



## ARTICLE



# Impulsivity and risk-seeking as Bayesian inference under dopaminergic control

John G. Mikhael <sup>1,2</sup>✉ and Samuel J. Gershman <sup>3,4</sup>

© The Author(s), under exclusive licence to American College of Neuropsychopharmacology 2021

Bayesian models successfully account for several of dopamine (DA)'s effects on contextual calibration in interval timing and reward estimation. In these models, tonic levels of DA control the precision of stimulus encoding, which is weighed against contextual information when making decisions. When DA levels are high, the animal relies more heavily on the (highly precise) stimulus encoding, whereas when DA levels are low, the context affects decisions more strongly. Here, we extend this idea to intertemporal choice and probability discounting tasks. In intertemporal choice tasks, agents must choose between a small reward delivered soon and a large reward delivered later, whereas in probability discounting tasks, agents must choose between a small reward that is always delivered and a large reward that may be omitted with some probability. Beginning with the principle that animals will seek to maximize their reward rates, we show that the Bayesian model predicts a number of curious empirical findings in both tasks. First, the model predicts that higher DA levels should normally promote selection of the larger/later option, which is often taken to imply that DA decreases 'impulsivity,' and promote selection of the large/risky option, often taken to imply that DA increases 'risk-seeking.' However, if the temporal precision is sufficiently decreased, higher DA levels should have the opposite effect—promoting selection of the smaller/sooner option (higher impulsivity) and the small/safe option (lower risk-seeking). Second, high enough levels of DA can result in preference reversals. Third, selectively decreasing the temporal precision, without manipulating DA, should promote selection of the larger/later and large/risky options. Fourth, when a different post-reward delay is associated with each option, animals will not learn the option-delay contingencies, but this learning can be salvaged when the post-reward delays are made more salient. Finally, the Bayesian model predicts correlations among behavioral phenotypes: Animals that are better timers will also appear less impulsive.

*Neuropsychopharmacology*; <https://doi.org/10.1038/s41386-021-01125-z>

## INTRODUCTION

The neuromodulator dopamine (DA) has been repeatedly associated with choice impulsivity, the tendency to prioritize short-term over long-term reward. Impulsive behaviors characterize a number of DA-related psychiatric conditions [1], such as attention-deficit/hyperactivity disorder [2–6], schizophrenia [7, 8], addiction [9, 10], and dopamine dysregulation syndrome [11, 12]. Furthermore, direct pharmacological manipulation of tonic DA levels in humans [13, 14] and rodents [15, 16] has corroborated a relationship between DA and impulsivity. The standard approach to measuring impulsive choice is the intertemporal choice task (ITC), in which subjects choose between a small reward delivered soon and a large reward delivered later [17]. A subject's preference for the smaller/sooner option is often taken as a measure of their impulsivity, or the extent to which they discount future rewards [18–21].

In the majority of animal studies, higher DA levels have been found to promote selection of the larger/later option (inhibiting impulsivity) [15, 22–28]. However, the inference that DA agonists inhibit impulsivity has been challenged in recent years, in part because, when ITCs are administered to humans,

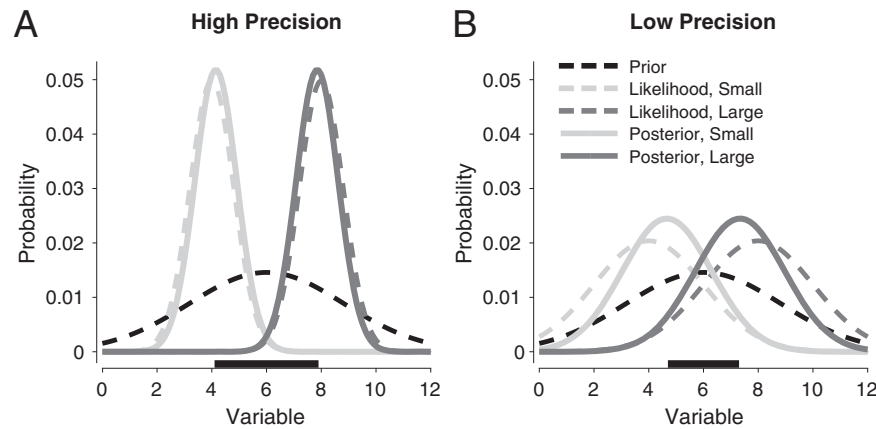
DA agonists seem to *promote* impulsivity [29]. Perhaps relevant to this contrast is that, while impulsive choices in humans are assessed through hypothetical situations ('Would you prefer \$1 now or \$10 in one month?'), ITCs in animals more closely resemble reinforcement learning tasks involving many trials of experienced rewards and delays. Complicating this picture further, the effect of DA, even within animal studies, is not straightforward. While in most studies, DA appears to decrease impulsivity, DA has been found to systematically increase impulsivity under some conditions [30–32], such as when the delay period is uncued [16] or when different delays for the larger/later option are presented in decreasing order across training blocks [33].

The relationship of DA with impulsive choice finds a parallel in its relationship with risk-seeking. Disruptions in risk preferences feature prominently in a number of DA-related conditions [1], including Parkinson's disease [34–37], schizophrenia [38, 39], and attention-deficit/hyperactivity disorder [40]. Moreover, direct manipulation of DA levels in Parkinson's patients [41], healthy humans [42], and rodents [29] has further established a link between DA and risk preferences.

<sup>1</sup>Program in Neuroscience, Harvard Medical School, Boston, MA, USA. <sup>2</sup>MD-PhD Program, Harvard Medical School, Boston, MA, USA. <sup>3</sup>Department of Psychology and Center for Brain Science, Harvard University, Cambridge, MA, USA. <sup>4</sup>Center for Brains, Minds and Machines, Massachusetts Institute of Technology, Cambridge, MA, USA. ✉email: john\_mikhael@hms.harvard.edu

Received: 1 November 2020 Revised: 17 July 2021 Accepted: 21 July 2021

Published online: 10 August 2021



**Fig. 1 Contextual influence is stronger when the encoding precision is low.** Distributions for two signals, one small and the other large. **A** When the encoding precision is high compared to the prior precision, the posteriors do not deviate significantly from the likelihood. **B** As the encoding precision decreases, the posteriors migrate toward the prior. The horizontal black segments illustrate the difference in posterior means under high vs. low precision.

Risk-seeking can be formalized as the tendency to prioritize uncertain rewards over less uncertain rewards of equal average value. For example, a risk-seeker will preferentially select an option yielding a reward of magnitude 10 on 50% of trials and no reward in the remaining trials, over an option yielding a reward of magnitude 5 on 100% of trials. A standard measure of risk-seeking is the probability discounting task (PD) [43–45], where subjects choose between a small reward delivered with complete certainty and a large reward that is only delivered with some probability. Subjects that are more likely to select the large/risky option than other subjects—regardless of reward probability—are labeled as being more risk-seeking. Though studies involving direct pharmacological manipulation have highlighted a key role for DA in setting this preference, the directionality of DA's effect has remained unclear: Whereas St Onge and Floresco [46] have found that increasing the DA level promotes risk-seeking in PDs, follow-up work has shown that DA may have exactly the opposite effect, depending on how the training blocks are ordered [47], a variable whose relevance for risk preferences is not immediately obvious.

Animal behavior in ITCs and PDs can be reinterpreted from a reinforcement learning perspective. With repeated trials of the same task, an optimal agent can learn to maximize its total accumulated rewards by estimating the reward rate for each option (reward magnitude divided by total trial duration) and choosing the option with the higher reward rate. Thus if the larger/later option has a sufficiently large reward or sufficiently short delay, it will be the optimal choice. However, if its reward were sufficiently small or its delay sufficiently long, the smaller/sooner option may be the superior choice instead, without any assumption of 'discounting.' Under this view, animals do not necessarily discount rewards at all, but rather make choices based on a reward-rate computation. The notions of true impulsivity in ITCs and risk-seeking in PDs have persisted, however, because animals tend to choose the smaller/sooner and large/risky options even when they objectively yield fewer rewards over many trials.

To address the question of whether animals compare reward rates, a body of theoretical and experimental work focused on impulsivity has demonstrated that the suboptimal tendency to choose the smaller/sooner option is better explained by *temporal* biases than by biases of choice [48–50] (see also [51]). This work has shown that animals behave in a way consistent with maximizing their reward rates, but they underestimate the elapsed time—and in particular, the periods after receiving the reward and before beginning the next trial. Thus animals estimate the reward rates for each option based largely on the pre-reward delays. This bias disproportionately benefits the

smaller/sooner option, which has a much shorter pre-reward delay. As a result, the animals make choices that can be interpreted as impulsive. Said differently, animals disproportionately underestimate the total trial duration for the smaller/sooner option compared to the larger/later option, making the former more appealing. While this discounting-free view derives animal behavior from a normative framework (maximizing reward rates), how and why DA modulates choice preferences remains the subject of much speculation.

In this paper, we build on recent theoretical work that cast DA in a Bayesian light [52, 53]. Here, DA controls the precision with which cues are internally represented, which in turn controls the extent to which the animal's estimates of the cues are influenced by context. In Bayesian terms, which we discuss below, DA controls the precision of the likelihood relative to that of the prior (the context). This framework predicts a well-replicated result in the interval timing literature, referred to as the 'central tendency' effect: When temporal intervals of different lengths are reproduced under DA depletion (e.g., in unmedicated Parkinson's patients), shorter intervals tend to be overproduced and longer intervals tend to be underproduced, and DA repletion rescues accurate timing [54–56]. We recently extended this framework to the representation of reward estimates [57]. In this case, the Bayesian framework predicts that DA should tip the exploration-exploitation balance toward exploitation, in line with empirical findings [58–60] (but see [61, 62]).

We show here that, under the Bayesian theory, higher DA levels should promote behaviors consistent with lower impulsivity in the standard ITC (selection of the larger/later option), but should have the opposite effect when the temporal precision of the delay period is selectively and sufficiently reduced. In both cases, high enough levels of DA should elicit preference reversals, and not only an amplification of the current preference. Furthermore, in manipulations of temporal precision, if animals are more likely to select the larger/later option at baseline, DA administration will tend to reverse that preference (promote the smaller/sooner option), and vice versa. We show that animals should not learn the contingencies between options and their post-reward delays, but that this learning can be salvaged if the post-reward delays are made more salient. We show that animals that display more precise behaviors in interval timing tasks should also appear less impulsive. Finally, we reproduce this analysis for the case of risk-seeking and PDs: Depending on the relative balance between the uncertainty about the reward magnitude and the uncertainty about the reward probability, we show that DA can either promote or suppress selection of the large/risky option.

## METHODS

### The Bayesian theory of dopamine

An agent wishing to encode information about some cue must contend with noise at every level, including the information source (which is seldom deterministic), storage (synapses are noisy), and signaling (neurons are noisy) [63]. We can formalize the noisy encoding as a mapping from an input signal (e.g., experienced reward) to a distribution over output signals (e.g., firing rates). For the purposes of this paper, we will remain agnostic about the specific neural implementation of the mapping, and instead discuss it in abstract terms. Thus a noisy encoding of some variable can be represented by a distribution over values: Tight distributions correspond to encodings with low noise (Fig. 1A), whereas wide distributions correspond to encodings with high noise (Fig. 1B).

Consider, then, a scenario in which an animal must estimate the average yield of a reward source from noisy samples. Because of the animal's uncertainty about the average yield (the encoding distribution has non-zero spread), its final estimate can be improved by utilizing other sources of information. For example, if the nearby reward sources tend to yield large rewards, then the animal should form an optimistic estimate of the reward source's average yield. Similarly, if nearby reward sources yield small rewards, then the animal should form a pessimistic estimate. Formally, the contextual information can be used to construct a prior distribution over average yield, and the encoding distribution can be used to construct a likelihood function for evaluating the consistency between the encoded information and a hypothetical average yield. Bayes' rule stipulates that the animal's final probabilistic estimate should reflect the product of the likelihood and prior:

$$p(\mu|m) \propto p(m|\mu)p(\mu), \quad (1)$$

referred to as the posterior distribution. Here,  $\mu$  is the variable being estimated (the reward yield),  $m$  is the stored value,  $p(m|\mu)$  is the likelihood, and  $p(\mu)$  is the prior. For simplicity, we take these distributions to be Gaussian throughout. Under standard assumptions for Gaussian distributions, the estimate  $\hat{\mu}$  corresponds to the posterior mean:

$$\hat{\mu} = \left( \frac{\lambda_0}{\lambda_0 + \lambda} \right) \mu_0 + \left( \frac{\lambda}{\lambda_0 + \lambda} \right) \mu. \quad (2)$$

Here,  $\mu_0$ ,  $\lambda_0$ ,  $\mu$ , and  $\lambda$  represent the prior mean, prior precision, likelihood mean, and encoding precision, respectively. In words, the agent takes a weighted average of the prior mean  $\mu_0$  and the likelihood mean  $\mu$ —weighted by their respective precisions  $\lambda_0$  and  $\lambda$  after normalization—to produce its estimate, the posterior mean  $\hat{\mu}$ . Intuitively, the tighter each distribution, the more it pulls the posterior mean in its direction.

The Bayesian theory of DA asserts that increasing the DA level increases the encoding precision  $\lambda$ , where the prior here represents the distribution of stimuli (i.e., the context). Thus when DA is high, the estimate  $\hat{\mu}$  does not heavily depend on contextual information, whereas when it is low, Bayesian migration of the estimate to the prior is strong (compare Fig. 1A and B). Shi et al. [56] have applied this theory to interval timing and shown that it predicts DA's effects on the central tendency: Parkinson's patients who are on their medication will have high  $\lambda$ , qualitatively corresponding to Fig. 1A. Then the temporal estimates for the short and long durations will be very close to their true values (here, 4 and 8 s). On the other hand, patients who are off their medication will have low  $\lambda$ , corresponding to Fig. 1B. Thus the estimates for both durations will migrate toward the prior mean, or the average of the two durations. In other words, the estimate for the short duration will be overproduced, and the estimate for the long duration will be underproduced, as observed [54, 55].

The Bayesian model can also be applied to reward magnitudes [64, 65]. Imagine a bandit task in which an agent samples from two reward sources, one yielding small rewards on average and the other yielding large rewards on average. Under lower levels of DA, the central tendency should decrease the difference between the two reward estimates (compare lengths of black segments on the x-axis in Fig. 1A and B). Under standard models of action selection, animals are more likely to choose the large option when the difference between the two estimates is large, and become more and more likely to sample other options as the difference decreases (see next section). This means that lower levels of DA should promote selection of the smaller reward, often taken to indicate a drive to 'explore,' as empirically observed [58–60] (but see [61, 62]). Thus, as previously proposed, DA may be interpreted as controlling the exploration-exploitation trade-off [66]. This is in line with the 'gain control' theory of DA, in which high DA levels have been hypothesized to amplify the difference between reward estimates during decision making [66–70]. The

Bayesian theory of DA subsumes the gain control view (see next section), but importantly, under this theory, animals do not become intrinsically more explorative or exploitative under different DA levels, but rather modify their behaviors to match the estimated difference in rewards.

We can also compare the degree of the central tendency in temporal and reward estimation, which will be important in the Results. Empirically, the central tendency in temporal tasks is normally weak. While it can be unmasked in healthy subjects [71–75] and animals [76], it is most evident in unmedicated Parkinson's patients [54], in whom the DA deficiency is profound. This implies a significant asymmetry at baseline: While decreasing the DA levels will have a strong behavioral signature (the central tendency), the effect of increased DA levels will be small (due to a 'ceiling effect,' in which the central tendency will continue to be weak). On the other hand, both increases and decreases to the DA level substantially affect the exploration-exploitation trade-off [24, 58–60, 77]. This suggests a more significant central tendency for rewards at baseline, which can be amplified or mitigated by DA manipulations. Below we will find that DA's effect in ITCs and PDs will depend on its relative contribution to each of the reward estimates and temporal estimates at baseline. Driven by the empirical observations, we take the baseline central tendency to be weaker in the domain of timing than in the domain of rewards.

Finally, it will be useful to distinguish between an animal's true precision and its estimated precision (or what it perceives its precision to be). True precision refers to the precision with which the animal actually encodes the signal. Estimated precision, on the other hand, determines how heavily to weigh the previously encoded signal against the context, as in Eq. (2). The above treatment assumes perfect calibration between encoding and decoding, so that the animal weighs a signal (during decoding) in perfect accordance with its true encoding precision. However, if the neural substrate of precision is the DA level, then it should be possible to elicit certain biases by selectively manipulating the DA level during decoding. Our main predictions will indeed involve tasks in which the DA level was pharmacologically manipulated after training but immediately before testing.

### Decision making under the Bayesian theory

Having estimated the relevant parameters in the task, how does the animal actually use these parameters to make decisions? Under standard models of action selection, the probability of selecting arm  $A_i$  with expected reward  $\hat{\mu}_i$  follows a softmax function [78, 79]:

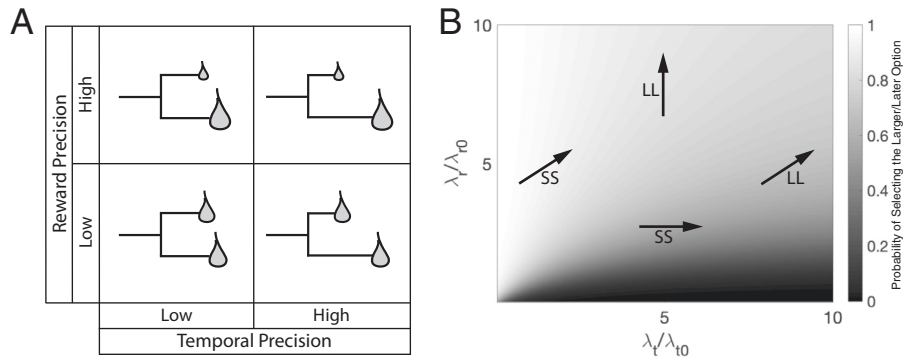
$$p(A_i) = \frac{e^{\beta \hat{\mu}_i}}{\sum_j e^{\beta \hat{\mu}_j}}, \quad (3)$$

where  $\beta$  is the inverse temperature parameter, which controls choice stochasticity. The studies examined in the Results all involve choices between two options; thus, we can restrict our analysis to the case of two arms,  $A_1$  and  $A_2$ , yielding large and small reward, respectively. Furthermore, each arm not only carries a different reward magnitude but also a different delay period between rewards. Thus, the animal must estimate the arms' reward rates  $\bar{R}_1$  and  $\bar{R}_2$  (or ratios of reward magnitude to delay), respectively, in order to maximize its total accumulated reward. Eq. (3) can then be written as

$$p(A_1) = \frac{1}{1 + e^{-\beta(\bar{R}_1 - \bar{R}_2)}}. \quad (4)$$

Notice here that the probability of selecting the option yielding the large reward rate depends on the difference between the reward estimates: As the quantity  $(\bar{R}_1 - \bar{R}_2)$  increases,  $p(A_1)$  increases. Furthermore, by controlling the encoding precisions and thus the central tendencies (either in the temporal or reward domain), DA modulates the estimated difference in posterior means (see horizontal black segments in Fig. 1). A number of authors have argued that DA implements gain control on the values  $\hat{\mu}_i$  in reinforcement learning tasks, possibly by controlling  $\beta$  [60, 66, 69]. The Bayesian theory subsumes the gain control theory by modulating the estimated difference directly.

We have made two important assumptions here. First, our choice rule, though conventional, disregards the contributions of the posterior precisions (i.e., their uncertainties). Recent studies have shown that human behavior in certain bandit tasks is better described by augmented models that incorporate random and directed exploration strategies, both of which make use of the posterior precisions [80–83]. We discuss the augmented model in Supplementary Text 1 and examine its implications for the Bayesian theory.



**Fig. 2 Behavior in ITCs depends on the relative change in reward precision compared to temporal precision.** **A** Schematic illustrating the reward and temporal estimates for each of the smaller/sooner and larger/later options under different reward and temporal precisions. Selectively increasing the reward precision (bottom cells to top cells) masks the reward central tendency, making the difference in reward estimates larger. According to Eq. (4), this promotes selection of the larger/later option. On the other hand, selectively increasing the temporal precision (left cells to right cells) masks the temporal central tendency, making the difference in temporal estimates larger. This promotes selection of the smaller/sooner option. **B** Isolines representing pairs of relative precisions that yield the same probability of selecting the larger/later option under Eq. (4). Note that these isolines have different concavities: In the top left, the isolines are concave up (or convex), whereas in the bottom right, the isolines are concave down. Selectively increasing the reward precision promotes the larger/later option (top arrow), whereas selectively increasing the temporal precision promotes the smaller/sooner option (bottom arrow). Based on empirical findings, we assume that the temporal precision at baseline is high, compared to the baseline reward precision (each normalized by its prior precision). This means that DA's net effect is to promote the larger/later option (right arrow). If, however, the temporal precision is sufficiently reduced, DA's net effect will be to promote the smaller/sooner option (left arrow). Plotted on each axis is the ratio of encoding and prior precisions, which determines the central tendency:  $w = \frac{\lambda}{\lambda + \lambda_0} = (1 + (\frac{\lambda_0}{\lambda})^{-1})^{-1}$ . For illustration, we have chosen  $\mu_r = 1$  and  $4$ , and  $\mu_t = 2$  and  $6$ , for the smaller/sooner and larger/later options, respectively, and  $\beta = 10$ . LL increase in probability of selecting the larger/later option, SS increase in probability of selecting the smaller/sooner option,  $\lambda_t$  temporal encoding precision,  $\lambda_{t0}$  temporal prior precision,  $\lambda_r$  reward encoding precision,  $\lambda_{r0}$  reward prior precision.

Our second assumption is about the *shape* of the posterior distributions. As mentioned above, we have assumed Gaussians throughout, with fixed encoding noise. These assumptions are for convenience: Indeed, the uncertainty about an arm should be lower when the arm is more frequently sampled (more well-learned), higher when it is sampled further in the past, and higher for larger-magnitude stimuli [84–87]. Our results will not depend strongly on either the Gaussian assumption or the absolute magnitude of the encoding uncertainty, but only that the central tendency be sufficiently reduced or strengthened under sufficiently high and low precisions, respectively.

## RESULTS

### Dopamine and intertemporal choice

ITCs involve choosing between a small reward delivered soon, and a large reward delivered later. In these tasks, the smaller/sooner delay is held fixed (and is often zero, resulting in immediate reward), while the larger/later delay is varied across blocks. When the delays are equal, animals will overwhelmingly choose the larger option, but as the delay for the larger option gets longer, animals become more and more likely to choose the smaller/sooner option (Fig. 3). This shift toward the smaller/sooner option has traditionally been explained in terms of reward discounting: The promise of a future reward is less valuable than that same reward delivered immediately, and becomes even less valuable as the delay increases. In other words, future rewards are discounted in proportion to the delay required to receive them. Previous computational models have shown this reward discounting to be well-described by a hyperbolic (or quasi-hyperbolic) function [21, 88].

A competing line of thought is that animals seek to maximize their reward rates (or equivalently, the total accumulated rewards in the task) [48, 49, 51], but are limited by a significant underestimation of the post-reward delays in the task [50]. On this view, animals compute the reward rate for each option—i.e., the undiscounted reward magnitude divided by the total trial time—but base the trial time largely on the pre-reward delay. This

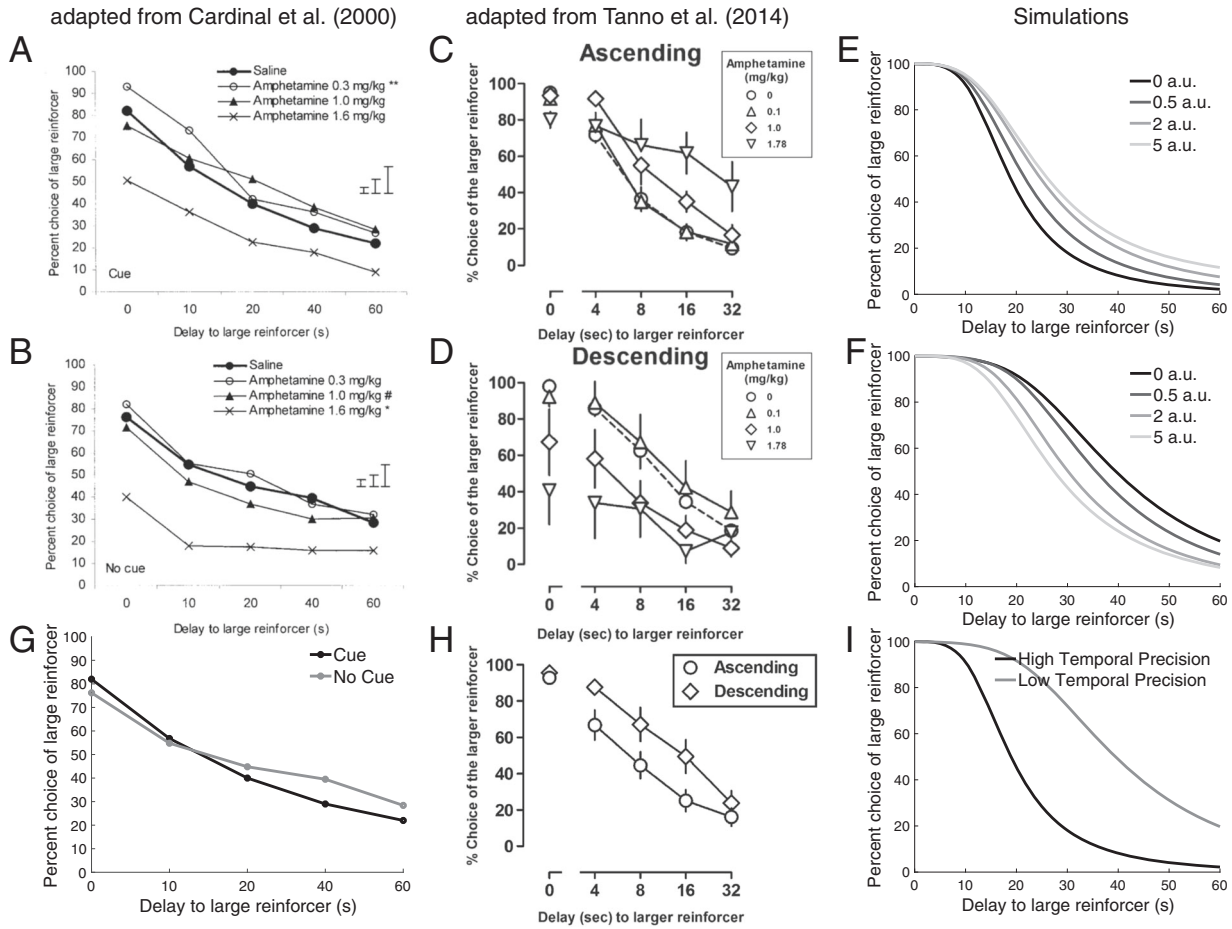
causes the reward rate for the smaller/sooner option to be disproportionately overestimated compared to that of the larger/later option. This view, much like the discounting view, predicts that animals will choose the larger/later option when its delay is short, but will gradually begin to prefer the smaller/sooner option as the delay is increased. Furthermore, the smaller/sooner option will be preferred in some cases even when it yields a lower reward rate, although this is due to a temporal bias (underestimation of post-reward delays), rather than a choice bias (reward discounting). Note here that we use the term ‘discounting’ to refer to the psychological principle that future rewards are valued less than immediate rewards by virtue of the need to wait for them. Thus, even though in the reward-rate view, rewards are divided by their temporal interval, they are not ‘discounted.’

While the reward-rate interpretation can accommodate the aspects of the data explained by the discounting model, it also captures aspects of animal behavior where the discounting model fails. In particular, Blanchard et al. [50] examined the effect of post-reward delays on behavior. Under the discounting model, behavior depends only on the reward magnitudes and *pre*-reward delays (over which the discounting occurs), and thus should be invariant to changes in the post-reward delays. The authors, however, found that monkeys modified their choices in line with a reward-rate computation, which must take into account both pre- and post-reward delays when computing the total trial time. Interestingly, the best fit to the data required that the post-reward delays be underestimated by about a factor of four, consistent with a bias of timing rather than a bias of choice in explaining animal behavior in ITCs. In what follows, we adopt the reward-rate interpretation in examining DA's role in ITCs.

Given DA's effects on reward estimates and durations, it is not surprising that DA would influence behavior in ITCs, where the agent's task is to maximize the ratio of these two, the reward rate  $\bar{R}$ :

$$\bar{R} = \frac{w_r \mu_r + (1 - w_r) \mu_{r0}}{w_t \mu_t + (1 - w_t) \mu_{t0}}, \quad (5)$$





**Fig. 3** DA agonists promote selection of the larger/later option when the temporal precision is high and the smaller/sooner option when the temporal precision is low. **A** Cardinal et al. [16] trained rats on an ITC in which the animals must choose between a reward of magnitude 1 delivered immediately and a reward of magnitude 4 delivered after a delay that varied across blocks. After training, the authors administered DA agonists and examined changes in the animals' behaviors. When a cue was present during the delay period, the authors found that the animals seemed less impulsive under DA agonists, or discounted future rewards less. The 0.3 mg/kg dose, but not the other doses, reached statistical significance (\*\* $p < 0.01$ , main effect of DA agonist). **B** However, when a cue was absent during the delay period, the animals appeared more impulsive with higher doses, i.e., discounted future rewards more strongly ( $*p < 0.05$ , main effect of DA agonist;  $\#p < 0.05$ , agonist-delay interaction). For **A** and **B**, vertical bars denote the standard error of the difference between means for 0.3, 1.0, and 1.6 mg/kg relative to saline, from left to right. **C** Tanno et al. [33] administered a similar task, but varied the order in which the delays were presented. When the delays were presented in an ascending order, the rats seemed less impulsive with higher doses of DA agonists. **D** However, when the delays were presented in a descending order, the rats seemed more impulsive with higher doses. **E** Our model recapitulates these effects: Under high temporal precision, such as in the presence of a visual cue during the delay (cue condition) or as suggested empirically by measuring response variability (ascending condition; Supplementary Text 3), DA's effect on the reward estimates will dominate in ITCs, which promotes selection of the larger/later option. **F** On the other hand, under sufficiently low temporal precision, DA's effect on the temporal estimates will dominate, which promotes selection of the smaller/sooner option. **G** At baseline, responses in the no-cue condition are biased toward the larger/later option compared to the cue condition. Note that any zero-delay difference cannot be due to a difference in the cues, since the tasks are identical in the absence of a delay. It is not clear whether these differences are statistically significant, as error bars were not provided for the saline conditions (although when the conditions were tested immediately before drug administration began, the difference was not statistically significant). Panel reproduced from the saline conditions in **A** and **B**. **H** Similarly, at baseline, responses in the descending condition are biased toward the larger/later option compared to the ascending condition. **I** Our model recapitulates these effects: Selective decreases to the temporal precision promote the larger/later option. For **E**, **F**, and **I**, see Supplementary Text 3 for simulation details. *a.u.* arbitrary units of DA.

which follows from Eq. (2). Here,  $w_r = \frac{\lambda_r}{\lambda_r + \lambda_{r'}}_t$  and  $\mu_{r0}$ ,  $\lambda_{r0}$ ,  $\mu_r$  and  $\lambda_r$  in the numerator represent the prior mean, prior precision, encoding distribution mean, and encoding distribution precision in the domain of rewards, respectively, and similarly for the domain of time in the denominator. Increasing the DA level increases both encoding precisions,  $\lambda_r$  and  $\lambda_r$ .

To understand DA's overall effect on the ratio  $\bar{R}$ , it will be useful to examine manipulations of reward and temporal precision separately. First, let us hold the temporal precisions (and thus the temporal estimates in the denominator) constant. A strong central tendency for the estimated rewards causes an overestimation of

the smaller reward and an underestimation of the larger reward, thus promoting selection of the smaller/sooner option compared to baseline. Because increasing DA masks the central tendency, its effect on the reward estimates in the numerator is to promote selecting the larger/later option (Fig. 2A and B, top arrow). Now let us hold the reward precisions constant. In the denominator, a stronger central tendency for the estimated durations causes an overestimation of the sooner duration and an underestimation of the later duration, thus promoting selection of the larger/later option. Because increasing DA masks the central tendency, its effect on the temporal estimates in the denominator is to promote

selecting the smaller/sooner option—the opposite of its effect in the numerator (Fig. 2A and B, bottom arrow). Thus the ultimate effect of DA will depend on its relative contribution to the reward and temporal estimates (see Supplementary Text 2 for an analytical derivation).

As discussed in the previous section, the central tendency at baseline DA levels is stronger for reward estimates than temporal estimates. It follows that the central tendency in the numerator dominates DA's influence in ITCs (Fig. 2B, right arrow). Under normal conditions, then, the framework predicts that increasing DA will promote the larger/later option, or behavior consistent with lower impulsivity under higher DA levels (Fig. 3E).

This prediction matches well with empirical findings, as the majority of studies have found administering DA agonists to decrease impulsivity in ITCs [15, 22–28] (see [29] for a recent review). For instance, Cardinal et al. [16] trained rats on an ITC involving a small reward delivered immediately and a large reward delivered after a delay that varied across blocks. After training, the authors administered DA agonists and tested the animals on the task. While the effect is smaller than in other studies (e.g., compare with Fig. 3C), the authors found that the DA agonists promoted selection of the larger/later option when a visual cue was present throughout the trial (Fig. 3A).

This prediction is based on the empirically motivated result that DA's effect on the reward estimate dominates its overall effect in ITCs. However, it should be possible to elicit exactly the opposite result—an increased preference for the smaller/sooner option with DA—under conditions where the central tendency of temporal estimates dominates. For instance, timing precision is affected by the inclusion of temporally-informative cues [84, 89, 90] as well as manipulations of the interval salience, presumably due to changes in the animal's alertness [91]. Then removing these cues and selectively decreasing the salience during the delay period should promote the temporal central tendency and, if significant enough, overwhelm the central tendency of rewards in the numerator (Fig. 2, left arrow). Cardinal et al. [16] examined exactly this manipulation: The authors found that DA, on average, promoted selection of the larger/later option only when a salient, temporally-informative visual cue was selectively available during the delay period. Otherwise, DA uncharacteristically had the opposite effect (Fig. 3B), as predicted when the temporal precision is sufficiently reduced (Fig. 3F).

It is important to note that DA manipulations can mediate preference reversals, which is captured by our model. For example, for the 8- and 16-second delays in Fig. 3C, the animal at baseline prefers the smaller/sooner option (chosen more than 50% of the time). But with high enough doses of DA agonists, it eventually comes to prefer the larger/later option. This empirical finding is important because it rules out hypotheses in which DA simply amplifies or mitigates existing preferences. For instance, and as mentioned above, a number of authors have proposed that DA serves a 'gain control' function on the action values during decision making [67–70]. This would predict that preferences should become more extreme with higher DA levels: Preferences above the indifference (50%) line should increase, and those below the indifference line should decrease, which is inconsistent with the empirical results.

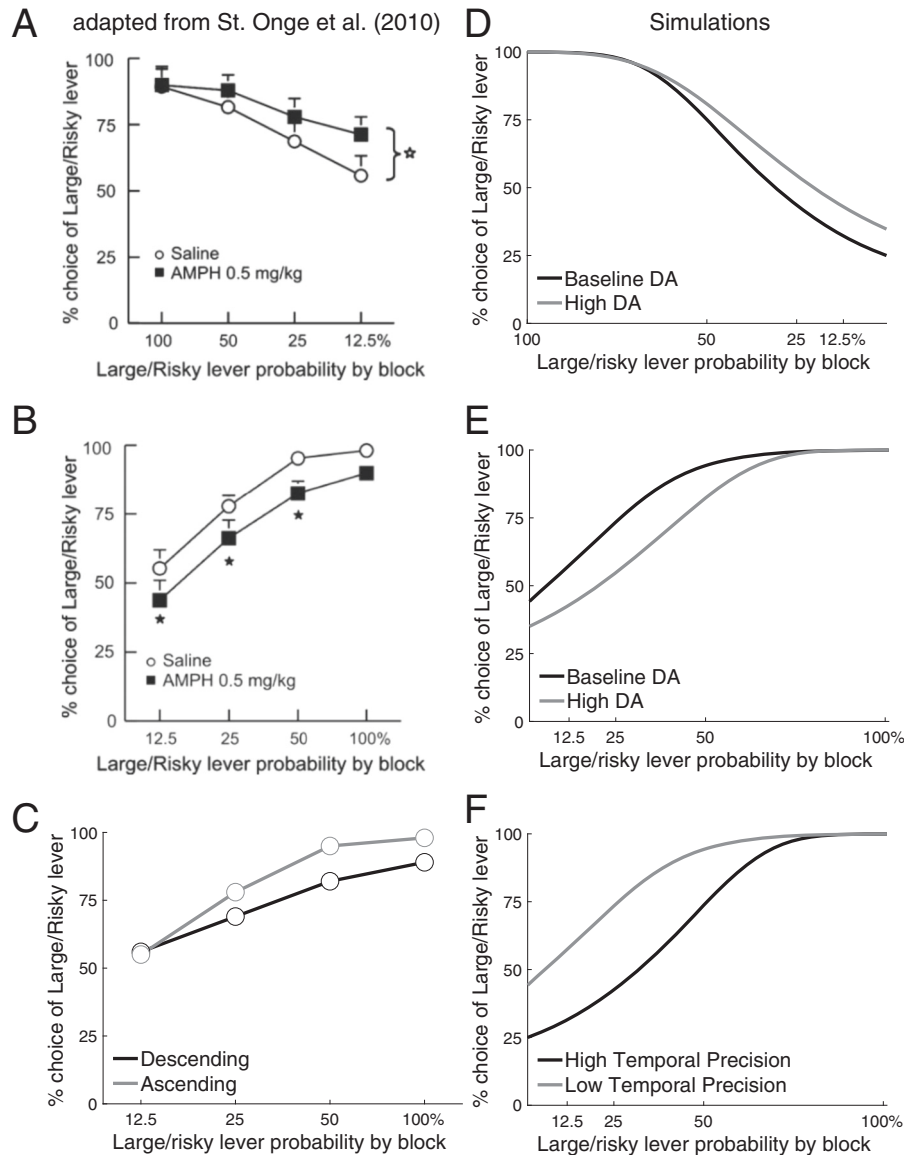
Though the majority of studies have found behaviors consistent with a negative correlation between impulsivity and DA, Cardinal et al. [16] found the opposite effect when the cue was selectively absent during the delay period, and we showed that the Bayesian framework captures this effect. We are aware of one other manipulation that may cause this opposite effect: In tasks where animals are trained on different delays for the larger/later option, Tanno et al. [33] have reported that DA's effect depends on the ordering of the delays. In particular, they found that DA agonists seemed to promote choosing the larger/later option, in line with most other studies, when the delays were presented in an

ascending order. However, if the delays were presented in a descending order, DA agonists had the opposite effect (see also [92]). This finding would be consistent with our framework, if the temporal precision in the ascending case were higher than that in the descending case (Fig. 3C, D). This may indeed be the case, as when learned in an ascending order, the animals' temporal behavior (i.e., the timing of the animal's lever press) was less variable than when learned in a descending order (Supplementary Text 3). An important limitation here is that this result does not control for the animal's motivational state: It is possible that the ordering effect influences not the animal's temporal precision but its motivation, leading to less temporally precise behavior. It is not clear *why* such an ordering effect exists, although one possibility is that this arises from a primacy effect in the inference about the temporal sequence [93–96], as the initial temporal precision is higher for the short delays due to Weber's law [84, 89, 90], and potentially also due to the incentive structure (the animal is more incentivized to attend to blocks with higher reward rates, i.e., those with short delays) [57, 97].

Third, the Bayesian framework makes a counterintuitive prediction about the relationship between baseline performance in ITCs and the effect of DA. According to our model, selectively increasing the temporal precision promotes the smaller/sooner option. However, DA's effect, when the temporal precision is already high, is to promote the *larger/later* option (compare bottom and right arrows in Fig. 2). This implies that conditions in which DA agonists promote the larger/later option will be conditions in which animals are, at baseline, more likely to select the smaller/sooner option. The authors of both studies above indeed observed this relationship: For both the cue and ascending conditions, animals were more likely to select the smaller/sooner option at baseline, compared to the no-cue and descending conditions, respectively (Fig. 3G, H), as predicted (Fig. 3I). Note, however, that this effect may also be due to baseline differences in the speed of the 'internal clock': There is some evidence to suggest that slower clocks are correlated with lower temporal precision [57]. This means that, in tasks with low temporal precision, the animal may perceive the interval to be shorter than in tasks with high temporal precision, which may make it more appealing (shorter intervals result in larger reward rates). We examine this point at length in Supplementary Text 4. Interestingly, combined with the correlation between clock speed and temporal precision, the Bayesian theory makes another prediction: that temporal intervals that are underestimated (slow internal clock) should not be well-learned (low encoding precision, resulting in strongly overlapping posteriors due to the central tendency). As mentioned previously, Blanchard et al. [50] indeed observed that post-reward delays are underestimated and not well-learned, but that this learning can be salvaged if the temporal intervals are made salient, as our model predicts (Supplementary Text 4).

It should be noted that, while our model is concerned with the main effect of DA manipulations, animal response profiles seem also to profoundly diverge in the descending task for the smaller delays (note splaying of response profiles in Fig. 3D). Our model can accommodate this result: Because of noisy learning, the encoding of current estimates will be biased toward estimates from previous blocks, a form of 'within-arm' contextual influence. This effect should be more apparent under higher DA levels, because of a silencing of the central tendency. We expand on this point in Supplementary Text 5.

Finally, having examined DA's effects on behaviors in interval timing and measures of impulsivity, we can also examine how the behavioral phenotypes covary with each other. Our model predicts that—due to natural differences in DA levels within a species—animals that are more precise timers should also appear less impulsive in ITCs, as has indeed been observed [98, 99] (Supplementary Text 6).



**Fig. 4** DA agonists promote selection of the large/risky option when the temporal precision is high and the small/safe option when the temporal precision is low. **A** St Onge et al. [47] trained rats on a PD in which the animals chose between 1 pellet delivered on 100% of trials and 4 pellets delivered with some probability that varied across blocks. After training, the authors administered the DA agonist amphetamine. They found that, when the probabilities were presented in a descending order, amphetamine induced an increase in the tendency to select the large/risky option. **B** However, when the probabilities were presented in an ascending order, amphetamine had the opposite effect—inducing a decrease in the tendency to select the large/risky option. **C** At baseline, responses in the ascending condition are biased toward the large/risky option compared to the descending condition. Panel reproduced from the saline conditions in **A** and **B**. As per Fig. 3G, it is unclear for which doses the differences in choice behavior are statistically significant, although visual inspection suggests a main effect of block order. **D** Our model recapitulates these effects: Under high temporal precision, DA's effect on the reward estimates will dominate, which promotes selection of the large/risky option. **E** On the other hand, under sufficiently low temporal precision, DA's effect on the temporal estimates will dominate, which promotes selection of the small/safe option. **F** Selective decreases to the temporal precision promote the large/risky option. See Supplementary Text 3 for simulation details.

We have sought to highlight here a link between temporal precision and the effect of DA in ITCs. Under the Bayesian theory, reward and temporal estimates normally regress to their contextual means in inverse proportion to the encoding precisions, and increasing the DA level mitigates this regression. This increases both the estimated cost (delay) and benefit (reward) of the larger/later option. When the temporal precision is already high (negligible regression to the temporal mean), the increase in benefit dominates DA's effect, and the animal becomes more likely to select the larger/later option. When the temporal precision is sufficiently low (strong regression), the increase in cost dominates, and the animal shifts its preference toward the

smaller/sooner option. Note here that our focus on temporal, rather than reward, precision is driven by the experimental paradigm: Of the four estimated parameters—small reward, short duration, large reward, and long duration—only the long duration is varied across blocks, so it is not surprising that variations of the ITC would be characterized by different temporal precisions. In principle, a similar analysis can be conducted for manipulations of the reward precisions.

#### Dopamine and probability discounting

We now extend the Bayesian analysis to risk-seeking and PDs. We follow a similar outline to the previous section in showing that,

like for the case of impulsivity and ITCs, the findings of DA in PDs are fully captured by the Bayesian theory, and do not necessarily reflect changes in risk preferences at all.

St Onge and Floresco [46] tested rats on a task in which they had to choose between an arm yielding 1 pellet with 100% probability (the small/safe option) and another arm yielding 4 pellets but with a probability that was varied across blocks (the large/risky option). As the probability decreased, the rats became less likely to choose the large/risky option. After the rats achieved stable performance in each block, the DA agonist amphetamine was administered. The authors found that the DA agonist induced a tendency to select the large/risky option, which was taken to suggest that increasing the DA levels may promote risk-seeking behaviors. However, in a follow-up study [47], the authors found that the ordering of the blocks mattered: The DA agonist only promoted the large/risky option when the probability of reward decreased with blocks. When the order was reversed, the DA agonist had the opposite effect (Fig. 4A, B).

The Bayesian theory predicts this finding. At a conceptual level, this is because the task sets up a trade-off between two attributes (reward magnitude and risk), whose central tendencies push the reward rate in opposite directions: An animal that learns the reward magnitudes well and largely disregards the risks will view the large/risky option as superior, whereas an animal that learns the risks (probabilities) well but disregards the reward magnitudes will prefer the small/safe option.

To evaluate the theory's predictions more concretely, we note that there are a number of ways for an animal to learn the reward rate, and therefore a number of ways to set up the two central tendencies. For instance, the animal may estimate the reward magnitude and reward probability separately and take their product. Alternatively, the animal may estimate the reward magnitude and average delay time between two rewards for each option, and take their ratio, as in Eq. (5). While both approaches can accommodate the empirical results, we will assume the latter for two reasons: First, this approach allows for a direct comparison with ITCs, as delay discounting and probability discounting elicit similar behaviors [100, 101], similar types of intolerance [102, 103], and a common neural substrate [44, 104]. Second, after Tanno et al. [33], if the animal is indeed computing the ratio of reward magnitude to the delay, then the block order manipulation makes a clear prediction about the temporal precision: A block order involving sequentially longer delays (here, the descending condition) results in higher temporal precision than an order involving sequentially shorter delays (here, the ascending condition). As was the case for the ITC, this will be the key variable determining DA's overall effect.

Thus, to maximize its reward rate, we assume an animal estimates the reward magnitude and average delay between rewards for each option, and computes their ratio. For the reward magnitudes, the central tendency promotes selection of the small/safe option. Increasing DA silences this effect, and thus promotes selection of the large/risky option. On the other hand, for the temporal interval, the central tendency promotes selection of the large/risky option, which involves a longer waiting time between rewards. Increasing DA therefore promotes selection of the small/safe option here—the opposite of its effect for reward magnitudes. Thus, once again, DA's overall effect depends on its relative contribution to each term. Following the previous section, the temporal precision is predicted to be low in the ascending condition (which involves sequentially shorter delays) compared to the descending condition (which involves sequentially longer delays). It follows that the temporal central tendency dominates in the ascending condition and the reward central tendency dominates in the descending condition, thus predicting the empirical findings (Fig. 4D, E).

Our theory also predicts that selectively increasing the temporal precision, without manipulating DA, should also make animals

more likely to select the small/safe option. Indeed, by examining the baseline (saline) task for each of the ascending and descending conditions, we find that animals were more likely to select the small/safe option in the descending condition at baseline compared to the ascending condition (Fig. 4C), as predicted (Fig. 4F).

Finally, having examined temporal and probability discounting separately, we briefly mention the 'rat gambling task,' a four-armed bandit task in which both the delay periods and reward probabilities (in addition to the reward magnitudes) are varied across arms. A well-replicated finding in this task has been that the arm yielding the largest reward magnitude is the second most frequently selected arm, even though it yields the lowest reward rate. In addition, amphetamine administration tends to disrupt selection of the arm yielding the highest reward rate in favor of that yielding the second highest reward rate [105]. The Bayesian theory can accommodate both findings (Supplementary Text 7).

## DISCUSSION

We have shown here that DA's effects in ITCs and PDs are well-described by a Bayesian framework in which animals maximize their reward rates. Under this view, DA controls the relative influence of context in computing the reward and temporal estimates, whose ratio forms the reward rate. Notably, the discounting-free model successfully predicts that DA agonists should promote selection of the larger/longer and large/risky options under high temporal precision, but should have exactly the opposite effects when the temporal precision is sufficiently low. The Bayesian view thus provides a principled framework for why DA would appear to inhibit impulsive and risky choices in some paradigms but promote them in others.

We have followed previous theoretical and experimental work in adopting a discounting-free model of choice behavior. However, our results do not necessarily rule out temporal or probability discounting more generally, nor a role for DA in these processes. For instance, and as mentioned in the Introduction, humans tend to prefer smaller/sooner options even in the absence of repeated trials that make reward-rate computations meaningful. But why discount future rewards in the first place? One influential hypothesis from economics is that future rewards are discounted because of the risks involved in the delay [106]. For example, a competitor may reach the reward first, or a predator may interfere in the animal's plans to collect the reward. As the delay increases, these alternative events become more likely, and the expected reward (the average over all alternatives) decreases. Another idea is that subjects respond *as if* they will have repeated opportunities to engage in the same task [107], thus mimicking the reinforcement learning problem that defines the animal variant of ITCs. More recently, Gabaix and Laibson [108] have argued that reward discounting may be due to the simulation noise involved in mentally projecting into the future: With later rewards, subjects must mentally simulate further into the future, so the simulation noise increases, and the precision decreases. Assuming a Bayesian framework with a prior centered at zero, the reward estimates will be closer to zero when rewards are more distant in the future, i.e., rewards are discounted with time (see also [109] for an extension of this hypothesis).

Interestingly, as mentioned in the Introduction, DA seems to have the opposite effect in the human variant of the task than in the majority of animal experiments, with a promotion of the smaller/sooner option with higher DA levels. That DA may serve a qualitatively different function in the human variant is not completely unexpected, given the substantial differences in the experimental paradigms. Notably, in the human variant, (1) the subject does not actually experience the pre-reward delay, (2) there is no post-reward delay, (3) the subject does not necessarily receive an actual reward, (4) the subject may experience a single



trial of this task, whereas animals are trained on many trials, and (5) the hypothetical delay is on the order of days (or months) and not seconds. Experience and repetitions may prove critical for our reinforcement learning task, and delays on the order of days engage different timing mechanisms than those on the order of seconds-to-minutes [110], which is the duration over which DA's central tendency effect has been observed. Nonetheless, the human findings may still be reconcilable with our framework under the 'repeated opportunities' hypothesis of Myerson and Green [107] mentioned above: It is possible that the temporal uncertainty surrounding durations that are not experienced, and that are on the order of days, is large and thereby dominates DA's central tendency effects. Thus DA agonists would be predicted to promote the smaller/sooner option.

Our framework leaves open a number of theoretical and empirical questions. First, our model takes DA to control the encoding precision, a property inherited from the Bayesian timing model of DA and further motivated by theories of DA as overcoming the cost of attention [57, 97]. However, our results only require that DA control the ratio of the encoding precision to the prior precision but not necessarily the encoding precision itself. Instead, it is certainly possible that increasing DA decreases the prior precision, as some authors have proposed [53]. Interestingly, this ambiguity is not specific to theories of DA, and has been a point of debate for some Bayesian theories of autism as well (compare weak priors [111] with strong likelihoods [112]).

A second open question concerns our assumption that estimates of the reward magnitude are biased by a central tendency effect. Thus far, this has been inferred mainly from exploration-exploitation paradigms (see [64] for a more direct examination), but a dopaminergic modulation of reward estimates has not, to our knowledge, been observed directly. Driven by the experimental literature, we have therefore focused our simulations on manipulations of *temporal* precision. Our work then opens the door to a fruitful line of experiments with novel predictions: For instance, one can develop ITCs and PDs where the large reward is varied rather than the delay or risk. Our framework predicts that DA agonists will promote the larger/late and large/risky options only when *reward* precision is low at baseline, and the smaller/sooner and small/safe options when reward precision is high. On the other hand, selectively increasing the reward precision will always promote the larger/late and large/risky options (Fig. 2). Thus, once again, by simply controlling the central tendency, DA agonists will appear to inhibit impulsivity and risk-seeking under some conditions, but promote them in others.

Third, having adopted an algorithmic view of DA's function, it remains for future work to ask how the Bayesian theory is actually implemented neurobiologically. Notably, DA exerts different effects depending on the postsynaptic receptor subtype: In reinforcement learning studies, midbrain DA neurons project to the striatum onto neurons primarily expressing either D1 or D2 receptors, which segregate anatomically into largely separate basal ganglia pathways [113] and seem to serve opposite purposes [114, 115], both in their fast- [116] and slow-timescale [69, 117, 118] activities. Asymmetries in receptor-mediated effects extend into interval timing studies (compare D1-mediated [119–123] with D2-mediated [119, 120] effects), and DA's effects also depend on enzymatic activity [124, 125] and projection site [126, 127]. Bridging the algorithmic and implementational levels of the Bayesian theory will be a necessary next step toward a more complete theory of DA.

Finally, we have examined in this work how manipulations of DA affect behavior in the ITC and PD, but it is interesting to ask what variables determine the DA level in the first place. An influential proposal has been that the tonic DA level is set by the average reward availability in an environment [128], as has also

been suggested empirically [129]. One unifying interpretation of the average reward and Bayesian theories is that, in high reward-rate environments, animals are more incentivized to attend to a task, and thus encode the relevant parameters with higher precision. In this manner, DA connects the encoding stage (learning the parameters) with the decoding stage (combining the learned parameters with contextual information, the focus of this paper) [57]. It will remain for future work to build upon and experimentally validate coherent theories of DA within an encoding-decoding framework.

To our knowledge, this is the first framework that can accommodate the seemingly conflicting effects of DA in measures of impulsive choice and risk-seeking across experimental conditions. Nonetheless, our aim throughout this work is not to rule out a role for DA in true impulsivity and risk-seeking, but rather to show how a single Bayesian framework can accommodate a wide range of otherwise perplexing behavioral and pharmacological phenomena.

## INSTITUTIONAL REVIEW BOARD

This is a theoretical study which does not describe any new data.

## CODE AVAILABILITY

Source code for all simulations can be found at [www.github.com/jgmikhael/bayesiantheory](https://www.github.com/jgmikhael/bayesiantheory).

## REFERENCES

1. American Psychiatric Association. Diagnostic and statistical manual of mental disorders (DSM-5®). American Psychiatric Pub, 2013. <https://doi.org/10.1176/appi.books.9780890425596>.
2. Tannock R, Schachar RJ, Carr RP, Chajczyk D, Logan GD. Effects of methylphenidate on inhibitory control in hyperactive children. *J Abnorm Child Psychol*. 1989;17:473–491.
3. Gillberg C, Melander H, von Knorring A-L, Janols L-O, Thernlund G, Hägglöf B, et al. Long-term stimulant treatment of children with attention-deficit hyperactivity disorder symptoms: a randomized, double-blind, placebo-controlled trial. *Arch Gen Psychiatry*. 1997;54:857–864.
4. Findling RL, Dogin JW. Psychopharmacology of ADHD: children and adolescents. *J Clin Psychiatry*. 1998;59:42–49.
5. Solanto MV. Neuropsychopharmacological mechanisms of stimulant drug action in attention-deficit hyperactivity disorder: a review and integration. *Behavioural Brain Res*. 1998;94:127–152.
6. Shiels K, Hawk LW Jr, Reynolds B, Mazzullo RJ, Rhodes JD, Pelham WE Jr, et al. Effects of methylphenidate on discounting of delayed rewards in attention deficit/hyperactivity disorder. *Exp Clin Psychopharmacol*. 2009;17:291.
7. Heerey EA, Robinson BM, McMahon RP, Gold JM. Delay discounting in schizophrenia. *Cogn Neuropsychiatry*. 2007;12:213–221.
8. Gold JM, Waltz JA, Prentice KJ, Morris SE, Heerey EA. Reward processing in schizophrenia: a deficit in the representation of value. *Schizophrenia Bull*. 2008;34:835–847.
9. Volkow ND, Fowler JS, Wang G-J, Swanson JM, Telang F. Dopamine in drug abuse and addiction: results of imaging studies and treatment implications. *Arch Neurol*. 2007;64:1575–1579.
10. Bickel WK, Jarmolowicz DP, Mueller ET, Koffarnus MN, Gatchalian KM. Excessive discounting of delayed reinforcers as a trans-disease process contributing to addiction and other disease-related vulnerabilities: emerging evidence. *Pharmacol Therapeutics*. 2012;134:287–297.
11. Dagher A, Robbins TW. Personality, addiction, dopamine: insights from Parkinson's disease. *Neuron*. 2009;61:502–510.
12. O'Sullivan SS, Evans AH, Lees AJ. Dopamine dysregulation syndrome. *CNS Drugs*. 2009;23:157–170.
13. de Wit H, Enggasser JL, Richards JB. Acute administration of d-amphetamine decreases impulsivity in healthy volunteers. *Neuropsychopharmacology*. 2002;27:813–825.
14. Pine A, Shiner T, Seymour B, Dolan RJ. Dopamine, time, and impulsivity in humans. *J Neurosci*. 2010;30:8888–8896.
15. Wade TR, de Wit H, Richards JB. Effects of dopaminergic drugs on delayed reward as a measure of impulsive behavior in rats. *Psychopharmacology*. 2000;150:90–101.

16. Cardinal RN, Robbins TW, Everitt BJ. The effects of d-amphetamine, chlordiazepoxide,  $\alpha$ -flupenthixol and behavioural manipulations on choice of signalled and unsignalled delayed reinforcement in rats. *Psychopharmacology*. 2000;152:362–375.
17. Jeffrey R Stevens and David W Stephens. The adaptive nature of impulsivity. American Psychological Association, 2010. <https://doi.org/10.1037/12069-013>.
18. Rachlin H, Green L. Commitment, choice and self-control 1. *J Exp Anal Behav*. 1972;17:15–22.
19. Ainslie G. Specious reward: a behavioral theory of impulsiveness and impulse control. *Psychological Bull*. 1975;82:463.
20. Tobin H, Logue AW. Self-control across species (*Columba livia*, *Homo sapiens*, and *Rattus norvegicus*). *J Comp Psychol*. 1994;108:126.
21. Howard Rachlin. The science of self-control. Harvard University Press, 2000. <https://books.google.com/books?id=D1gwEAAAQBAJ&dq=+Howard+Rachlin.+The+science+of+self-control.+Harvard+University+Press,+2000.&lr=>.
22. Winstanley CA, Dalley JW, Theobald DEH, Robbins TW. Global 5-HT depletion attenuates the ability of amphetamine to decrease impulsive choice on a delay-discounting task in rats. *Sychopharmacology*. 2003;170:320–331.
23. Denk F, Walton ME, Jennings KA, Sharp T, Rushworth MFS, Bannerman DM. Differential involvement of serotonin and dopamine systems in cost-benefit decisions about delay or effort. *Psychopharmacology*. 2005;179:587–596.
24. van Gaalen MM, van Koten R, Schoffelmeer ANM, Vanderschuren LJM. Critical involvement of dopaminergic neurotransmission in impulsive decision making. *Biol Psychiatry*. 2006;60:66–73.
25. Winstanley CA, LaPlant Q, Theobald DEH, Green TA, Bachtell RK, Perrotti LI, et al.  $\Delta$  FosB induction in orbitofrontal cortex mediates tolerance to cocaine-induced cognitive dysfunction. *J Neurosci*. 2007;27:10497–10507.
26. Floresco SB, Maric TL, Ghods-Sharifi S. Dopaminergic and glutamatergic regulation of effort-and delay-based decision making. *Neuropsychopharmacology*. 2008;33:1966–1979.
27. Koffarnus MN, Newman AH, Grundt P, Rice KC, Woods JH. Effects of selective dopaminergic compounds on a delay discounting task. *Behavioural Pharmacol*. 2011;22:300.
28. Paterson NE, Wetzler C, Hackett A, Hanania T. Impulsive action and impulsive choice are mediated by distinct neuropharmacological substrates in rat. *Int J Neuropsychopharmacol*. 2012;15:1473–1487.
29. D'Amour-Horvat V, Leyton M. Impulsive actions and choices in laboratory animals and humans: effects of high vs. low dopamine states produced by systemic treatments given to neurologically intact subjects. *Front Behav Neurosci*. 2014;8:432.
30. Charrier D, Thiebot MH. Effects of psychotropic drugs on rat responding in an operant paradigm involving choice between delayed reinforcers. *Pharmacol Biochem Behav*. 1996;54:149–157.
31. Evenden JL, Ryan CN. The pharmacology of impulsive behaviour in rats: the effects of drugs on response choice with varying delays of reinforcement. *Psychopharmacology*. 1996;128:161–170.
32. Helms CM, Reeves JM, Mitchell SH. Impact of strain and D-amphetamine on impulsivity (delay discounting) in inbred mice. *Sychopharmacology*. 2006;188:144–151.
33. Tanno T, Maguire DR, Henson C, France CP. Effects of amphetamine and methylphenidate on delay discounting in rats: interactions with order of delay presentation. *Psychopharmacology*. 2014;231:85–95.
34. Molina J, Sáinz-Artiga MJ, Fraile A, Jiménez-Jiménez FJ, Villanueva C, Orti-Pareja M, et al. Pathologic gambling in Parkinson's disease: a behavioral manifestation of pharmacologic treatment? *Mov Disord*. 2000;15:869–872.
35. Mimura M, Oeda R, Kawamura M. Impaired decision-making in Parkinson's disease. *Parkinsonism Relat Disord*. 2006;12:169–175.
36. Gallagher DA, O'Sullivan SS, Evans AH, Lees AJ, Schrag A. Pathological gambling in Parkinson's disease: risk factors and differences from dopamine dysregulation. An analysis of published case series. *Mov Disord*. 2007;22:1757–1763.
37. Clark CA, Dagher A. The role of dopamine in risk taking: a specific look at Parkinson's disease and gambling. *Front Behav Neurosci*. 8, 2014;22:196.
38. Cheng GLF, Tang JCY, Li FWS, Lau EYY, Lee TMC. Schizophrenia and risk-taking: impaired reward but preserved punishment processing. *Schizophrenia Res*. 2012;136:122–127.
39. Reddy LF, Lee J, Davis MC, Althuler L, Glahn DC, Miklowitz DJ, et al. Impulsivity and risk taking in bipolar disorder and schizophrenia. *Neuropsychopharmacology*. 2014;39:456–463.
40. Grall-Bronnec M, Wainstein L, Augy J, Bouju G, Feuillet F, Vénisse J-L, et al. Attention deficit hyperactivity disorder among pathological and at-risk gamblers seeking treatment: a hidden disorder. *Eur Addiction Res*. 2011;17:231–240.
41. Voon V, Gao J, Brezing C, Symmonds M, Ekanayake V, Fernandez H, et al. Dopamine agonists and risk: impulse control disorders in Parkinson's disease. *Brain*. 2011;134:1438–1446.
42. Sevy S, Hassoun Y, Bechara A, Yechiam E, Napolitano B, Burdick K, et al. Emotion-based decision-making in healthy subjects: short-term effects of reducing dopamine levels. *Psychopharmacology*. 2006;188:228–235.
43. Mobini S, Chiang T-J, Ho M-Y, Bradshaw CM, Szabadi E. Effects of central 5-hydroxytryptamine depletion on sensitivity to delayed and probabilistic reinforcement. *Psychopharmacology*. 2000;152:390–397.
44. Cardinal RN, Howes NJ. Effects of lesions of the nucleus accumbens core on choice between small certain rewards and large uncertain rewards in rats. *BMC Neurosci*. 2005;6:37.
45. Adriani W, Laviola G. Delay aversion but preference for large and rare rewards in two choice tasks: implications for the measurement of self-control parameters. *BMC Neurosci*. 2006;7:1–11.
46. St Onge JR, Floresco SB. Dopaminergic modulation of risk-based decision making. *Neuropsychopharmacology*. 2009;34:681–697.
47. St Onge JR, Chiu YC, Floresco SB. Differential effects of dopaminergic manipulations on risky choice. *Psychopharmacology*. 2010;211:209–221.
48. Kacelnik A. Normative and descriptive models of decision making: time discounting and risk sensitivity. In CIBA foundation symposium, pages 51–70, Wiley Online Library; 1997. <https://books.google.com/books?hl=en&lr=&id=fw4GboUy7aAC&oi=fnd&pg=PA51&dq=Alex+Kacelnik.+Normative+and+descriptive+models+of+decision+making:+time+discounting+and+risk+sensitivity.+In+CIBA+foundation+>.
49. Daw ND, Touretzky DS. Behavioral considerations suggest an average reward TD model of the dopamine system. *Neurocomputing*. 2000;32:679–684.
50. Blanchard TC, Pearson JM, Hayden BY. Postreward delays and systematic biases in measures of animal temporal discounting. *Proc Natl Acad Sci*. 2013;110:15491–15496.
51. K Nambodiri VM, Mihalas S, Marton T, Hussain Shuler MG. A general theory of intertemporal decision-making and the perception of time. *Front Behav Neurosci*. 2014;8:61.
52. Friston KJ, Shiner T, FitzGerald T, Galea JM, Adams R, Brown H, et al. Dopamine, affordance and active inference. *PLoS Comput Biol*. 2012;8:e1002327.
53. Costa VD, Tran VL, Turchi J, Averbeck BB. Reversal learning and dopamine: a Bayesian perspective. *J Neurosci*. 2015;35:2407–2416.
54. Chara M, Rakitin B, Levy R, Meck WH, Deweer B, Dubois B, et al. Coupled temporal memories in Parkinson's disease: a dopamine-related dysfunction. *J Cogn Neurosci*. 1998;10:316–331.
55. Malapani C, Deweer B, Gibbon J. Separating storage from retrieval dysfunction of temporal memory in Parkinson's disease. *J Cogn Neurosci*. 2002;14:311–322.
56. Shi Z, Church RM, Meck WH. Bayesian optimization of time perception. *Trends Cogn Sci*. 2013;17:556–564.
57. Mikhael JG, Lai L, Gershman SJ. Rational inattention and tonic dopamine. *PLoS Comput Biol*. 2021;17:e1008659.
58. Eisenegger C, Naef M, Linssen A, Clark L, Gandamaneni PK, Müller U, et al. Role of dopamine D2 receptors in human reinforcement learning. *Neuropsychopharmacology*. 2014;39:2366.
59. Lee E, Seo M, Dal Monte O, Averbeck BB. Injection of a dopamine type 2 receptor antagonist into the dorsal striatum disrupts choices driven by previous outcomes, but not perceptual inference. *J Neurosci*. 2015;35:6298–6306.
60. Cinotti F, Fresno V, Aklil N, Coutureau E, Girard B, Marchand AR, et al. Dopamine blockade impairs the exploration-exploitation trade-off in rats. *Sci Rep*. 2019;9:6770.
61. Beeler JA, Daw ND, Frazier CRM, Zhuang X. Tonic dopamine modulates exploitation of reward learning. *Front Behav Neurosci*. 2010;4:170.
62. Hamid AA, Pettibone JR, Mabrouk OS, Hetrick VL, Schmidt R, Vander Weele CM, et al. Mesolimbic dopamine signals the value of work. *Nat Neurosci*. 2016;19:117–126.
63. Faisal AA, Selen LPJ, Wolpert DM. Noise in the nervous system. *Nat Rev Neurosci*. 2008;9:292–303.
64. Gershman SJ, Niv Y. Novelty and inductive generalization in human reinforcement learning. *Top Cogn Sci*. 2015;7:391–415.
65. Stojić H, Schulz E, Analytis PP, Speekenbrink M. It's new, but is it good? How generalization and uncertainty guide the exploration of novel options. *J Exp Psychol Gen*. 2020;149:1878–1907.
66. Humphries MD, Khamassi M, Gurney K. Dopaminergic control of the exploration-exploitation trade-off via the basal ganglia. *Front Neurosci*. 2012;6:9.
67. Cohen JD, Servan-Schreiber D. Context, cortex, and dopamine: a connectionist approach to behavior and biology in schizophrenia. *Psychological Rev*. 1992;99:45.
68. Braver TS, Cohen JD, Servan-Schreiber D. A computational model of prefrontal cortex function. In *Advances in Neural Information Processing Systems*, pages 141–148, 1995. <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.1018.3084&rep=rep1&type=pdf>.

69. Collins AGE, Frank MJ. Opponent actor learning (OpAL): modeling interactive effects of striatal dopamine on reinforcement learning and choice incentive. *Psychological Rev.* 2014;121:337.
70. Averbeck BB, Costa VD. Motivational neural circuits underlying reinforcement learning. *Nat Neurosci.* 2017;20:505.
71. Jazayeri M, Shadlen MN. Temporal context calibrates interval timing. *Nat Neurosci.* 2010;13:1020.
72. Acerbi L, Wolpert DM, Vijayakumar S. Internal representations of temporal statistics and feedback calibrate motor-sensory interval timing. *PLoS Comput Biol.* 2012;8:e1002771.
73. Bausenhardt KM, Dyjas O, Ulrich R. Temporal reproductions are influenced by an internal reference: Explaining the Vierordt effect. *Acta Psychologica.* 2014;147:60–67.
74. Mayer KM, Di Luca M, Ernst MO. Duration perception in crossmodally-defined intervals. *Acta psychologica.* 2014;147:2–9.
75. Roach NW, McGraw PV, Whitaker DJ, Heron J. Generalization of prior information for rapid Bayesian time estimation. *Proc Natl Acad Sci USA.* 2017;114:412–417.
76. De Corte BJ, Matell MS. Temporal averaging across multiple response options: insight into the mechanisms underlying integration. *Anim Cognition.* 2016;19:329–342.
77. Pessiglione M, Seymour B, Flandin G, Dolan RJ, Frith CD. Dopamine-dependent prediction errors underpin reward-seeking behaviour in humans. *Nature.* 2006;442:1042.
78. Shepard RN. Stimulus and response generalization: tests of a model relating generalization to distance in psychological space. *J Exp Psychol.* 1958;55:509.
79. Luce RD. Individual choice behavior: a theoretical analysis. John Wiley and Sons; 1959. [https://books.google.com/books?hl=en&lr=&id=ERQsKkPIkkC&oi=fnd&pg=PP1&dq=R.+Duncan+Luce.+Individual+Choice+Behavior:+a+Theoretical+Analysis&ots=2jou\\_s\\_b7m&sig=6lCKcwp8EQlibJ7jCPKlwNpw2w#v=onep](https://books.google.com/books?hl=en&lr=&id=ERQsKkPIkkC&oi=fnd&pg=PP1&dq=R.+Duncan+Luce.+Individual+Choice+Behavior:+a+Theoretical+Analysis&ots=2jou_s_b7m&sig=6lCKcwp8EQlibJ7jCPKlwNpw2w#v=onep).
80. Wilson RC, Geana A, White JM, Ludvig EA, Cohen JD. Humans use directed and random exploration to solve the explore–exploit dilemma. *J Exp Psychol.* 2014;143:2074.
81. Gershman SJ. Deconstructing the human algorithms for exploration. *Cognition.* 2018;173:34–42.
82. Gershman SJ, Tzovaras BG. Dopaminergic genes are associated with both directed and random exploration. *Neuropsychologia.* 2018;120:97–104.
83. Gershman SJ. Uncertainty and exploration. *Decision.* 2019;6:277.
84. Gibbon J. Scalar expectancy theory and Weber's law in animal timing. *Psychological Rev.* 1977;84:279.
85. Killeen PR, Cate H, Tran T. Scaling pigeons' choice of feeds: bigger is better. *J Exp Anal Behav.* 1993;60:203–217.
86. Bateson M, Kacelnik A. Accuracy of memory for amount in the foraging starling, *Sturnus vulgaris*. *Anim Behav.* 1995;50:431–443.
87. Kacelnik A, Bateson M. Risky theories—the effects of variance on foraging decisions. *Am Zool.* 1996;36:402–434.
88. Mazur JE. An adjusting procedure for studying delayed reinforcement. *Commons, ML; Mazur, JE; Nevin, JA*, pages 55–73, 1987. <https://books.google.com/books?hl=en&lr=&id=1q5mAgAAQBAJ&oi=fnd&pg=PA55&dq=James+E+Mazur.+An+adjusting+procedure+for+studying+delayed+reinforcement.+Commons,+ML.%3B+Mazur,+JE.%3B+Nevin,+JA,+>
89. Church RM, Meck M. A concise introduction to scalar timing theory. In *Functional and Neural Mechanisms of Interval Timing*, pages 3–22, 2003.
90. Staddon JER. Some properties of spaced responding in pigeons. *J Exp Anal Behav.* 1965;8:19–28.
91. Li Q, Liu P, Huang S, Huang X. The effect of phasic alertness on temporal precision. *Atten Percept Psychophys.* 2018;80:262–274.
92. Maguire DR, Henson C, France CP. Effects of amphetamine on delay discounting in rats depend upon the manner in which delay is varied. *Neuropharmacology.* 2014;87:173–179.
93. Sanborn AN, Griffiths TL, Navarro DJ. Rational approximations to rational models: alternative algorithms for category learning. *Psychological Rev.* 2010;117:1144.
94. Levy RP, Reali F, Griffiths TL. Modeling the effects of memory on human online sentence processing with particle filters. In *Advances in Neural Information Processing Systems*, pages 937–944, 2009. <http://173.236.226.255/tom/papers/sentencep1.pdf>.
95. Abbott JT, Griffiths TL. Exploring the influence of particle filter parameters on order effects in causal learning. In *Proceedings of the Annual Meeting of the Cognitive Science Society*, volume 33, 2011. <https://escholarship.org/content/qt3k26g6wn/qt3k26g6wn.pdf>.
96. Thaker P, Tenenbaum JB, Gershman SJ. Online learning of symbolic concepts. *J Math Psychol.* 2017;77:10–20.
97. Manohar SG, Chong TT-J, Apps MAJ, Batla A, Stamelou M, Jarman PR, et al. Reward pays the cost of noise reduction in motor and cognitive control. *Curr Biol.* 2015;25:1707–1716.
98. Marshall AT, Smith AP, Kirkpatrick K. Mechanisms of impulsive choice: I. individual differences in interval timing and reward processing. *J Exp Anal Behav.* 2014;102:86–101.
99. McClure J, Podos J, Richardson HN. Isolating the delay component of impulsive choice in adolescent rats. *Front Integr Neurosci.* 2014;8:3.
100. Richards JB, Zhang L, Mitchell SH, De Wit H. Delay or probability discounting in a model of impulsive behavior: effect of alcohol. *J Exp Anal Behav.* 1999;71:121–143.
101. Ho M-Y, Mobini S, Chiang T-J, Bradshaw CM, Szabadi E. Theory and method in the quantitative analysis of "impulsive choice" behaviour: implications for psychopharmacology. *Psychopharmacology.* 1999;146:362–372.
102. Mazur JE. Conditioned reinforcement and choice with delayed and uncertain primary reinforcers. *J Exp Anal Behav.* 1995;63:139–150.
103. Rachlin H, Raineri A, Cross D. Subjective probability and delay. *J Exp Anal Behav.* 1991;55:233–244.
104. Cardinal RN, Cheung THC. Nucleus accumbens core lesions retard instrumental learning and performance with delayed reinforcement in the rat. *BMC Neurosci.* 2005;6:9.
105. Zeeb FD, Robbins TW, Winstanley CA. Serotonergic and dopaminergic modulation of gambling behavior as assessed using a novel rat gambling task. *Neuropsychopharmacology.* 2009;34:2329–2343.
106. Samuelson PA. A note on measurement of utility. *Rev Econ Stud.* 1937;4:155–161.
107. Myerson J, Green L. Discounting of delayed rewards: models of individual choice. *J Exp Anal Behav.* 1995;64:263–276.
108. Gabaix X, Laibson D. Myopia and discounting. Technical Report, National Bureau of Economic Research, 2017.
109. Gershman SJ, Bhui R. Rationally inattentive intertemporal choice. *Nat Commun.* 2020;11:1–8.
110. Buhusi CV, Meck WH. What makes us tick? Functional and neural mechanisms of interval timing. *Nat Rev Neurosci.* 2005;6:755–765.
111. Pellicano E, Burr D. When the world becomes 'too real': a Bayesian explanation of autistic perception. *Trends Cogn Sci.* 2012;16:504–510.
112. Lawson RP, Rees G, Friston KJ. An aberrant precision account of autism. *Front Hum Neurosci.* 2014;8:302.
113. Graybiel AM. The basal ganglia. *Curr Biol.* 2000;10:R509–R511.
114. Albin RL, Young AB, Penney JB. The functional anatomy of basal ganglia disorders. *Trends Neurosci.* 1989;12:366–375.
115. DeLong MR. Primate models of movement disorders of basal ganglia origin. *Trends Neurosci.* 1990;13:281–285.
116. Shen W, Flajolet M, Greengard P, Surmeier DJ. Dichotomous dopaminergic control of striatal synaptic plasticity. *Science.* 2008;321:848–851.
117. Dreyer JK, Herrik KF, Berg RW, Hounsgaard JD. Influence of phasic and tonic dopamine release on receptor activation. *J Neurosci.* 2010;30:14273–14283.
118. Frank M. Dynamic dopamine modulation in the basal ganglia: a neurocomputational account of cognitive deficits in medicated and nonmedicated Parkinsonism. *Cogn Neurosci, J.* 2005;17:51–72.
119. Drew MR, Fairhurst S, Malapani C, Horvitz JC, Balsam PD. Effects of dopamine antagonists on the timing of two intervals. *Pharmacol Biochem Behav.* 2003;75:9–15.
120. Drew MR, Simpson EH, Kellendonk C, Herzberg WG, Lipatova O, Fairhurst S, et al. Transient overexpression of striatal D2 receptors impairs operant motivation and interval timing. *J Neurosci.* 2007;27:7731–7739.
121. Coull JT, Cheng R-K, Meck WH. Neuroanatomical and neurochemical substrates of timing. *Neuropsychopharmacology.* 2011;36:325.
122. Narayanan NS, Land BB, Solder JE, Deisseroth K, DiLeone RJ. Prefrontal D1 dopamine signaling is required for temporal control. *Proc Natl Acad Sci USA.* 2012;109:20726–20731.
123. De Corte BJ, Wagner LM, Matell MS, Narayanan NS. Striatal dopamine and the temporal control of behavior. *Behav Brain Res.* 2019;356:375–379.
124. Frank MJ, Doll BB, Oas-Terpstra J, Moreno F. Prefrontal and striatal dopaminergic genes predict individual differences in exploration and exploitation. *Nat Neurosci.* 2009;12:1062.
125. Balcí F, Wiener M, Çavdaroglu B, Coslett HB. Epistasis effects of dopamine genes on interval timing and reward magnitude in humans. *Neuropsychologia.* 2013;51:293–308.
126. Bjorklund A, Lindvall O. Dopamine-containing systems in the CNS, 1984. [https://scholar.google.com/scholar?hl=en&as\\_sdt=0%2C22&q=Bjorklund+Lindvall.+Dopamine-containing+systems+in+the+CNS&btnG=](https://scholar.google.com/scholar?hl=en&as_sdt=0%2C22&q=Bjorklund+Lindvall.+Dopamine-containing+systems+in+the+CNS&btnG=)
127. Le Moal M, Simon H. Mesocorticolimbic dopaminergic network: functional and regulatory roles. *Physiological Rev.* 1991;71:155–234.

128. Niv Y, Daw ND, Joel D, Dayan P. Tonic dopamine: opportunity costs and the control of response vigor. *Psychopharmacology*. 2007;191:507–520.
129. R St Onge J, Ahn S, Phillips AG, Floresco SB. Dynamic fluctuations in dopamine efflux in the prefrontal cortex and nucleus accumbens during risk-based decision making. *J Neurosci*. 2012;32:16880–16891.

### ACKNOWLEDGEMENTS

The authors are grateful to Rahul Bhui for comments on an earlier draft of the paper.

### AUTHOR CONTRIBUTIONS

JGM and SJG developed the model and contributed to the writing of the paper. JGM analyzed and simulated the model, made the figures, and wrote the first draft.

### FUNDING INFORMATION

The project described was supported by National Institutes of Health grants T32GM007753 (JGM), T32MH020017 (JGM), and U19 NS113201-01 (SJG). The content is solely the responsibility of the authors and does not necessarily represent the official views of the National Institutes of Health. The funders had no role in study

design, data collection and analysis, decision to publish, or preparation of the manuscript.

### Compliance with ethical standards

### COMPETING INTERESTS

The authors declare no competing interests.

### ADDITIONAL INFORMATION

**Supplementary information** The online version contains supplementary material available at <https://doi.org/10.1038/s41386-021-01125-z>.

**Correspondence** and requests for materials should be addressed to J.G.M.

**Reprints and permission information** is available at <http://www.nature.com/reprints>

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.