrated by the first panel of Fig 4D. The time

449    constants $\tau_\delta$ and $\tau_T$ are set to realistic values taken from the literature (see Methods for details). A

450    simulation of the system is shown in Fig 4D. We find that $\delta$ settles to its equilibrium value quite quickly

451    (after tens of milliseconds) and without oscillations. This is likely due to the difference in time

452    constants—the thalamic activity changes much faster than the striatal dopamine concentration. Our results

453    suggests that even a simple system as the one in Fig 4C can compute scaled prediction errors through a

454    feedback loop.

455

23

# Discussion

456

457 Above, we presented a new model of error-driven learning: the SPE model. We tested it in simulations

458 and compared it with neural data. Now, we will discuss the new model more broadly. First, we will

459 summarize our key findings. Then, we will present several testable predictions that follow from the

460 model. Finally, we will discuss how the SPE model relates to other models from neuroscience and from

461 artificial intelligence.

462 ## *Summary*

463 This work introduces the SPE model, which describes how an organism might adapt its learning

464 mechanism to changing levels of reward observation noise $\sigma$. First, we proposed the SPE learning rules,

465 which can track the mean and standard deviation of a reward signal. We then tested the performance of

466 the new rules. Comparing SPE learning with RW learning, we found that the new learning rules can

467 improve performance when a learner faces unknown or varying levels of reward observation noise. Next,

468 we reviewed empirical evidence relating to SPE learning. On the neural level, we found that SPE learning

469 describes dopamine responses better than conventional models in several studies. We further showed how

470 the basal ganglia pathways might implement the learning rules of the SPE model, and how scaled

471 prediction errors could be computed in a dopaminergic feedback loop.

472 ## *Experimental predictions*

473 Our model makes several predictions on different levels of analysis. First, SPE learning can be

474 distinguished from other types of learning on the level of behavior. This is because according to SPE, the

475 learning rate (and hence the speed of learning) should depend on the stochasticity of the reinforcement

476 that drives learning. From this, different predictions follow.

477 For example, reward magnitude should ***not*** affect instrumental learning speed if animals were exposed to

478 the reward prior to the instrumental phase. This is because SPE learning would allow the animal to

479 normalize the reward magnitude through the adaptive scaling of prediction errors, as in the experiment by

24

480 Tobler et al. (10). The dopamine-guided instrumental learning through normalized prediction errors

481 would then not be affected by the overall scale of the rewards. Concretely, this could be tested in a

482 decision-making task with two conditions. In one condition, correct choices are reinforced with high

483 rewards (for example 0.6 ml of juice). In the other condition, low rewards are provided, for example 0.2

484 ml of juice. Conditions must be cued, perhaps by two different visual stimuli that precede the decision.

485 According to the RW model, we expect to see a steeper learning curve in the high reward condition, as in

486 the experiment of (28). However, SPE theory predicts that this difference between conditions should

487 vanish if an appropriate pretraining is applied. For example, the cues could be associated with the

488 different reward sizes through Pavlovian conditioning. According to SPE theory, the pretraining should

489 establish a condition-specific normalization cued by the stimuli. This normalization should then lead to

490 normalized learning curves in the instrumental phase. Ultimately, the difference between the learning

491 curves should vanish.

492 Furthermore, SPE learning predicts that learning rates should change if reward stochasticity is changed.

493 This could be tested by having participants track and predict a drifting reward signal which shows

494 different levels of stochasticity at different times. If participants use SPE, the learning rate should

495 decrease with increasing stochasticity, which would lead to invariant update magnitudes. On the other

496 hand, if participants do not use SPE, increasing reward stochasticity would not affect the learning rate,

497 and hence lead to larger updates.

498 On the neural level, SPE learning predicts trial-by-trial changes of how the dopaminergic prediction error

499 is normalized to a given reward signal. Neural recordings from the relevant brain areas during the

500 learning phase could be compared with simulations of the SPE model to test the theory. In particular, SPE

501 predicts that if reward stochasticity increases slowly, the scale of the corresponding dopamine bursts

502 should stay invariant. Standard theory, on the other hand, would predict that the scale of dopamine bursts

503 grows proportional to the scale of rewards, as prediction errors are a linear function of rewards.

504    Since the SPE model is closely related to the AU model (11), it inherits a prediction on how the activity in

505    the basal ganglia pathways should depend on reward stochasticity: if reward stochasticity is high, the sum

506    of the activity in the direct and the indirect pathway should be high—if reward stochasticity is low, the

507    sum should be low as well. It has been indeed observed that the neural activity in striatum increases with

508    reward uncertainty (29, 30). Cell-type specific imaging techniques such as photometry could be used to

509    further assert whether the uncertainty is encoded in the sum of activity of striatal neurons on the direct

510    and indirect pathways.

## *Relation to models in neuroscience*
511

### *The AU model*
512

513    The SPE model is closely related to the AU model of Mikhael and Bogacz (11)—both models describe

514    how the basal ganglia pathways track reward uncertainty; they also share the distributed encoding of

515    reward statistics. It is thus not surprising that the learning rules of the two models have similarities.

516    However, the SPE model differs from the AU model in several important aspects.

517    Of course, the scaled prediction error itself is the key new feature that drives most of the interesting

518    effects we investigated in this work. It is through the scaling of the prediction error that our new model

519    puts its estimate of the reward observation noise to good use. The AU model tracks reward noise as well

520    but does not use its estimate to improve learning performance (or for anything else). In contrast, the SPE

521    model explains not only **how** to track $\sigma$, but also **why**.

522    Further, the AU model assumes that there are two separate dopamine signals that modulate activity and

523    plasticity of striatal neurons, namely that the tonic level modulates activity, while the phasic bursts trigger

524    plasticity. However, it has been recently demonstrated that even a brief, burst-like activation of

525    dopaminergic neurons changes the activity levels of striatal neurons (31). Additionally, it has been shown

526    that reward prediction errors modulate the tendency to make risky choices (22), and risk attitudes are

527    known to depend on the balance between the direct and indirect pathways (32, 33). In this paper, we

26

528  demonstrated that a more realistic assumption, that the dopamine signal encoding prediction error also

529  changes the activity levels in striatum, enables scaling of prediction errors by uncertainty.

530  *The Kalman-TD model*

531  SPE learning is not the only model that addresses the scaling of dopamine responses. One recent theory—

532  Kalman-TD—explained those responses, as well as other phenomena such as preconditioning, as a

533  consequence of volatility tracking (5). Kalman-TD applies the Kalman filter method to the computational

534  problem of TD learning: reward prediction in the time domain. The resulting model features vector-

535  valued learning rates that constantly adapt to observations and outcomes. It elegantly describes how

536  covariances between cues and cue-specific uncertainties might modulate learning and can be shown to

537  explain several empirical phenomena. However, the Kalman-TD theory does not address the tracking of

538  observation noise (the theory focuses on process noise). It also does not discuss how prediction error

539  scaling might be implemented. We may thus view it as a complement rather than a competitor to the

540  theory presented above.

541  *The reward taxis model*

542  Another model was recently proposed to explain the effects reported by Tobler et al. (10) and other

543  phenomena. The model is called ***reward taxis*** (34), and explains the dopaminergic range adaptation using

544  a logarithm: if both rewards and reward expectations were transformed by a logarithmic function,

545  prediction errors would be given by $\delta = \log r - \log m = \log \frac{r}{m}$. In the experiment of Tobler et al. (10)

546  rewards were given in 50 % of the trials. For a reward of size $r$, the expected reward would then be $m = \frac{r}{2}$

547  , and the prediction error would be $\delta = \log \frac{r}{\frac{r}{2}} = \log 2$, i.e., independent of reward size. Reward taxis can

548  hence explain the results of Tobler et al. (10) quite elegantly.

549  However, that explanation breaks down as we look at other experiments. We have already mentioned the

550  experiment by Rothenhoefer et al. (21), which featured two reward distributions with equal means and

551    ranges but different standard deviations. We show those distributions in Fig 3C. Rothenhoefer et al. (21)

552    first used Pavlovian conditioning in a way similar way to Tobler et al. (10), pairing the two reward

553    distributions with two different cues. They then recorded the dopamine responses at reward delivery, for

554    all reward sizes of each distribution. We reproduce their data in Fig 3D (first panel). The responses to the

555    middle reward are similar for both distributions, but the responses to the extreme rewards differ: they

556    seem scaled up for the normal distribution.

557    What would the reward taxis theory predict for the responses in this experiment? Both distributions have

558    the same mean; reward taxis hence predicts similar responses for both distributions. The experimental

559    data thus falsifies the reward taxis model in this experiment. In contrast, the SPE model predicts different

560    responses for the two distributions—we show this in Fig 3D (last panel). Overall, it appears as if

561    dopamine responses to reward distributions with variable width are better captured by the SPE model than

562    the reward taxis model.

563    *Free energy models*

564    Finally, we want to discuss the relation of our model to free-energy models: the scaled reward prediction

565    errors in this work are formally related to the precision weighted prediction errors of the free-energy

566    approach, especially when the recognition density (the learner's model of the world) is taken to be

567    Gaussian (35-37). In that case, the prediction errors that drive inference and learning in free energy

568    models are often weighted by precisions, i.e., inverse variances. The connection to scaled reward

569    prediction errors becomes very close when the free energy approach is applied to reward prediction,

570    dopamine and the basal ganglia system, as has been done in the DopAct framework (38). This framework

571    integrates several theoretical ideas (free energy, reinforcement learning, habits without values and active

572    inference), and suggests that dopaminergic prediction errors drive both learning and action planning.

573    Precision weighted prediction errors encoded by dopamine transients feature in one variant of that model,

574    but they are not the focus of the theory, and possible implementations or empirical consequences of these

575    weighted prediction errors have not been investigated so far. Furthermore, it is important to note that

28

576    precision, or inverse variance, scales differently to standard deviation, and might hence not explain

577    classical observations such as those reported by Tobler et al. (10).

### *Relation to models in artificial intelligence*

579    Scaled reward prediction errors have been explored outside of neuroscience as well: in the field of AI-

580    type reinforcement learning, it was noticed that normalizing reward prediction errors can enable an agent

581    to learn effectively across several different tasks (39). This is consistent with our conclusions: different

582    tasks come with different levels of reward observation noise, and adaptive scaling can normalize

583    performance across tasks without requiring the need for fine-tuning. However, the rules for scaling

584    prediction errors in AI are different from the SPE learning rules and have not been designed with the

585    intention to model learning in biological systems. Further, Hessel et al. (39) have focused on typical

586    benchmark tasks of AI-type reinforcement learning (i.e. Atari games and others), while we have explored

587    the types of tasks that are used in neuroscience and psychology.

588    Prediction error scaling also occurs at a more basic level of AI, inside the optimization algorithms that are

589    used to improve the parameters of neural networks. A very prominent example is the Adam optimizer

590    (40), which implements a variant of gradient descent in which all updates are normalized using an

591    estimate of the second moment of the gradient distribution. By making gradient descent effective across

592    different gradient magnitudes, adaptive optimizers such as Adam contribute to the spectacular successes

593    of deep learning. This supports the main idea of this work—that scaling prediction errors can be

594    beneficial for learning. However, here we only looked at the scaling of *reward* prediction errors. Adam-

595    style optimization in machine learning, as well as free-energy models in computational neuroscience

596    suggest that there might be similar mechanisms for other neural error signals as well. Therefore, scaling

597    of prediction errors may be a fundamental and common mechanism in the brain. While the mechanisms

598    and evidence presented in this work focus on reward prediction errors and the basal ganglia system, it

599    would also be an interesting direction for future work to investigate scaled prediction errors in other

600    systems within the brain.

601 # Methods

602 ## *Reward prediction performance*

603 In Fig 2, we compare the performance of the RW model with the performance of the SPE model. The

604 SPE model was defined by the learning rules given in Eq. 2 – 4. The RW model was defined by

605
$$\delta = r - m_t$$

606 Eq. 16

607
$$m_t = m_{t-1} + \alpha\delta$$

608 Eq. 17

609 With a constant learning rate $\alpha$.

610 For Fig 2A, the RW model was parametrized with $\alpha = 0.5$ and the SPE learning rules were parametrized

611 with $\alpha_m = 1$ and $\alpha_s = 0.1$. Rewards were sampled from a normal distribution with drifting mean. The

612 process noise was fixed at $\nu = 1$. We used three different observation noise levels (1, 5 and 15).

613 For Fig 2B, we used 10 different learning rates for the RW model (ranging from 0.007 to 0.993), the same

614 parameters as above for the SPE model, and 100 different levels of observation noise, evenly distributed

615 on a logarithmic scale from 0.1353 to 1096.6. The process noise was fixed at $\nu = 1$ as above. For each

616 combination, we simulated $10^5$ trials and computed the average squared difference between the model

617 predictions $m$ and the true mean $\mu$ across all trials.

618 ## *Simulations of the task of Tobler et al.*

619 To simulate the relevant parts of the experiment reported by Tobler et al. (10), we modelled Pavlovian

620 conditioning with three different stimuli, which were associated with three different reward magnitudes (

621 $r = 0.05, r = 0.15, r = 0.5$). The stimuli were followed by the associated reward in one half of the trials

622 and by no reward in the other half. The rewarded trials were selected pseudorandomly, such that there

623 were two rewarded and two non-rewarded trials in every four successive trials.

624 We simulated 2000 trials per stimulus and extracted prediction errors from the last 1500. Discarding the

625 first 500 trials accounts for the substantial pretraining of Tobler et al. (10).

626 We used two models: an RW model and a SPE model. The learning rules of the RW model are given in

627 Eq. 16 and Eq. 17. The rules were used with $\alpha_m = 0.0067$ and $m_0 = 0$. The learning rules of the SPE

628 model are given in Eq. $2 - 4$. These rules were used with $\alpha_m = \alpha_s = 0.0067$ and $m_0 = 0$, $s_0 = 1$.

629 To compare our simulations to the experimental data from Tobler et al. (10), we extracted the prediction

630 errors $\delta$ from the simulations and averaged them for each model, outcome and condition separately. There

631 were three conditions (corresponding to the three reward sizes) with two outcomes (reward or no nothing)

632 each, resulting in a total of six combinations per model. Finally, we normalized the six averaged

633 prediction errors by their standard deviation for each model.

634 ### *A dynamical model of the basal ganglia*

635 The differential equations Eq. 11 and Eq. 12 were solved using MATLAB's ***ode15s***, from $t = -200\ ms$

636 until $t = 500\ ms$. As inputs, we used step functions

637
$$r(t) = \theta(t)r_{step}$$

638 Eq. 18

639
$$G(t) = \theta(t)G_{step}$$

640 Eq. 19

641
$$N(t) = \theta(t)N_{step}$$

642 Eq. 20

31

643 with $\theta(t) = 1$ for $t > 0$ and $\theta(t) = 0$ for $t < 0$, and $G_{step} = 10$, $N_{step} = 6$, and $r_{step} = 4$. These inputs

644 correspond to a learned mean $m = 2$ and a learned standard deviation $s = 8$ for $\lambda = 1$.

645 The time constant $\tau_T$ of the thalamic population was set to 10 ms, based on the measurement of the

646 membrane time constant of thalamic neurons reported by Paz et al. (41). The time constant $\tau_\delta$ for striatal

647 dopamine was set to 300 ms, based on figure 2C of Montague et al. (42): the dopamine transient in that

648 figure decays to $\exp -1$ of its peak value in about 300 ms.

# 649 Supporting information captions

650 Appendix S1. Derivation from Bayesian learning

651 Fig S1. The mode-matching method

652 Appendix S2. The high noise limit of the steady state Kalman filter

653 Fig S2. The learning rate of the steady state Kalman filter

654

# Supporting information
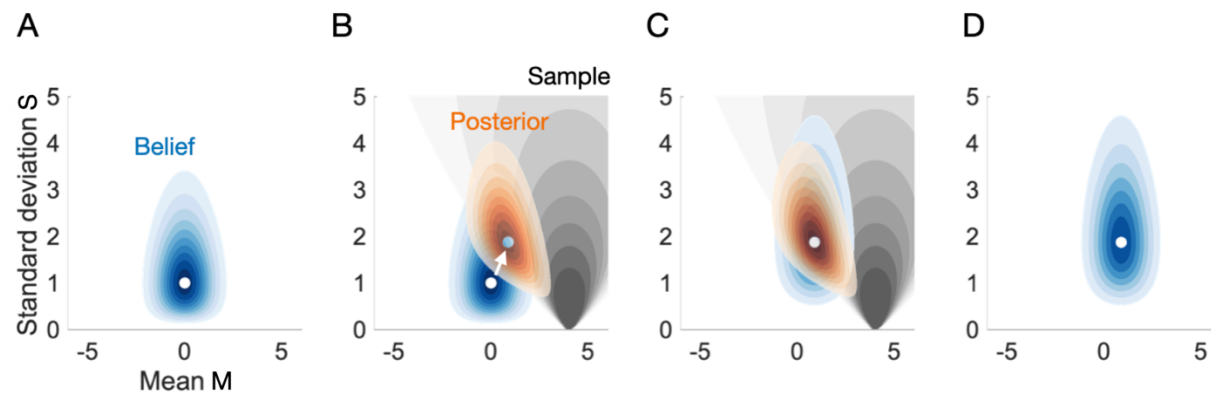
## *Appendix S1. Derivation from Bayesian learning*

One way to derive the scaled prediction error learning rules is the Bayesian mode-matching method, which is a novel (as far as we know) method to approximate Bayesian learning. We first introduce this method. We then apply the method to the problem of tracking the mean and standard deviation of a signal and thus find a new set of learning rules.

### *The mode-matching method*

The mode-matching method is based on Bayesian principles. Let us consider the problem of learning the mean and standard deviation of a signal. A fully Bayesian learner would always maintain a belief about the values of the mean and standard deviation, encoded as a probability distribution over all possible pairs of values. It would also maintain a generative model of the signal. When the learner is provided with new information (say another sample of the signal), it applies Bayes' law to combine its current belief (now the prior) and the likelihood of the observation (computed using the generative model) into a posterior distribution, which encodes its belief after observing the sample. This process is then repeated ad infinitum, with the posterior after one sample turning into the prior for the next.

Now, consider a learner that cannot encode arbitrary belief distributions. Instead, it can only adapt a few of the parameters of a belief distribution with otherwise fixed shape. For example, it might encode a belief using a normal distribution with fixed width and update it by adapting the mean. How might such a learner—let us call it a fixed-shape learner—approximate a fully Bayesian learner best?

Here, we propose the mode-matching method: after observing a new sample, the fixed-shape learner should change the parameters of its belief distribution such that the maximum of the distribution (its mode) is aligned with the maximum of the true posterior. We show this process schematically in Fig S1.

677

*Fig S1. The mode-matching method. The fixed-shape belief distribution is represented by a blue*

*shading; darker shades of blue indicate higher probabilities. Similarly, the true posterior distribution is*

*represented as an orange shading, and the likelihood of the sample is represented by a grey shading. The*

*axes represent the mean and standard deviation of a signal. Units are arbitrary. **A** The fixed-shape belief*

*distribution encodes the learner's knowledge before a new observation is made. **B** A new observation is*

*made. The likelihood of the observed value is indicated by the grey shading. Using the fixed-shape belief*

*as a prior, a posterior can be computed. The posterior's mode is different from the mode of the fixed-*

*shape belief; therefore, an update (indicated by a white arrow) is required. **C** The fixed-shape belief has*

*been modified such that its mode aligns with the mode of the true posterior. **D** The modified fixed-shape*

*belief represents the learner's knowledge after the new observation has been taken into account. The*

*distributions were computed using the densities given in Eq. S7 – S9.*

689

After the update, the fixed-shape learner's belief is still different from the true posterior. This is because

the shape of the true posterior is generally not the same as the fixed belief shape that the learner uses.

Hence, mode-matching is only an approximation of Bayesian learning, and some features are lost in this

approximation.

694   Mode-matching is formally related to the variational Bayes scheme (36, 37, 43, 44), which works by

695   minimizing the Kullback-Leibler divergence between the true posterior and a fixed belief shape (usually a

696   multivariate normal distribution). However, mode-matching does not minimise the Kullback-Leibler

697   divergence; instead, it minimises the distance between the modes of the distributions.

698   The learning rules that can be derived with the mode-matching method are not as precise as those derived

699   from variational Bayes, let alone fully Bayesian learning. What makes mode-matching interesting is that

700   it can be used to derive relatively simple, tractable learning rules, as we shall see in the next section.

701   *New learning rules via mode-matching*

702   Let us consider a situation in which an organism tracks the size of a reward associated with some

703   behavior. By engaging in that behavior, it samples the reward size $r$. Using these samples, it attempts to

704   estimate the mean reward $\mu$ that can be expected from performing the behavior at any given time.

705   To derive the learning rules for this situation, we start with a generative model for the reward process, and

706   the learner's fixed-shape belief distributions over the process variables. We model rewards as normally

707   distributed around a mean $\mu$, with a standard deviation $\sigma$, as defined in Eq. 1. Note that $\sigma$ quantifies trial-

708   by-trial fluctuations, and therefore observation noise. The distribution in Eq. 1 is stationary; this means

709   that the environment is modelled as stable.

710   We further assume that the learner maintains beliefs $M$ and $S$ about $\mu$ and $\sigma$, in form of a normal

711   distribution over possible values of $\mu$ and a gamma distribution over possible value of $\sigma$:

712   $$M \sim N(m, \sigma_m)$$

713                                                                                                                        Eq. S1

714   $$S \sim \Gamma(a, b)$$

715                                                                                                                        Eq. S2

716    The learner can change its beliefs by adapting the mean (and hence mode) $m$ of the normal distribution,

717    and the mode $\frac{a-1}{b}$ of the gamma distribution. The standard deviation $\sigma_m$ and the rate parameter $b$ stay

718    fixed.

719    How should we interpret this belief encoding? Allowing $m$ and $s$ to vary implies that the learner considers

720    both the mean reward $\mu$ and the observation noise $\sigma$ as unknown—it can adapt its beliefs about these

721    variables. Fixing $\sigma_m$ and $b$ implies that the learner's uncertainty about the mean and the standard deviation

722    of the signal are kept constant—it cannot adapt those. The learner will thus not become more certain

723    about either the mean or the standard deviation as it gathers more and more data. Fixing $\sigma_m$ and $b$ keeps

724    the resulting learning rules simple. An additional advantage of this design arises when the environment

725    fails to be stationary—then, high certainty about the tracked variables would prevent the learner from

726    adapting to new situations. The model of the reward generating process and the learner's belief system

727    form our central assumptions—the rest follows. The learner we are about to derive will interpret all

728    rewards it sees as being sampled from a normal distribution with fixed mean and variance, and it will

729    make its inferences accordingly.

730    Now, let us use mode-matching to derive learning rules from our assumptions. To find out how the

731    learner should update $m$ and $s$ after sampling a reward $r$, we must first find the mode of the true posterior

732    distribution. For this, we can use a well-known way to simplify calculations. Bayes' theorem states that

733    
$$P(x\,|\,\theta) = P(x|\theta)P(\theta)/P(x)$$

734    Eq. S3

735    with $\theta$ the parameters that are to be inferred and $x$ the data that is observed. Now we notice that

736    
$$\log P(\theta|x) = \log P(x\,|\,\theta) + \log P(\theta) - \log P(x)$$

737    Eq. S4

738    with the last term independent of the parameters $\theta$. We can define the function

36

739
$$E = \log P(x \mid \theta) + \log P(\theta)$$

740
Eq. S5

741 and it is easy to see that the parameters $\theta_{max}$ that maximize the function $E$ also maximise the posterior

742 distribution $P(\theta|x)$ (this is true because the logarithm is strictly monotonic). The function $E$, often called

743 *energy* in analogy to statistical physics, is related to the famous of free energy function which plays a key

744 role in many contemporary theories of brain function (38, 45, 46). In the case at hand, the function $E$ is

745 given as

746
$$E = \log(P(r \mid M,S)) + \log(P(M \mid m,\sigma_m)P(S \mid a,b))$$

747
$$= \log\left(S^{-1}\exp\left(-\frac{1}{2}\frac{(r-M)^2}{S^2}\right)\exp\left(-\frac{1}{2}\frac{(M-m)^2}{\sigma_m^2}\right)S^{a-1}\exp(-Sb)\right) + C$$

748
Eq. S6

749 with $C$ a term that does not depend on $M$ or $S$, and

750
$$P(r \mid M,S) = (2\pi S^2)^{-\frac{1}{2}}\exp\left(-\frac{1}{2}\frac{(r-M)^2}{S^2}\right)$$

751
Eq. S7

752
$$P(M \mid m,\sigma_m) = (2\pi\sigma_m^2)^{-\frac{1}{2}}\exp\left(-\frac{1}{2}\frac{(M-m)^2}{\sigma_m^2}\right)$$

753
Eq. S8

754
$$P(S \mid a,b) = \frac{b^a}{\Gamma(a)}S^{a-1}\exp(-Sb)$$

755
Eq. S9

37

756     the probability density functions associated with the distributions Eq. 1, Eq. S1 and Eq. S2. To find the

757     maximum of $E$ with respect to $M$ and $S$, and hence the mode of the posterior, we can investigate the

758     gradient $\left(\frac{\partial E}{\partial M}, \frac{\partial E}{\partial S}\right)$ of $E$, which vanishes at the maximum. Evaluation the conditions $\frac{\partial E}{\partial M} = 0$ and $\frac{\partial E}{\partial S} = 0$,

759     we find

760
$$M_{max} - m = \frac{\sigma_m^2}{S_{max}^2}(r - M_{max})$$

761                                                                    Eq. S10

762
$$S_{max} - s = \frac{1}{b}\left(\left(\frac{r - M_{max}}{S_{max}}\right)^2 - 1\right)$$

763                                                                     Eq. S11

764     for the location $(M_{max}, S_{max})$ of the maximum of E. In Eq. S11, $s = \frac{a-1}{b}$ is the mode of the gamma

765     distribution.

766     To interpret these equations, note that the left-hand side yields the distance of $S_{max}$ and $M_{max}$ from the

767     mode of their respective prior distributions. The right-hand side quantifies the mismatch between what

768     was expected based on $M_{max}$ and $S_{max}$ and what actually happened: based on $M_{max}$ and $S_{max}$, the reward

769     $r$ was expected to be close to $M_{max}$ and $(r - M_{max})^2$ was expected to be close to $S_{max}^2$. The mismatches

770     are weighted with a measure of prior narrowness, $\sigma_m^2$ in Eq. S10 and $\frac{1}{b}$ in Eq. S11.

771     We now must solve these equations for $M_{max}$ and $S_{max}$ to find the mode of the true posterior. We could

772     try and find the exact solutions, but considering that the equations are nonlinear, we would have to expect

773     complicated expressions. Here we will not choose that route: we shall restrict ourselves to approximate

774     solutions.

775     We focus on the scenario in which the priors of both $M_{max}$ and $S_{max}$ are very narrow. Formally, this

776     corresponds to $\sigma_m^2 \ll 1$ and $\frac{1}{b} \ll 1$, or equivalently $\sigma_m^2 \sim \epsilon$ and $\frac{1}{b} \sim \epsilon$ with $\epsilon \ll 1$. To derive an approximate

777     solution for this regime, we use expansions

778 $$M_{max} = M_{max,\,0} + M_{max,1} + O(2)$$

779                               Eq. S12

780 $$S_{max} = S_{max,\,0} + S_{max,1} + O(2)$$

781                               Eq. S13

782     for the variables we want to solve for. Here, $M_{max,1} \sim \epsilon$ and $S_{max,1} \sim \epsilon$ are first order terms with respect

783     to the small constants $\sigma_m^2$ and $\frac{1}{b}$. These expansions can be thought of as Taylor expansions of the variables

784     of interest, keeping only terms up to first order. To determine the zeroth and first order terms, we insert

785     these expansions into Eq. S10 and Eq. S11 and collect all terms of a certain order. Using this procedure,

786     we obtain

787 $$M_{max,0} = m$$

788                               Eq. S14

789 $$S_{max,0} = s$$

790                               Eq. S15

791     for the zeroth order and

792 $$M_{max,1} = \frac{\sigma_m^2}{S_{max,0}^2}(r - M_{max,0}) = \frac{\sigma_m^2}{s^2}(r - m)$$

793                               Eq. S16

794
$$S_{max,1} = \frac{1}{b}\left(\left(\frac{r - M_{max,0}}{S_{max,0}}\right)^2 - 1\right) = \frac{1}{b}\left(\left(\frac{r-m}{s}\right)^2 - 1\right)$$

795 Eq. S17

796  for the first order, where the zeroth order results in Eq. S14 and Eq. S15 were already used. Reinserting

797  these contributions into Eq. S12 and Eq. S13, we find that the mode of the posterior is approximately at

798
$$M_{max} = m + \frac{\sigma_m^2}{s^2}(r - m) + O(2)$$

799 Eq. S18

800
$$S_{max} = s + \frac{1}{b}\left(\left(\frac{r-m}{s}\right)^2 - 1\right) + O(2)$$

801 Eq. S19

802  where $O(2)$ reminds us that we have neglected terms of second or higher order in $\frac{1}{b}$ and $\sigma_m^2$. The mode of

803  the posterior is now found—at least approximately. The final step of the mode-matching method consists

804  in updating the mode of the fixed-shape belief distribution—which is $(m, s)$—by aligning it with the

805  maximum of the true posterior, which is (approximately) given by $(M_{max}, S_{max})$ in Eq. S18 and Eq. S19.

806  If we were just looking for computationally lightweight learning rules that approximate Bayesian

807  learning, we could stop here. However, we are ultimately interested in modelling learning in biological

808  systems, in particular the basal ganglia system. We must hence consider that changes in synaptic strength

809  can only depend on local information (such as pre- and postsynaptic potentials) and low-dimensional

810  global feedback signals (such as dopamine release in the striatum). We can achieve this here by applying

811  yet another set of approximations. First, we identify certain factors as learning rates: $\frac{\sigma_m^2}{s^2}$ is replaced by $\alpha_m$,

812  and $\frac{1}{b}$ by $\alpha_s$. Then, we simplify the equations by making the learning rates constant; we hence omit the $s$-

813  dependence of $\alpha_m$. With these changes, we arrive at the learning rules specified in Eq. $2 - 4$. These rules

814     feature a global feedback signal $\delta$ and track the mean reward $m$ as well as the observation noise $s$. Both

815     $m$ and $s$ are fed back into the learning system as they enter what we will call the ***scaled*** prediction error $\delta$.

816     *Appendix S2. The high noise limit of the steady state Kalman filter*

817     Here, we show that the SPE learning rules approximate the one-dimensional steady-state Kalman filter in

818     the limit of high observation noise. We start by defining the Kalman filter model. We then derive the

819     steady-state Kalman filter, and finally take the high-noise limit.

820     *The definition of the Kalman filter*

821     The Kalman filter is a computational method for state estimation and prediction (3). It can be derived

822     from Bayesian principles and is optimal for tracking signals with certain characteristics. Here, we focus

823     on a one-dimensional Kalman filter which is used for predicting rewards, following Piray and Daw (2).

824     The rules they use are

825
$$m_t = m_{t-1} + k_t(r_t - m_{t-1})$$

826      Eq. S20

827
$$k_t = \frac{w_{t-1} + v^2}{w_{t-1} + v^2 + \sigma^2}$$

828      Eq. S21

829
$$w_t = (1 - k_t)(w_{t-1} + v^2)$$

830      Eq. S22

831     where $r_t$ is the reward, $k_t$ the learning rate or Kalman gain and $w_t$ the posterior variance in trial in trial $t$.

832     Note that our notation differs slightly from that of Piray and Daw (2), for the sake of consistency within

833     this work.

834    The above rules can be shown to be optimal for tracking signals such as those we used above, i.e., signals

835    that consist of samples drawn from a normal distribution with a drifting mean (3).

836    *The steady-state Kalman filter*

837    The Kalman filter has several variables that must be updated on every trial. If one requires a simpler

838    model with almost similar properties, one option is to use a Kalman filter in the limit $t \to \infty$: as for $t \to \infty$,

839    the posterior variance $w_t$ and the Kalman gain $k_t$ converge to limits $w_\infty$ and $k_\infty$.

840    Eq. S20 with $k_\infty$ instead of $k_t$ is called a **steady-state Kalman filter**. By construction, the normal Kalman

841    filter becomes more similar to the steady-state Kalman filter the more trials pass. In practice, performance

842    often does not differ much between the two (3).

843    What are the limits $w_\infty$ and $k_\infty$? One may use Eq. S21 and Eq. S22 to determine them. By setting $k_t =$

844    $k_{t-1}$ and $w_t = w_{t-1}$, we find

845
$$w_\infty = \frac{v^2}{2}\left(\sqrt{4\frac{\sigma^2}{v^2} + 1} - 1\right)$$

846                                                                                             Eq. S23

847
$$k_\infty = \frac{\sqrt{4\frac{\sigma^2}{v^2} + 1} + 1}{\sqrt{4\frac{\sigma^2}{v^2} + 1} + 1 + 2\frac{\sigma^2}{v^2}}$$

848                                                                                             Eq. S24

849    To use the steady-state Kalman filter, one just needs to compute $k_\infty$ and plug it into Eq. S20. One can

850    then use this single equation to track the signal, with no other computations required. The steady-state

851    Kalman filter is thus equivalent to the RW model in Eq. 14 and Eq. 15, parametrized with an optimal

852    learning rate (that is to say, optimal for a signal with statistics $v^2$ and $\sigma^2$).

853 *The high-noise limit*

854 The steady-state Kalman filter is less complex than the full Kalman filter. However, its learning rate $k_\infty$ is

855 still a complex function of the signal statistics $v$ and $\sigma$. Can it be simplified? Let us consider of a signal

856 with high observation noise, i.e., with $\sigma^2$ much larger than $v^2$. Using Eq. S24, we find that

857
$$k_\infty \to \frac{v}{\sigma}$$

858 Eq. S25

859 for $\frac{\sigma^2}{v^2} \to \infty$. This means that a steady-state Kalman filter with gain $\frac{v}{\sigma}$ is approximately optimal for signals

860 with $\sigma^2 \gg v^2$. In Fig S2, we compare the optimal steady-state learning rate $k_\infty$ with the approximately

861 optimal learning rate $v/\sigma$ for different levels of $\sigma$, with $v$ fixed at $v = 1$.
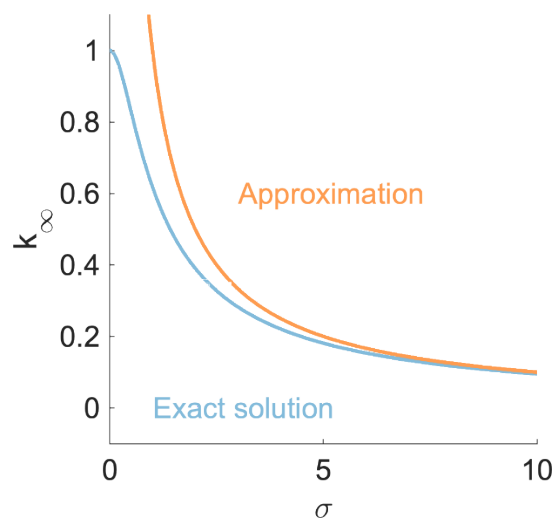
862



863

864 **Fig S2. The learning rate of the steady state Kalman filter.** *We show the learning rate $k_\infty$ of the steady*

865 *state Kalman filter as a function of the observation noise $\sigma$. We provide the exact value (blue line) and the*

866 *approximation $k_\infty \approx \frac{v}{\sigma}$ (orange line).*

867

868     We find that the approximation becomes very close very quickly—for $\frac{\sigma}{\nu} > 2$, the relative difference

869     between the optimal learning rate and its approximation is already less than 30 %. Fig S2 further suggests

870     that the approximation breaks down as $\frac{\nu}{\sigma}$ approaches unity—the optimal learning rate for signals with

871     $\sigma = 0$ is one; any higher learning rate will be detrimental for the performance.

872     In summary, we find the learning rule

873
$$m_t = m_{t-1} + \frac{\nu}{\sigma}(r_t - m_{t-1})$$

874                                                                      Eq. S26

875     to be approximately optimal for $\sigma \ll \nu$ and large $t$. The rule Eq. 26 bears striking resemblance to one of

876     the SPE learning rules: the rule in Eq. 3. The difference between the two rules is just how the scaling is

877     attributed: in the Kalman filter, one would perhaps speak of a scaled learning rate, while in the SPE

878     model, one attributes the scaling to the error term. Mathematically, both formulations are equivalent.

879     A real difference between the Kalman filter and the SPE model is that the latter has a mechanism to track

880     $\sigma$. No such mechanism exists in the Kalman filter. Both models require $\nu$ as an external input (for the SPE

881     model, the corresponding parameter is $\alpha_m$).

882     We conclude that the SPE model can be viewed as an implementation of approximately optimal one-

883     dimensional state estimation, equipped with a mechanism to supply some of the required parameters—the

884     observation noise $\sigma$. Other models have been proposed to track the process noise $\nu$, for example by Piray

885     and Daw (2). A combination of these approaches might be an interesting direction for future research.

886

# References

1. Schultz W, Dayan P, Montague PR. A neural substrate of prediction and reward. Science. 1997;275(5306):1593-9.

2. Piray P, Daw ND. A simple model for learning in volatile environments. PLoS computational biology. 2020;16(7):e1007963.

3. Simon D. Optimal state estimation: Kalman, H infinity, and nonlinear approaches: John Wiley & Sons; 2006.

4. Grewal MS, Andrews AP. Applications of Kalman filtering in aerospace 1960 to the present [historical perspectives]. IEEE Control Systems Magazine. 2010;30(3):69-78.

5. Gershman SJ. Dopamine, inference, and uncertainty. Neural Computation. 2017;29(12):3311-26.

6. Szirtes G, Póczos B, Lőrincz A. Neural kalman filter. Neurocomputing. 2005;65:349-55.

7. Wolpert DM. Computational approaches to motor control. Trends in cognitive sciences. 1997;1(6):209-16.

8. Kakei S, Tanaka H, Ishikawa T, Tomatsu S, Lee J. The Input-Output Organization of the Cerebrocerebellum as Kalman Filter. Cerebellum as a CNS Hub: Springer; 2021. p. 391-411.

9. Piray P, Daw ND. Unpredictability vs. volatility and the control of learning. bioRxiv. 2020.

10. Tobler PN, Fiorillo CD, Schultz W. Adaptive coding of reward value by dopamine neurons. Science. 2005;307(5715):1642-5.

11. Mikhael JG, Bogacz R. Learning reward uncertainty in the basal ganglia. PLoS computational biology. 2016;12(9):e1005062.

12. Dabney W, Kurth-Nelson Z, Uchida N, Starkweather CK, Hassabis D, Munos R, et al. A distributional code for value in dopamine-based reinforcement learning. Nature. 2020;577(7792):671-5.

13. Gerfen CR, Engber TM, Mahan LC, Susel Z, Chase TN, Monsma FJ, et al. D1 and D2 dopamine receptor-regulated gene expression of striatonigral and striatopallidal neurons. Science. 1990;250(4986):1429-32.

14. Surmeier DJ, Ding J, Day M, Wang Z, Shen W. D1 and D2 dopamine-receptor modulation of striatal glutamatergic signaling in striatal medium spiny neurons. Trends in neurosciences. 2007;30(5):228-35.

15. Gerfen CR, Surmeier DJ. Modulation of striatal projection systems by dopamine. Annual review of neuroscience. 2011;34:441-66.

16. Collins AG, Frank MJ. Opponent actor learning (OpAL): Modeling interactive effects of striatal dopamine on reinforcement learning and choice incentive. Psychological review. 2014;121(3):337.

17. Frank MJ, Seeberger LC, O'reilly RC. By carrot or by stick: cognitive reinforcement learning in parkinsonism. Science. 2004;306(5703):1940-3.

18. Sommer MA. The role of the thalamus in motor control. Current opinion in neurobiology. 2003;13(6):663-70.

19. Redgrave P, Prescott TJ, Gurney K. The basal ganglia: a vertebrate solution to the selection problem? Neuroscience. 1999;89(4):1009-23.

20. Rescorla RA, Wagner AR. A theory of Pavlovian conditioning: Variations in the effectiveness of reinforcement and nonreinforcement. Classical conditioning II: Current research and theory. 1972;2:64-99.

21. Rothenhoefer KM, Hong T, Alikaya A, Stauffer WR. Rare rewards amplify dopamine responses. Nature neuroscience. 2021;24(4):465-9.

22. Moeller M, Grohn J, Manohar S, Bogacz R. An association between prediction errors and risk-seeking: Theory and behavioral evidence. PLoS computational biology. 2021;17(7):e1009213.

23. Shen W, Flajolet M, Greengard P, Surmeier DJ. Dichotomous dopaminergic control of striatal synaptic plasticity. Science. 2008;321(5890):848-51.

933   24.  Fisher SD, Robertson PB, Black MJ, Redgrave P, Sagar MA, Abraham WC, et al. Reinforcement
934   determines the timing dependence of corticostriatal synaptic plasticity in vivo. Nature communications.
935   2017;8(1):1-13.
936   25.  Dreyer JK, Herrik KF, Berg RW, Hounsgaard JD. Influence of phasic and tonic dopamine release on
937   receptor activation. Journal of Neuroscience. 2010;30(42):14273-83.
938   26.  Dodson PD, Dreyer JK, Jennings KA, Syed EC, Wade-Martins R, Cragg SJ, et al. Representation of
939   spontaneous movement by dopaminergic neurons is cell-type selective and disrupted in parkinsonism.
940   Proceedings of the National Academy of Sciences. 2016;113(15):E2180-E8.
941   27.  Möller M, Bogacz R. Learning the payoffs and costs of actions. PLoS computational biology.
942   2019;15(2):e1006285.
943   28.  Ferrucci L, Nougaret S, Brunamonti E, Genovesio A. Effects of reward size and context on learning in
944   macaque monkeys. Behavioural brain research. 2019;372:111983.
945   29.  Preuschoff K, Bossaerts P, Quartz SR. Neural differentiation of expected reward and risk in human
946   subcortical structures. Neuron. 2006;51(3):381-90.
947   30.  White JK, Monosov IE. Neurons in the primate dorsal striatum signal the uncertainty of object–
948   reward associations. Nature communications. 2016;7(1):1-8.
949   31.  Lahiri AK, Bevan MD. Dopaminergic transmission rapidly and persistently enhances excitability of
950   D1 receptor-expressing striatal projection neurons. Neuron. 2020;106(2):277-90. e6.
951   32.  St Onge JR, Floresco SB. Dopaminergic modulation of risk-based decision making.
952   Neuropsychopharmacology. 2009;34(3):681-97.
953   33.  Zalocusky KA, Ramakrishnan C, Lerner TN, Davidson TJ, Knutson B, Deisseroth K. Nucleus
954   accumbens D2R cells signal prior outcomes and control risky decision-making. Nature.
955   2016;531(7596):642-6.
956   34.  Karin O, Alon U. The dopamine circuit as a reward-taxis navigation system. bioRxiv. 2021.
957   35.  Friston KJ, Trujillo-Barreto N, Daunizeau J. DEM: a variational treatment of dynamic systems.
958   Neuroimage. 2008;41(3):849-85.
959   36.  Buckley CL, Kim CS, McGregor S, Seth AK. The free energy principle for action and perception: A
960   mathematical review. Journal of Mathematical Psychology. 2017;81:55-79.
961   37.  Bogacz R. A tutorial on the free-energy framework for modelling perception and learning. Journal of
962   mathematical psychology. 2017;76:198-211.
963   38.  Bogacz R. Dopamine role in learning and action inference. Elife. 2020;9:e53262.
964   39.  Hessel M, Soyer H, Espeholt L, Czarnecki W, Schmitt S, van Hasselt H, editors. Multi-task deep
965   reinforcement learning with popart. Proceedings of the AAAI Conference on Artificial Intelligence; 2019.
966   40.  Kingma DP, Ba J. Adam: A method for stochastic optimization. arXiv preprint arXiv:14126980. 2014.
967   41.  Paz JT, Chavez M, Saillet S, Deniau J-M, Charpier S. Activity of ventral medial thalamic neurons
968   during absence seizures and modulation of cortical paroxysms by the nigrothalamic pathway. Journal of
969   Neuroscience. 2007;27(4):929-41.
970   42.  Montague PR, McClure SM, Baldwin P, Phillips PE, Budygin EA, Stuber GD, et al. Dynamic gain
971   control of dopamine delivery in freely moving animals. Journal of Neuroscience. 2004;24(7):1754-9.
972   43.  Dayan P, Hinton GE, Neal RM, Zemel RS. The helmholtz machine. Neural computation.
973   1995;7(5):889-904.
974   44.  Kingma DP, Welling M. Auto-encoding variational bayes. arXiv preprint arXiv:13126114. 2013.
975   45.  Friston K. The free-energy principle: a unified brain theory? Nature reviews neuroscience.
976   2010;11(2):127-38.
977   46.  Gershman SJ. What does the free energy principle tell us about the brain? arXiv preprint
978   arXiv:190107945. 2019.