

Deliberation gated by opportunity cost adapts to context with urgency

Maximilian Puelma Touzel,^{1,2,*} Paul Cisek,³ and Guillaume Lajoie^{1,4,5}

¹*Mila, Quebec AI Institute*

²*Department of Computer Science and Operations Research, Université de Montréal*

³*Department of Neuroscience, Université de Montréal*

⁴*Department of Mathematics and Statistics, Université de Montréal*

⁵*Canada CIFAR AI Chair*

Abstract

Finding the right amount of deliberation, between insufficient and excessive, is a hard decision making problem that depends on the value we place on our time. Average-reward, putatively encoded by tonic dopamine, serves in existing reinforcement learning theory as the stationary opportunity cost of time, and of deliberation in particular. However, this cost often varies with environmental context that can change over time. Here, we introduce an opportunity cost of deliberation estimated adaptively on multiple timescales to account for non-stationary contextual factors. We use it in a simple decision-making heuristic based on average-reward reinforcement learning (AR-RL) that we call *Performance-Gated Deliberation* (PGD). We propose PGD as a strategy used by animals wherein deliberation cost is implemented directly as urgency, a previously characterized neural signal effectively controlling the speed of the decision-making process. We show PGD outperforms AR-RL solutions in explaining behaviour and urgency of non-human primates in a context-varying random walk prediction task and is consistent with relative performance and urgency in a context-varying random dot motion task. We make readily testable predictions for both neural activity and behaviour and call for an integrated research program in cognitive and systems neuroscience around the value of time.

Keywords: primate decision-making, reinforcement learning, urgency, opportunity cost

* puelmatm@mila.quebec

symbol	quantity
t	within-trial time
k	trial index
S_t	within-trial state at time t
\mathbf{S}_t	state sequence up to time t
R_k	reward of k th trial
T_k	duration of k th trial
t_k^{dec}	decision time of k th trial
C_t^{del}	within-trial opportunity cost of deliberation
r_{max}	maximum reward achievable in a trial
b_t	belief of correct report given \mathbf{S}_t
\bar{r}_t	expected reward for reporting at time t
C_t^{com}	within-trial opportunity cost of commitment
ρ	stationary reward rate
ρ^*	optimal stationary reward rate
α	context parameter
ρ_α	context-conditioned stationary reward rate
T_α	context-conditioned stationary average trial duration
$\hat{\rho}_k^\tau$	reward history filtered through a timescale, τ
τ_{long}	a long timescale over which to estimate ρ
τ_{context}	a context-specific timescale over which to estimate ρ_α
ν	tracking cost sensitivity
K	subjective reward scale factor
T_{block}	characteristic duration of a trial block
c	auxiliary deliberation cost rate
N_t	tokens difference
p	jump probability of random walk, $p \geq 1/2$

Table I. Symbol glossary. Highlighted in gray are parameters of the PGD model presented in this paper.

10

INTRODUCTION

11 Humans and other animals make a wide range of decisions throughout their daily lives.
 12 Any particular action usually arises out of a hierarchy of decisions involving a careful balance
 13 between resources, including one that is always limited: time. The cost of *spending* time
 14 depends on its value, a construct that relies on comparing against the alternative things
 15 an agent could potentially do with it. Estimating time's value is not straightforward for a
 16 number of reasons. There are alternative choices at multiple decision levels, e.g. moving on
 17 from a job and moving on from a career, and each level requires its own evaluation. Moreover,
 18 the value of alternatives needs to be tracked as they may change over time depending on the
 19 context in which a decision is made. For example, animals will learn to value a given food
 20 resource differently depending on whether it is encountered during times of plenty versus
 21 scarcity. The agent's knowledge of and ability to track context thus influences the value it
 22 assigns to possible alternatives.

23 These are significant, practical complications of making decisions contingent on *opportu-*

24 *nity costs* [1], the formal economic concept capturing the value of the alternatives lost by
25 committing a limited resource to a given use. The opportunity cost of time is nevertheless
26 well-studied in decision-making theory. It plays the role of a reference reward in defini-
27 tions of relative value, most notably as the average reward in average-reward reinforcement
28 learning (AR-RL) [2].

29 In neuroscience, AR-RL was first proposed to extend the reward prediction error hy-
30 pothesis for phasic dopamine to account also for the observed properties of tonic dopamine
31 levels [3]. It has since been used to emphasize the relative nature of reward-based decision-
32 making [4] in explanations of human and animal behaviour in foraging [5], free-operant
33 conditioning [6], perceptual decision-making [7, 8], cognitive effort/control [8, 9], and even
34 economic exchange [10].

35 Unlike the alternative discount-reward approach, AR-RL is a theoretically well-defined
36 and numerically stable formulation for long horizon decision problems [11], such as those
37 in *continuing environments* in which there is no definite end [12]. Solutions to AR-RL
38 problems maximize average reward, in contrast to traditional fixed accuracy criteria in
39 perceptual decision-making tasks that focus on maximizing trial reward alone [13]. The
40 solutions to AR-RL formulations of tasks of long sequence of trials are decision boundaries
41 in the state space of a trial. Determining this decision boundary requires maximizing the
42 relative value, defined using the opportunity cost of time. The resulting optimal decision
43 boundaries typically ‘collapse’ over a trial: they cut deliberation short, e.g. in tasks where
44 trial difficulty is variable [7, 14]. Up to now, however, AR-RL and most of its applications
45 have focused on fixed context and have used the stationary average reward as the fixed
46 opportunity cost of time, which ignores context-dependent performance variation. This is
47 perhaps not surprising given that in psychological and neuroscientific studies of decision-
48 making, we usually eliminate such contextual factors from the experimental design such
49 that our models describe stationary behaviour. However, the brain mechanisms under study
50 are adapted to a more diverse natural world in which changing environmental factors are
51 often relevant, hard to infer and vary over time [4].

52 We pursue a theory of approximate relative-value decision-making under uncertainty in a
53 setting relevant to decision-making neuroscience. We start by showing that value in AR-RL
54 can be expressed using the opportunity costs of deliberation and commitment. Here, the
55 commitment cost is the shortfall in reward (relative to the maximum possible in a trial)
56 that is expected to be lost when committing to a decision at a given time. Highlighting the
57 risk of value representations in non-stationary environments, we propose an approximation
58 to the AR-RL value-optimal solution, Performance-Gated Deliberation (PGD), that uses
59 the increasing opportunity cost of time in a trial to collapse the decision boundary directly,
60 by-passing the need to maximize relative value. PGD thus reduces decision-making to
61 estimating two opportunity costs: a commitment cost learned from the statistics of the
62 environment and a deliberation cost estimated from tracking one’s own performance in that
63 environment. It explains how an agent, without explicitly tracking context parameters or
64 storing a value function, can trade-off speed and accuracy according to performance at
65 the typically longer timescales over which context changes. We propose that deliberation
66 cost is then directly encoded as “urgency” in the neural dynamics underlying decision-
67 making [7, 15–17]. The theory is thus directly testable using both behaviour and neural
68 recordings.

69 To illustrate how PGD applies in a specific continuing decision-making task, and to make
70 the links to a neural implementation explicit, we analyze behavior and neural recordings

71 collected over eight years from two non-human primates (NHPs) [18, 19]. They performed
72 successive trials of the “tokens task”, a probabilistic guessing task in which information
73 about the correct choice is continuously changing within each trial, and a task parame-
74 ter controlling the incentive to decide early (the context) is varied over longer timescales.
75 Behavior in the task, in both humans [16] and monkeys [19], provides additional support
76 to an existing hypothesis about how neural dynamics implements time-sensitive decision-
77 making [15]. Specifically, neural recordings in monkeys suggest that the evidence needed
78 to make the decision predominates in dorsolateral prefrontal cortex [20]; a growing context-
79 dependent urgency signal is provided by the basal ganglia [21]; and the two are combined to
80 bias and time, respectively, a competition between potential actions that unfolds in dorsal
81 premotor and primary motor cortex [18]. Similar findings have been reported in other tasks -
82 for example, in the frontal eye fields during decisions about eye-movements [17]. We propose
83 PGD as a theoretical explanation for why decision-making mechanisms are organized in this
84 way. As an algorithm, it serves as a robust means to balance immediate rewards and the cost
85 of time across multiple timescales. As a quantitative model, it serves to explain concurrently
86 recorded behaviour and neural urgency in continuing decision-making tasks. From neural
87 recordings in non-human primates and and behaviour in human and non-human primates,
88 we show that it does so more accurately than AR-RL solutions. Adapting PGD to the
89 random dot motion task in which urgency was first characterized [17], we make quantitative
90 predictions about neural urgency is such tasks, which we validate on their data within error
91 bounds.

92 RESULTS

93 A. Theory of performance-gated deliberation

94 1. Opportunity costs of deliberation and commitment, and drawbacks of average-reward 95 reinforcement learning

96 We consider a class of tasks consisting of a long sequence of trials indexed by $k =$
97 $1, 2, \dots$ (see fig. 1a), each of which provides the opportunity to obtain some reward by choos-
98 ing correctly. In each trial, a finite sequence of states, S_t , $t = 0, \dots, t_{\max}$, is observed that
99 provide evidence for an evolving belief about the correct choice among a fixed set of options.
100 To keep notation simple, we suppress denoting the trial index, k , on quantities such as trial
101 state, S_t , that also depend on trial time, t . The time of decision, t_k^{dec} , and the chosen option
102 determine both the reward received, R_k , and the trial duration, $T_k \geq t_k^{\text{dec}}$. Importantly,
103 decision timing can affect performance because earlier decisions typically lead to shorter
104 trials (and thus more trials in a given time window), while later decisions lead to higher
105 accuracy. Effectively balancing such speed-accuracy trade-offs is central to performing well
106 in continuing episodic task settings. For a fixed strategy, the *stationary reward rate* (see
107 slope of dashed line in fig. 1a(right)) is
108

$$\rho := \lim_{k \rightarrow \infty} \frac{\sum_k R_k}{\sum_k T_k} . \quad (1)$$

109 For a stochastic environment, the definition of ρ includes an ensemble average. Free-operant
110 conditioning, foraging, and several perceptual decision-making tasks often fall into this class.

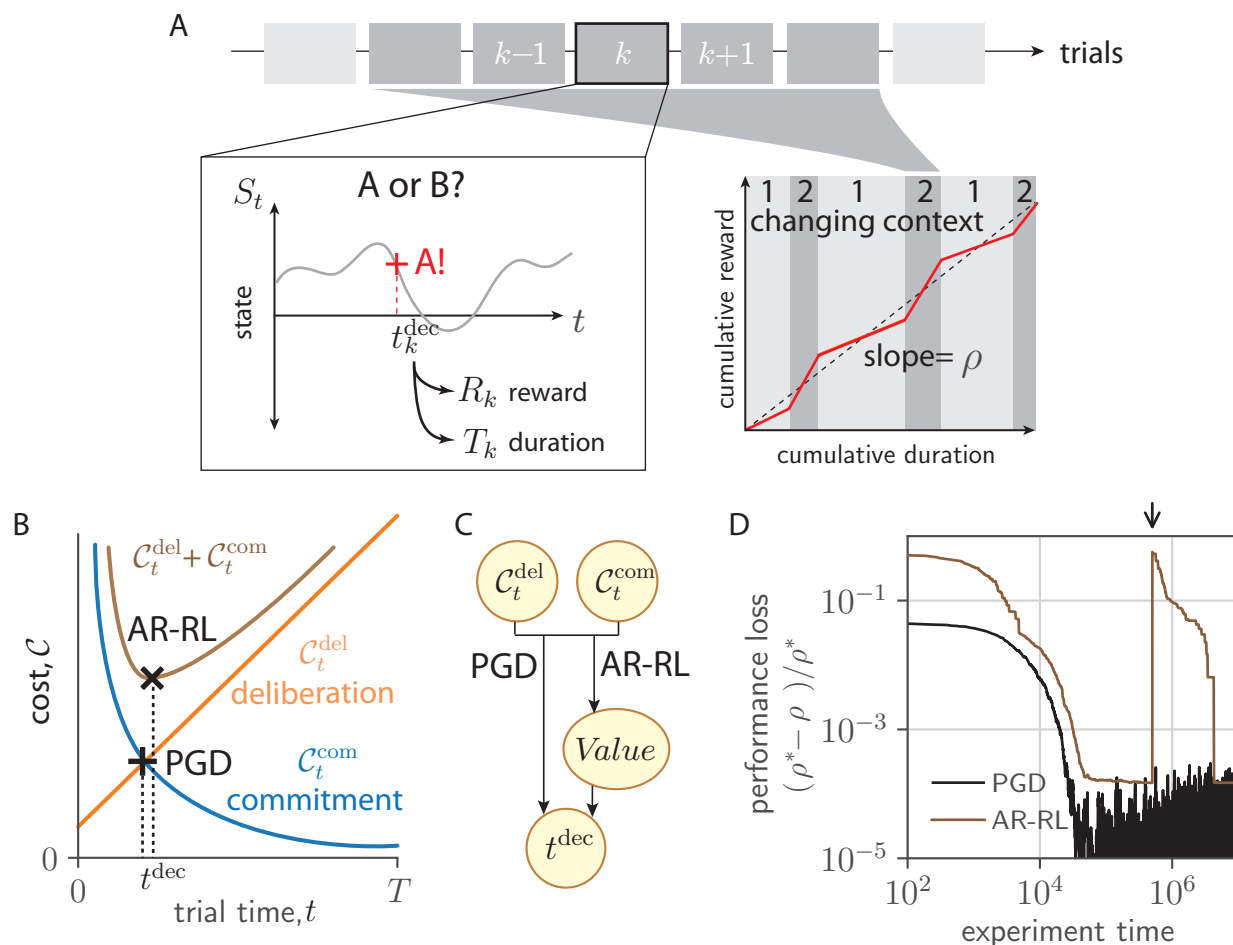


Figure 1. *AR-RL and Performance-Gated Deliberation*. (a) Task setting. Left: Within trial state, S_t evolves over trial time t in successive trials indexed by k . The decision ‘A’ is reported at the decision time t_k^{dec} (red cross), determining trial reward, R_k , and trial duration, T_k . Right: Sketch of cumulative reward versus cumulative duration. Context-conditioned reward rate (slope of red line), varies with alternating context (labelled 1 and 2) around average reward, ρ (dashed line). (b) Decision rules based on opportunity costs of commitment, C_t^{com} , and deliberation, C_t^{del} . The AR-RL rule (black ‘x’) finds t that minimizes $C_t^{\text{del}} + C_t^{\text{com}}$. The PGD rule (black cross) finds t^{dec} at which they intersect, $C_t^{\text{del}} = C_t^{\text{com}}$. (c) Schematic diagram of each algorithm’s dependency. PGD computes a decision time directly from the two opportunity costs, while AR-RL uses both to first estimate a value function, whose maximum specifies the decision time. (d) Loss (error in performance with respect to the optimal policy, $(\rho^* - \rho)/\rho^*$) over learning time in a patch-leaving task (AR-RL: brown, PGD: black). The arrow indicates when the state labels were randomly permuted.

111 Previous work [7, 22] has studied the belief of correct report for binary rewards, $b_t = P(R_k =$
 112 $1 | \mathbf{S}_t, t^{\text{dec}} = t)$, which also gives the expected trial reward, $\bar{r}_t = b_t \cdot 1 + (1 - b_t) \cdot 0 = b_t$ [7]
 113 (see [23] for more about the relationship between value-based and perceptual decisions). \mathbf{S}_t
 114 denotes the state sequence observed so far, (S_0, \dots, S_t) . We consider greedy strategies that
 115 report the choice with the largest belief at decision time. The decision problem is then about
 116 *when* to decide.

117 Average-reward reinforcement learning (AR-RL), first proposed in artificial intelli-

118 gence [24], was later incorporated into reward prediction error theories of dopamine sig-
119 nalling [3] and employed to account for the opportunity cost of time [6]. AR-RL was
120 subsequently used to study reward-based decision-making in neuroscience and psychol-
121 ogy [7, 8, 25, 26]. AR-RL centers around the average-adjusted future return, which penalizes
122 the passage of time according to the average reward. A reporting decision is associated with a
123 return that for trial-based tasks combines the remainder of the current trial and all future
124 trials, $\bar{r}_t - \rho(T_k - t) + \sum_{k' > k} (R_{k'} - \rho T_{k'})$, where ρ (*c.f.* eq. (1)) is either estimated online
125 or obtained self-consistently (see [Methods](#) for details). Value is defined as the future return
126 averaged over trial sequence realizations. This average of a sum of reward deviations into
127 the future converges on account of the decaying effects of the state at which the decision is
128 made. The AR-RL algorithms we consider aim to achieve the highest ρ by also maximizing
129 the average-adjusted value. We now provide an alternative, but equivalent definition of
130 average-adjusted trial return in terms of opportunity costs incurred by the agent.

131 We denote the opportunity cost of committing at time t within a trial as $\mathcal{C}_t^{\text{com}}$, defined
132 as the difference

$$\mathcal{C}_t^{\text{com}} = r_{\max} - \bar{r}_t, \quad (2)$$

133 where r_{\max} is the maximum trial reward possible *a priori*. Within a trial, an agent lowers
134 its commitment cost towards zero by accumulating more evidence, i.e. by waiting. Waiting,
135 however, incurs another opportunity cost: the reward lost by not acting. We denote this
136 opportunity cost of deliberation incurred up to a time t in a trial as $\mathcal{C}_t^{\text{del}}$. In AR-RL, the
137 constant opportunity cost rate of time is integrated so that for $T_k = t_k^{\text{dec}}$,

$$\mathcal{C}_t^{\text{del}} = \rho t. \quad (3)$$

138 With these definitions, the average-adjusted trial return for deciding at a time t can be
139 expressed as $r_{\max} - (\mathcal{C}_t^{\text{com}} + \mathcal{C}_t^{\text{del}})$. It is maximized by jointly minimizing $\mathcal{C}_t^{\text{del}}$ and $\mathcal{C}_t^{\text{com}}$ ([fig. 1b](#)),
140 giving the AR-RL optimal solution (see [Methods](#) for a formal statement and solution of the
141 AR-RL problem). Expressed in this way, the average-adjusted trial return emphasizes the
142 more general perspective that an agent's solution to the speed-accuracy trade-off is about
143 how it balances the decaying opportunity cost of commitment and the growing opportunity
144 cost of deliberation.

145 Despite their utility, value representations such as the average-adjusted trial return can
146 be a liability in real world tasks where task statistics are non-stationary. To illustrate this,
147 we consider the following foraging task. An foraging agent feeds among a fixed set of food
148 (e.g. berry) patches. Total berries consumed in a patch saturates with duration t according
149 to a given saturation profile, shared across patches, as the fewer berries left are harder to
150 find. Patches differ in their richness (e.g. berry density), which is randomly sampled and
151 fixed over the task. Denoting patch identity (serving as context) by s , the food return is
152 directly observed and deterministic given s . To perform well, the agent needs to decide when
153 to move on from depleting the current patch. Further details about the task and its solution
154 are given in the [Methods](#). For a broad class of online AR-RL algorithms, the agent learns the
155 average-adjusted trial return as a function of state and time. For a given patch, it then leaves
156 when this return is at its maximum (*c.f.* [fig. 1b](#)). In [fig. 1d](#), we show how the performance
157 (brown line) approaches that of the optimal policy in time as the estimation of the AR-RL
158 trial return improves with experience (see [Methods](#) for implementation details). However, if
159 the agent's environment undergoes a significant disturbance (e.g. a forest fire due to which
160 the patch locations are effectively re-sampled), the performance of this AR-RL algorithm can
161 drop back to where it started. We implement such a disturbance via random permutation

162 of the state labels at the time indicated by the arrow in [fig. 1d](#). This is true over a range of
163 learning rates and the number of patches ([fig. S8](#)). More generally, any approach that relies
164 on estimating state-value associations shares this drawback, including those approaches that
165 implicitly learn those associations by directly learning a policy instead [27]. Could context-
166 dependent decision times be obtained without having to associate value or action to state?
167 A means to do so is presented in the next section.

168 2. Performance-Gated Deliberation

169 We propose that instead of maximizing value as in AR-RL, which minimizes the sum of
170 the two opportunity costs, $\mathcal{C}_t^{\text{del}} + \mathcal{C}_t^{\text{com}}$, the agent simply takes as its decision criterion when
171 they intersect (shown as the black cross in [fig. 1b](#)).

$$t^{\text{dec}} := \min_t \{t \mid \mathcal{C}_t^{\text{del}} \geq \mathcal{C}_t^{\text{com}}\} \quad (\text{PGD decision rule}) \quad (4)$$

172 We call this heuristic rule at the center of our results *Performance-Gated Deliberation*
173 (PGD). Plotted alongside the AR-RL performance in [fig. 1d](#) for our example foraging task,
174 PGD (black line) achieves better performance than AR-RL overall. It is also insensitive to
175 the applied disturbance since PGD uses $\mathcal{C}_t^{\text{del}}$ and $\mathcal{C}_t^{\text{com}}$ directly when deciding, rather than
176 as input to problem of optimizing average-adjusted value as in AR-RL ([fig. 1c](#)).

177 We constructed the above task so that PGD is the AR-RL optimal solution. In general,
178 however, PGD is a well-motivated approximation to the optimal strategy, so we call it a
179 heuristic. In the more general stochastic setting where there is residual uncertainty in trial
180 reward at decision time, the PGD agent will have to learn the association between state
181 and expected reward, \bar{r}_t . This association is learned from within-trial correlations only. In
182 contrast, the opportunity cost of time as the basis for the deliberation cost depends on
183 across-trial correlations that together determine the overall performance. It is thus more
184 susceptible to non-stationarity. A typical task setting is when the value of the same low-level
185 action plan differs across context. From hereon, we will assume the agent has learned the
186 stationary opportunity cost of commitment and so focus on resolving the remaining problem:
187 how to learn and use an opportunity cost of deliberation that exhibits non-stationarity on
188 the longer timescales over which context varies.

189 3. Reward filtering for a dynamic opportunity cost of deliberation

190 The state disturbance in the toy example above altered task statistics at only a single
191 time point. In general, however, changes in task statistics over time can occur throughout
192 the task experience. A broader notion of deliberation cost beyond the static average reward
193 is thus needed—one that can account for extended timescales over which performance varies.
194 Such a cost serves as a dynamic reference in a relative definition of value based on a non-
195 stationary opportunity cost of time. We first address how performance on various timescales
196 can be estimated.

197 As a concrete example, we make use of the task that we will present in detail in the fol-
198 lowing section. This task has a context parameter, α , that can vary in time on characteristic
199 timescales longer than the moment-to-moment and can serve as a source of non-stationarity
200 in performance. Here, the context sequence, α_k , varies on a single timescale, e.g. through pe-
201 riodic switching between two values. The resulting performance ([fig. 2a\(top\)](#)) varies around

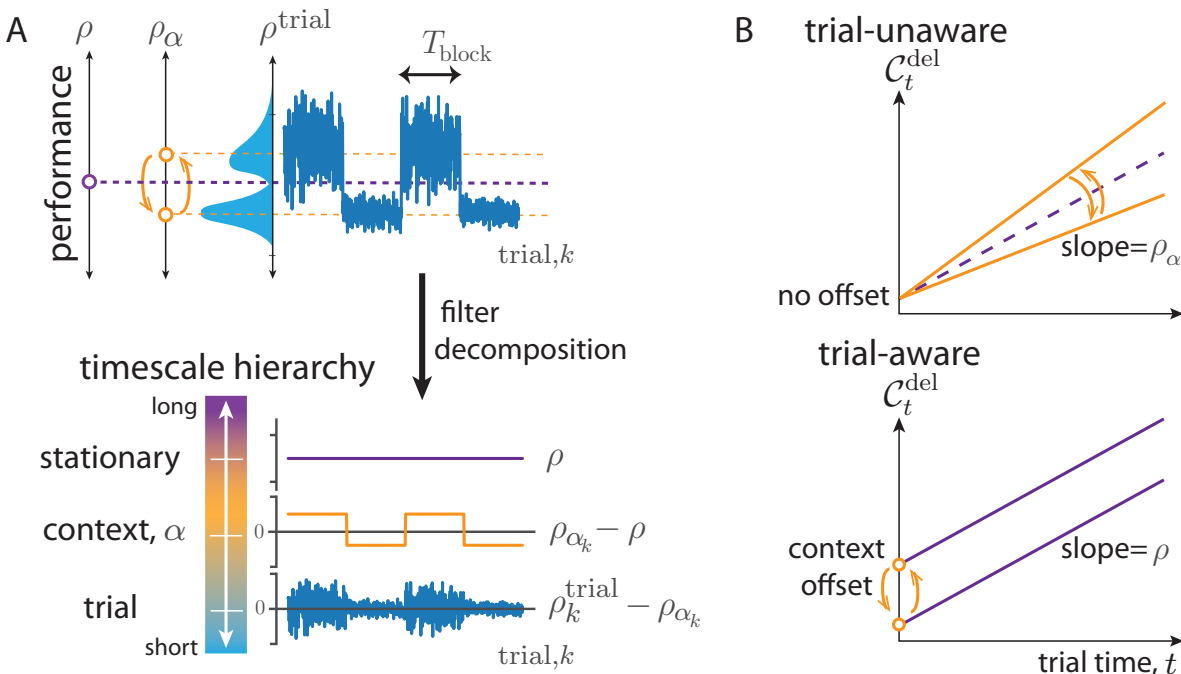


Figure 2. *Non-stationary opportunity cost.* (a) Top: Dynamics of trial performance ($\rho_k^{\text{trial}} := R_k/T_k$; blue) with its distribution as well as dynamics of between context-conditioned averages of performance ($\rho_\alpha = \langle \rho_k^{\text{trial}} \rangle_{k|\alpha}$; orange), and the effectively stationary average performance ($\rho \sim \langle \rho_k^{\text{trial}} \rangle_k$; purple). Bottom: these are decomposed into a hierarchy by filtering reward history on trial, context, and long timescales, respectively. (b) Two hypothetical forms for context-specific trial opportunity cost. Top: Trial-unaware cost in which context varies the slope around ρ . Bottom: Trial-aware cost in which context variation is through a bias (eq. (5)).

202 the stationary average, ρ (purple), with context variation due to the switching (orange), as
 203 well as context-conditioned trial-to-trial variation (blue). The decomposition of time-varying
 204 performance into these multiple, timescale-specific components can be achieved by passing
 205 the reward signal through parallel filters, each designed to retain the signal variation specific
 206 to that timescale (fig. 2a(bottom)). There are multiple approaches to this decomposition.
 207 We chose a heuristic approach in which the performance over a finite memory timescale can
 208 be estimated by filtering the sequence of rewards through a simple low-pass filter [8, 28].
 209 This filter is defined by an integration time, τ , tuned to trade off the bias and variance
 210 of the estimate in order to best capture the variation on the desired timescale (e.g. how
 211 performance varies over different contexts). We denote such an estimate $\hat{\rho}_k^\tau$, and show in
 212 the Methods that it approximates the average reward over the last τ time units. We discuss
 213 the question of biological implementation in the discussion, but note here that the number
 214 and values of τ needed to represent performance variation in a given task could be learned
 215 or selected from a more complete set in an online fashion during task learning. In an exper-
 216 imental setting, these learned values can in principle be inferred from observed behaviour
 217 and we developed such an approach in the analysis of data that we present in the following
 218 section.

219 Applying this heuristic decomposition here, the stationary reward rate, ρ , can be esti-
 220 mated to high precision by using a long integration time, τ_{long} , to the reward sequence R_k ,

221 producing the estimate $\hat{\rho}_k^{\tau_{\text{long}}}$. If α_k were a constant sequence, $C_t^{\text{del}} = \hat{\rho}_k^{\tau_{\text{long}}}t$, the station-
222 ary opportunity cost of deliberation eq. (3) of AR-RL. However, in this example context
223 varies on a specific timescale, to which the former is insensitive. Thus, a second filtered
224 estimate $\hat{\rho}_k^{\tau_{\text{context}}}$ is needed to estimate performance on this timescale. Unlike $\hat{\rho}_k^{\tau_{\text{long}}}$, this es-
225 timate tracks the effective instantaneous, context-specific performance, ρ_{α_k} . Its estimation
226 error arises from a trade-off, controlled by the integration time, τ_{context} , between its speed
227 of adaptation and its finite memory.

228 We consider two distinct hypotheses for how to extend AR-RL to settings where perfor-
229 mance varies over context. The first hypothesis, $C_t^{\text{del}} = \rho_{\alpha}t$, is the straightforward, *trial-*
230 *unaware* extension of eq. (3), shown in fig. 2b(top). Here, performance is tracked only
231 on a timescale sufficient to capture context variation and the corresponding cost estimate,
232 $\hat{\rho}_{k-1}^{\tau_{\text{context}}}$, is incurred moment-to-moment, neglecting the trial-based task structure. However,
233 this incorrectly lumps together two distinct opportunity costs: those incurred by moment-
234 by-moment decisions and those incurred as a result of the effective planning implied by
235 performance that varies over context. In particular, context is defined over trials not mo-
236 ments, and thus the context-specific component of opportunity cost of a trial is a sunken
237 cost paid at the outset of a trial. This inspires a second *trial-aware* hypothesis

$$C_t^{\text{del}} = \rho t + (\rho_{\alpha} - \rho)T_{\alpha} . \quad (\text{trial-aware opportunity cost}) \quad (5)$$

238 Equation (5) is plotted over trial time t in fig. 2b(bottom). Its first term is the AR-RL
239 contribution from the stationary opportunity cost of moment-to-moment decisions using
240 the stationary reward rate, ρ estimated with $\hat{\rho}_k^{\tau_{\text{long}}}$. The second, novel term in eq. (5) is a
241 context-specific trial cost deviation incurred at the beginning of each trial and computed as
242 the average deviation in opportunity cost accumulated over a trial from that context (T_{α}
243 is the average duration of a trial in context α). This deviation fills the cost gap made by
244 using the stationary reward rate ρ in the moment-to-moment opportunity cost instead of
245 the context-specific average reward, ρ_{α} . This baseline cost derived from the orange time
246 series in fig. 2a(bottom) vanishes in expectation, as verified through the mixed-context
247 ensemble average reward (e.g. $\rho \equiv \sum_{\alpha} \rho_{\alpha}T_{\alpha} / \sum_{\alpha} T_{\alpha}$ when the context is distributed evenly
248 among trials such that $\sum_{\alpha} (\rho_{\alpha} - \rho)T_{\alpha} = 0$). Thus, this opportunity cost reduces to that
249 used in AR-RL when ignoring context, and suggests a generalization of average-adjusted
250 value functions to account for non-stationary context. We estimate this baseline cost using
251 $(\hat{\rho}_{k-1}^{\tau_{\text{context}}} - \hat{\rho}_{k-1}^{\tau_{\text{long}}})T_{k-1}$, where we have used the sample T_{k-1} in lieu of the average T_{α} . See fig. S1
252 for a signal filtering diagram that produces this estimate of eq. (5) from reward history. A
253 main difference between the cost profiles from the two hypotheses is the cost at early times.
254 Both the behaviour and neural recordings we analyze below seem to favor the second, trial-
255 aware hypothesis eq. (5). We hereon employ that version in the main text, and show the
256 results for the trial-unaware hypothesis in fig. S7.

257

B. Neuroscience application: PGD in the tokens task

258 In this section, we apply the PGD algorithm to the “tokens task” [16]. We first give a
259 simulated example with periodic context dynamics. We then present an application to a
260 set of non-human primate experiments in which context variation was non-stationary [19].
261 For the latter, we used the decision time dynamics over trials to fit a model for each of the
262 two subjects. We then validated the models by assessing their ability to explain (1) the

263 concurrently recorded behaviour via their context-specific behavioural strategies and (2) the
264 neural activity in premotor cortex (PMd) via the temporal profile of the underlying neural
265 urgency signals.

266 In the tokens task, the subject must guess as to which of two peripheral reaching targets
267 will receive the majority of tokens that randomly jump, one by one every 200ms, from a
268 central pool initialized with a fixed number of tokens. Importantly, after the subject reports,
269 the interval between remaining jumps contracts to once every 150ms (the “slow” condition)
270 or once every 50ms (the “fast” condition), giving the subject the possibility to save time by
271 taking an early guess. The interval contraction factor, $1 - \alpha$, for slow ($\alpha = 1/4$) and fast
272 ($\alpha = 3/4$) condition is parametrized $\alpha \in [0, 1]$, the incentive strength to decide early, which
273 then serves as the task context.

274 In contrast to the patch leaving task example from Section A, the tokens task has many
275 within-trial states and the state dynamics is stochastic. With the t^{th} jump labelled $S_t \in$
276 $\{-1, 1\}$ serving as the state, for the purposes of prediction, the history of states can be
277 compressed into the tokens difference, $N_t = \sum_{i=1}^t S_i$, between the two peripheral targets
278 with $N_0 = 0$. The dynamics of N_t is an unbiased random walk (see [fig. 3a](#)), with its current
279 value sufficient to determine the belief of a correct report, b_t (computed in [Methods](#)). Since
280 for binary rewards, b_t is also the expected reward, N_t is also sufficient for determining the
281 opportunity cost of commitment, $\mathcal{C}_t^{\text{com}}$ ([eq. \(2\)](#)). We display this commitment cost dynamics
282 in [fig. 3b](#). It evolves on a lattice (gray), always starting at 0.5 (for $p = 1/2$) and ending at 0
283 for all p . We assume the agent has learned to track this commitment cost. The PGD agent
284 uses this commitment cost, along with the estimate of the trial-aware deliberation cost, to
285 determine when to stop deliberating and report its guess.

287 1. A simulated example for a regularly alternating context sequence

288 We first show the behaviour of the PGD algorithm in the simple case where α switches
289 back and forth every 300 trials (see [fig. 3](#)). We call such segments of constant α ‘trial blocks’,
290 with context alternating between slow ($\alpha = 1/4$) and fast ($\alpha = 3/4$) blocks. The decision
291 space in PGD is a space of opportunity costs, equivalent to the alternative decision space
292 formulated using beliefs [7]. In particular, one can think of the deliberation cost as the
293 decision boundary ([fig. 3b](#)). This boundary is dynamic (see [Supplemental video](#)), depending
294 on performance history via the estimates, $\hat{\rho}_k^{\tau_{\text{context}}}$ and $\hat{\rho}_k^{\tau_{\text{long}}}$, of the context-conditioned and
295 stationary average reward, respectively. The result of these dynamics is effective context
296 planning: the PGD algorithm sacrifices accuracy to achieve shorter trial duration in trials
297 of the fast block, achieving a higher context-conditioned reward rate compared to decisions
298 in the slow block (*c.f.* the slopes shown in the inset of [fig. S2d](#)). This behaviour can be
299 understood by analyzing the dynamics of $\hat{\rho}_k^{\tau_{\text{context}}}$ and $\hat{\rho}_k^{\tau_{\text{long}}}$, and their effect on the dynamics
300 of the decision time ensemble.

301 The two performance estimates behave differently from one another solely because of
302 their distinct integration times. Ideally, an agent would choose τ_{context} to be large enough
303 that it serves to average over trial-to-trial fluctuations in a context, but short enough to
304 not average over context fluctuations. In contrast, the value of τ_{long} would be chosen large
305 enough to average over context fluctuations. We apply those choices in this simulated
306 example, with rounded values chosen squarely in the range in which the values inferred
307 from the behaviour in the following application will lie. As a result of this chosen values,
308 the context estimate $\hat{\rho}_k^{\tau_{\text{context}}}$ relaxes relatively quickly after context switches to the context-

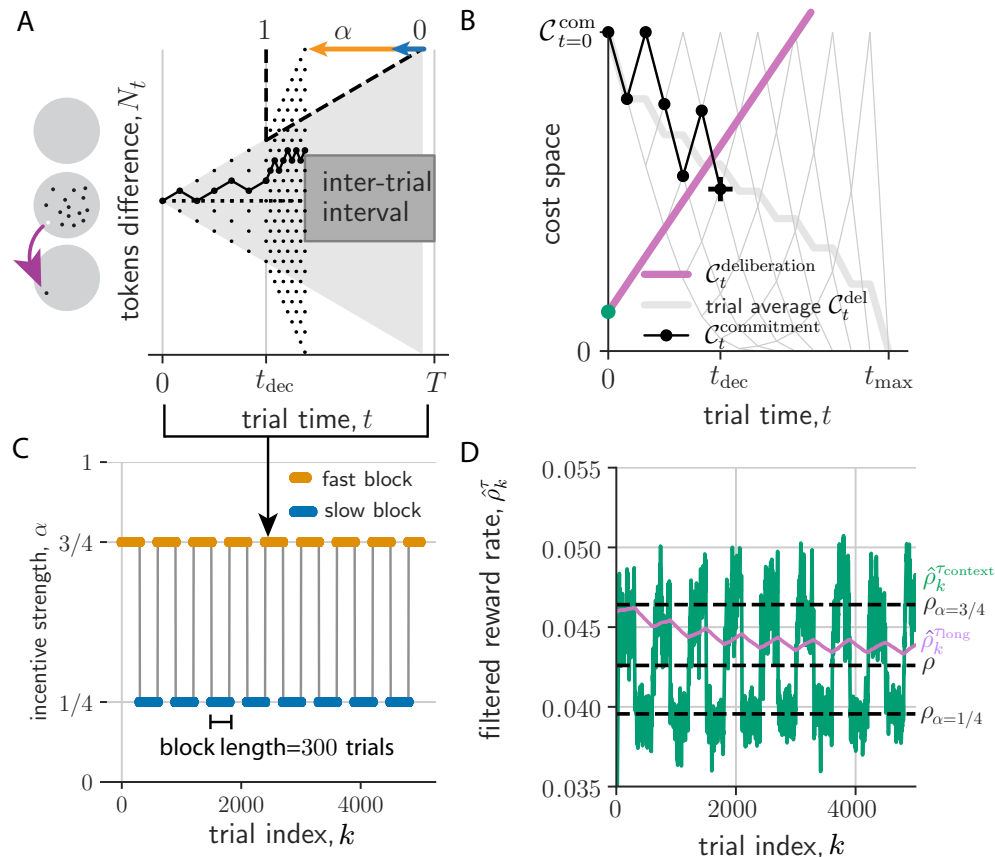


Figure 3. *PGD agent performs the tokens task for periodic context switching.* (a) A tokens task trial. Left: Tokens jump from a center to a peripheral region (gray circles). Right: The tokens difference, N_t , evolves as a random walk that accelerates according to α (here $3/4$) post-decision time, t^{dec} . The trial duration is T , which includes an inter-trial interval. (b) Decision dynamics in cost space obtained from evidence dynamics in (a). Commitment cost trajectories (gray lattice; thick gray: trial-averaged) start at $C_{t=0}^{\text{com}}$ and end at 0. Trajectory from (a) shown in black. t^{dec} (black cross) is determined by the crossing of the commitment and deliberation cost. (c) Incentive strength switches between two values every 300 trials. (d) Expected rewards filtered on τ_{long} ($\hat{\rho}_k^{\tau_{\text{long}}}$, purple) and τ_{context} ($\hat{\rho}_k^{\tau_{\text{context}}}$, green). Black dashed lines from bottom to top are $\rho_{\alpha=1/4}$, ρ , and $\rho_{\alpha=3/4}$.

309 conditioned stationary average performance (dashed lines in [fig. 3d](#)), but exhibits stronger
 310 fluctuations as a result. The estimate of the stationary reward, $\hat{\rho}_k^{\tau_{\text{long}}}$, on the other hand has
 311 relatively smaller variance. This variance results from the residual zigzag relaxation over
 312 the period of the limit cycle. Given the characteristic block duration, T_{block} , we can be more
 313 precise. In particular, when T_{block} is much less than τ_{long} ($T_{\text{block}}/\tau_{\text{long}} \ll 1$), the within-block
 314 exponential relaxation is roughly linear. Thus, the average unsigned deviation between $\hat{\rho}_k^{\tau_{\text{long}}}$
 315 and the actual stationary reward, ρ , can be approximated using $1 - \exp[-T_{\text{block}}/\tau_{\text{long}}] \approx$
 316 $T_{\text{block}}/\tau_{\text{long}} \ll 1$. This scaling fits the simulated data well ([fig. S2d](#): inset).

317 The dynamics of these two performance estimates drives the dynamics of the k -conditioned
 318 decision time ensemble via how they together determine the deliberation cost ([eq. \(5\)](#); [Sup-](#)
 319 [plemental video](#)). For example, the mean component of this ensemble relaxes after a context

320 switch to the context-conditioned average, while the fluctuating component remains strong
321 due to the sequence of random walk realizations (fig. S2c). In the case of periodic context,
322 the performance estimates and thus also the decision time ensemble relax into a noisy peri-
323 odic trajectory over the period of a pair of fast and slow blocks (fig. 3d). Over this period,
324 they exhibit some stationary bias and variance relative to their corresponding stationary
325 averages (distributions shown in fig. S2e).

326 2. *Fit to behavioural data from non-human primates and model validation*

327 Next, we fit a PGD agent to each of the two non-human primates' behaviour in the
328 tokens task experiments reported in [19] and compare to AR-RL solutions. As with the
329 above example (*c.f.* fig. 3), trials were structured in alternating blocks of $\alpha = 1/4$ and
330 $\alpha = 3/4$. Figure 4a shows context-switching α -sequence from these experiments, which, in
332 contrast to the above example exhibits large, irregular fluctuations in block size [29].

333 So far, PGD has only two free parameters: the two filtering time constants, τ_{long} and
334 τ_{context} . We anticipated only a weak dependence of the fit on the τ_{long} , so long as it exceeded
335 the average duration of a handful of trial blocks enabling a sufficiently precise estimate of
336 ρ . In contrast, the context filtering timescale, τ_{context} , is a crucial parameter as it dictates
337 where the PGD agent lies on a bias-variance trade-off in estimating ρ_{α_k} , the value of which
338 determines the context-specific contribution to the deliberation cost (eq. (2)). To facilitate
339 the model's ability to fit individual differences, we introduce a subjective reward bias factor,
340 K , that scales the rewards fed into the performance filters. We also add a tracking-cost sen-
341 sitivity parameter, ν , that controls τ_{context} to avoid wasting adaptation speed (see Methods
342 for details). The latter made it possible to fit the asymmetric switching behaviour observed
343 in the average decision time dynamics. With these four parameters, we quantitatively match
344 the baselines and exponential-like relaxation of the average decision time dynamics around
345 the two context switches (fig. 4b,c; see Methods for fitting details).

346 A comparison of the best-fitting parameter values over the two monkeys (fig. 4d-f) sug-
347 gests that the larger the reward bias, K (fig. 4e), the more hasty the context-conditioned
348 performance estimate (the smaller τ_{context}), and the lower the sensitivity to the tracking cost
349 (fig. 4f). This is consistent with the hypothesis that subjects withhold cognitive effort in
350 contexts of higher perceived reward [8]. Along with the correspondence in temporal statistics
351 of the behaviour (e.g. fig. S6), the fitted model parameters for the two subjects provides a
352 basis on which to interpret the subject differences in the results of the next section, in par-
353 ticular their separation on a speed-accuracy trade-off, as originating in the distinct reward
354 sensitivity shown here.

355 To better understand where both the data and the learned PGD agent lie in the space
356 of strategies for the tokens task, we computed reward-rate (AR-RL) optimal solutions for
357 a given fixed context, α (here $\alpha \in [0, 1]$), using the same approach as [7] (conventional
358 discount-reward value iteration achieved the same solution in the limit of the undiscounted
359 case; result not shown). In each of average-reward and discount-reward formulations, the
360 dynamic programming approach involves iterating Bellman's equation to obtain the optimal
361 value functions from which the optimal policy and its reward rate can be obtained (see
362 Methods for details). The optimal reward rate as a function α is shown in fig. 5a. The
363 strategies generating these reward rates interpolate from the wait-for-certainty strategy at
364 low α to the one-and-done strategy [30] at high α . The α -conditioned reward rates achieved
365 by the two primates with their corresponding PGD model, and a reference human [31]

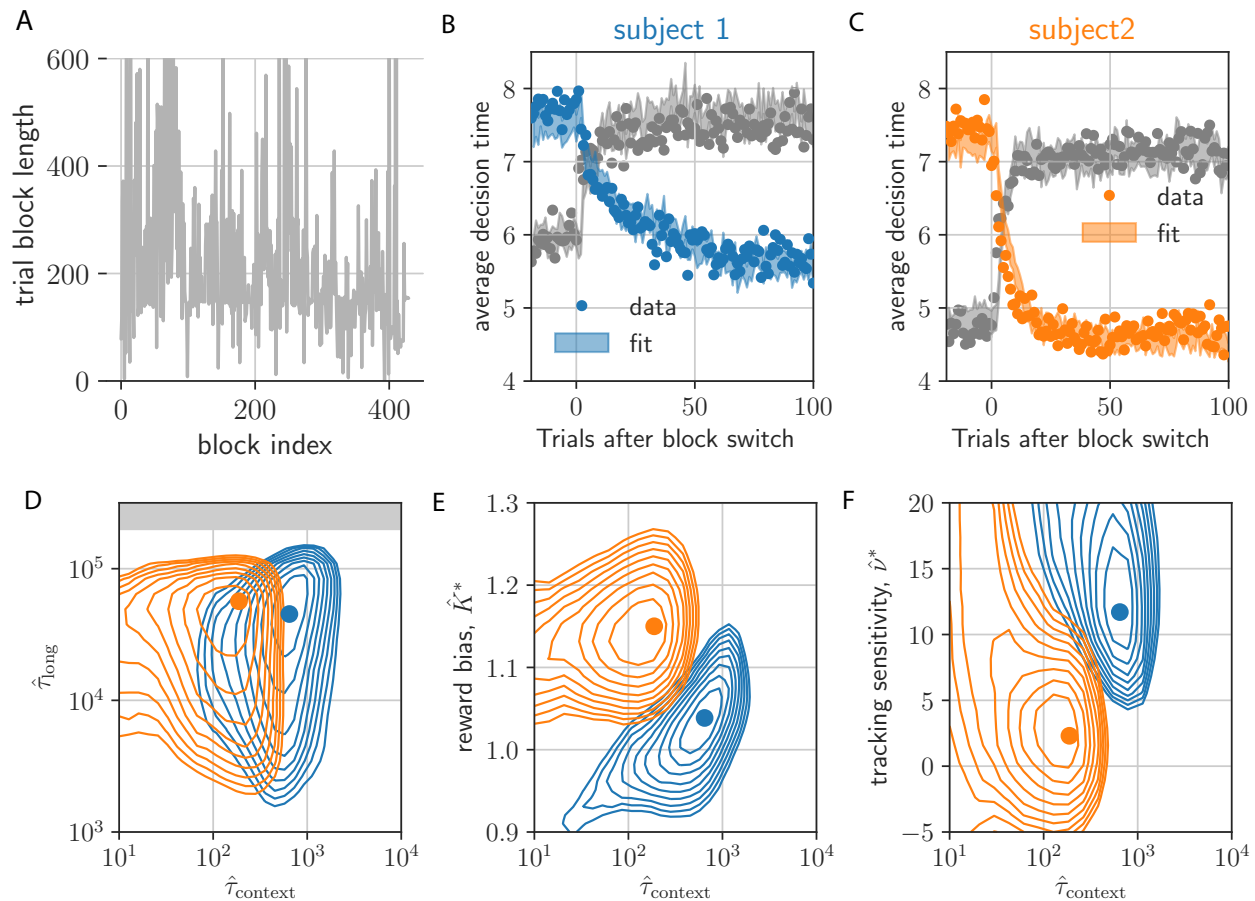


Figure 4. *PGD model fit to NHP behaviour for non-stationary α -dynamics reported in Ref. [19].* (a) Block length sequence used in the experiment. (b,c) decision times (dots) aligned on the context-switching event type (fast-to-slow in gray; slow-to-fast in color) and averaged. Shaded regions are the standard error bounds of the models' average decision times. (d) Error evaluated on a $(\hat{\tau}_{\text{context}}, \hat{\tau}_{\text{long}})$ -plane cut through the parameter space at the best-fitting $\nu = \hat{\nu}^*$ and $K = \hat{K}^*$ (gray area indicates timescales within an order of magnitude of the end of the experiment). Contours show the first 10 contours incrementing by 0.01 error from the minimum (shown as a circle marker). Colors refer to subject, as in (b) and (c). (e) Same for $(\hat{\tau}_{\text{context}}, \hat{K}^*)$ at $\hat{\tau}_{\text{long}} = \hat{\tau}_{\text{long}}^*$ and $\nu = \hat{\nu}^*$. (f) Same for $(\hat{\tau}_{\text{context}}, \hat{\nu}^*)$ at $\hat{\tau}_{\text{long}} = \hat{\tau}_{\text{long}}^*$ and $K = \hat{K}^*$.

367 are also shown in [fig. 5a](#). They clearly fall below the optimal strategy, and, as expected,
 368 above the strategy that picks one of the three actions (report left, report right, and wait) at
 369 random.

370 To confirm that this similarity in performance between PGD and the data arises from
 371 a better fit to the behaviour than AR-RL, we plotted the distribution of the differences
 372 between model and data decision times, $|\Delta t_{\text{dec}}|$, conditioned on the context ([fig. 5b,c](#)). For
 373 comparison with previous work [7] and to account for deliberation cost in AR-RL, we added
 374 to the AR-RL reward objective a constant auxiliary deliberation cost rate, c , incurred up
 375 to the decision time in each trial, and chose the cost rate, c^* , that gave the lowest mean
 376 difference. In both contexts, PGD exhibits lower error than this c^* AR-RL solution.

377 To reveal the source of this discrepancy in both performance and behaviour, we turned

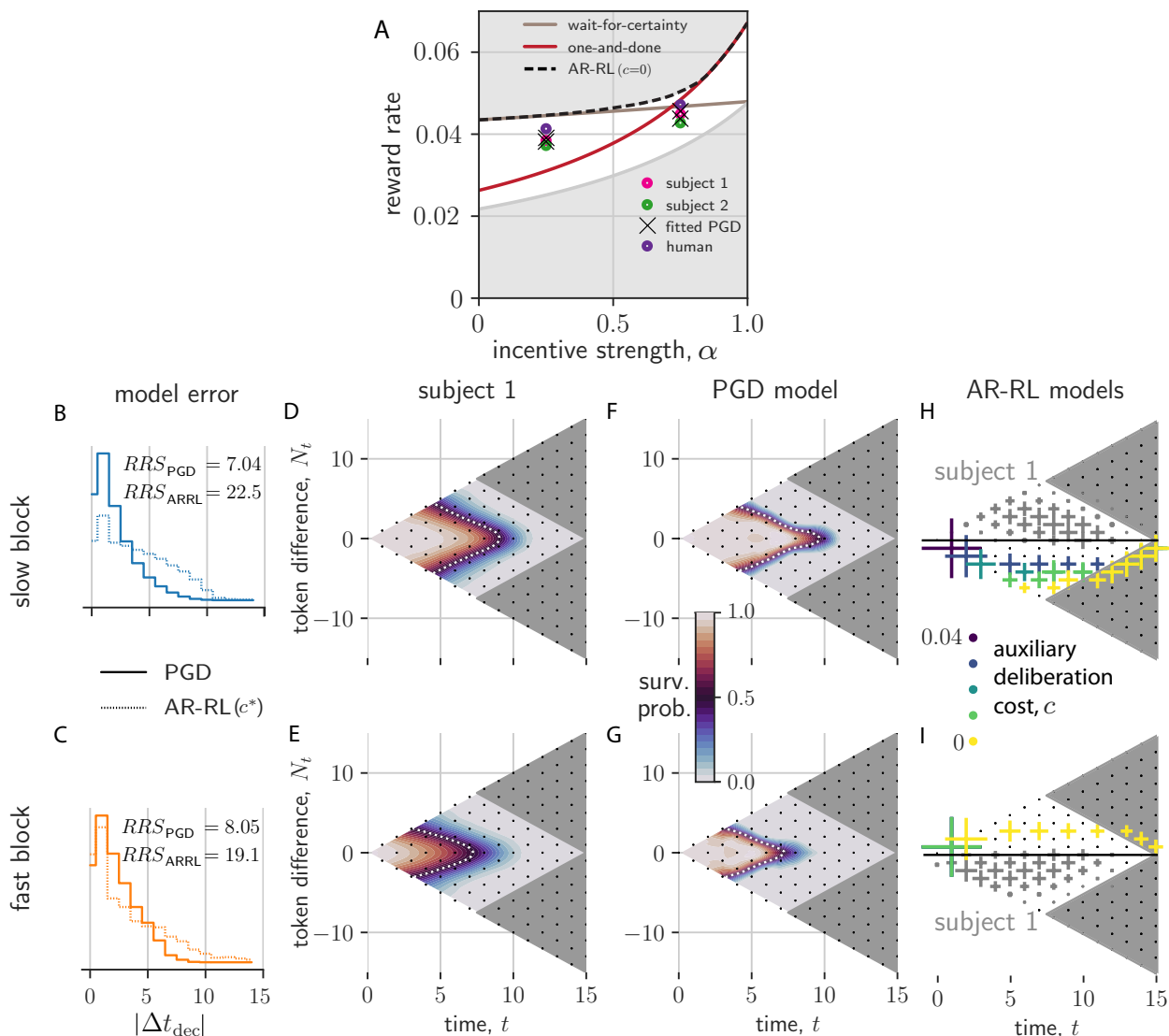


Figure 5. *Context-conditioned analysis of PGD and comparison to AR-RL models.* (a) Shown is the reward rate as a function of incentive strength, α (wait-for-certainty strategy shown in brown; one-and-done strategy shown in red). We additionally show the slow and fast context-conditioned reward rates for the two primates and the PGD model fitted to them, as well as a reference expert human. Reward rates for the human and non-human primates are squarely in between the best (black dashed) and uniformly random (gray) strategy. (b,c) The distribution over trials of differences in decision times between model and data, $|\Delta t_{\text{dec}}| = |t_{\text{dec,data}} - t_{\text{dec,model}}|$, conditioned on slow and fast block contexts. Solid lines are for PGD. Dotted lines are for the AR-RL solution using the cost rate, c^* , with the lowest mean error. The residual sum of squares (RRS) for each model/block combination is displayed. (d-g) Interpolated state-conditioned survival probabilities, $P(t^{\text{dec}} = t | N_t, t)$, over slow (d,f) and fast (e,g) blocks. White dotted lines show the $P(t^{\text{dec}} = t | N_t, t) = 0.5$ contour. (h,i) State-conditioned decision time frequencies (cross size) from AR-RL optimal decision boundaries across different values of the cost rate, c (colored crosses) for slow (h) and fast (i) conditions. Only samples with $N_t < 0$ and $N_t > 0$, respectively, are shown. For comparison, the reflected axes shows as gray crosses the state-conditioned decision time frequencies of the data.

378 to analyzing the corresponding policies of PGD and c -based AR-RL agents. A robust and
379 rich representation of the behavioural statistics is the state and time-conditioned survival
380 probability that a decision has not yet occurred. It serves as a summary of the action policy
381 associated with a stationary strategy (see [Methods](#) for its calculation from response times).
382 Applied equally to the decision times of both model and data, it can provide a means of
383 comparison even in this non-stationary setting. We give this conditional probability for
384 each of the two contexts for subject 1 and its fitted PGD model in [fig. 5d-g](#). We left the
385 many possible noise sources underlying the behaviour out of the model in order to more
386 clearly demonstrate the PGD algorithm. However, such noise sources would be necessary
387 to quantitatively match the variability in the data (e.g. added noise in the performance
388 estimates leads to larger variability in the location of the decision boundary and thus also
389 to larger spread in these survival probability functions (not shown)). In the absence of
390 these noise sources, we see the model underestimates the spread of probability over time
391 and tokens state. Nevertheless, the remarkably smooth average strategy is well captured by
392 the model (white dashed lines in [fig. 5d-g](#)). Specifically, policies approximately decide once
393 either of the peripheral targets receive a certain number of tokens. Comparing results across
394 context, we find that fast block strategies ([fig. 5e,g](#)) exhibit earlier decision times relative
395 to slow block strategies ([fig. 5d,f](#)) in both model and data. The strategies for subject 2
396 are qualitatively similar, but shifted to earlier times relative to subject 1 ([fig. S3](#)). Our
397 model explains this inter-individual difference as resulting from subject 2's larger reward
398 bias and faster context integration (*c.f.* [fig. 4e](#)). The correspondence between the PGD
399 model and data over the many token states in [fig. 5d-g](#) explains their similar performance
400 (*c.f.* [fig. 5a](#)). This similarity in policy is remarkable given that the model has essentially
401 only a single, crucial degree of freedom (τ_{context}), *a priori* unrelated to how decision times
402 depend on token state. Note that in both the fitted PGD model and the primate behaviour,
403 residual ambiguity ($N_t \approx 0$) is resolved at intermediate trial times ([fig. 5b-e](#)).

404 The AR-RL strategies are plotted across c in [fig. 5g,h](#). In contrast, they give no interme-
405 diate decision times at ambiguous ($N_t \approx 0$) states, invariably waiting until the ambiguity
406 resolves. This in fact holds over the entire (α, c) -plane (see [fig. S9](#) for the complete depen-
407 dence), and also under the addition of a movement cost, i.e. a constant cost incurred by
408 either of the reporting actions (data not shown). Thus, whereas AR-RL policies shift around
409 the edges of the relevant decision space as α or c is varied, the PGD policy lies squarely
410 in the bulk, tightly overlaying the policy extracted from the data. We conclude that the
411 context-conditioned strategies of the non-human primates in this task are well-captured by
412 PGD, while having little resemblance to the behaviour that would maximize reward rate
413 with or without a fixed deliberation cost rate. We address the additional freedom of a
414 time-varying cost rate in the discussion.

415 3. Neural urgency and context-dependent opportunity cost

416 So far, we have fit and analyzed the PGD model with respect to recorded behaviour. Here,
417 we take a step in the important direction of confronting the above theory of behaviour with
418 the neural dynamics that we propose drive it. The proposal for the tokens task mentioned
419 at the end of the introduction has evidence strength and urgency combining in PMd, whose
420 neural dynamics implements the decision process. In [fig. 6a](#), we restate in a schematic
421 diagram an implementation of this dynamics that includes a collapsing decision boundary.
422 In the one-dimensional belief space for the choice ([fig. 6a\(top\)](#)) [[7](#), [32](#)], the rising belief

423 collides with the collapsing boundary to determine the decision time. In the equivalent
424 commitment and deliberation cost formulation developed here (fig. 6a(middle)), the falling
425 commitment cost collides with the rising deliberation cost. The collapsing boundary in
426 belief space can be parametrized as $C - u_t$, where C is the initial strength of belief, e.g.
427 some desired confidence, that is lowered by a growing function of trial time $u_t > 0$. The
428 decision criterion is then $b_t > C - u_t$, where b_t is the belief, i.e. the probability of a correct
429 report. For AR-RL optimal policies, u_t emerges from value maximization and thus has a
430 complicated dependence on the opportunity cost sequence, $\mathcal{C}_t^{\text{del}}$. For PGD, in contrast, C
431 is interpreted as the maximum reward r_{max} and u_t is identically $\mathcal{C}_t^{\text{del}}$. For a linear neural
432 encoding model in which belief, rather than evidence, is encoded in neural activity, the sum
433 of the encoded belief \tilde{b}_t and the encoded collapsing boundary, \tilde{u}_t , evolve on a one-dimensional
434 choice manifold. According to the proposal, when this sum becomes sufficiently large (e.g.
435 $\tilde{b}_t + \tilde{u}_t > \tilde{C}$ for some threshold \tilde{C}), PMd begins to drive the activity in downstream motor
436 areas towards the associated response.

437 Neural urgency was computed from the PMd recordings of [19] in [33]. This computation
438 relies on the assumption that while a single neuron's contribution to \tilde{b}_t will depend on
439 its selectivity for choice (left or right report), the urgency \tilde{u}_t is a signal arising from a
440 population-level drive to all PMd neurons, irrespective of their selectivity. Thus, \tilde{u}_t can
441 be extracted from neural recordings by conditioning on zero-evidence states ($\tilde{b}_t = 0$) and
442 averaging over cells. In [33], error bars were computed at odd times via bootstrapping; data
443 at even times was obtained by interpolating between $N_t = \pm 1$; and data was pooled from
444 both subjects. We have excluded times at which firing rate error bars exceed the range
445 containing predictions from both blocks. To assess the correspondence of the components
446 of the deliberation cost developed here and neural urgency, in fig. 6b we replot their result
447 (c.f. fig.8b of [33]). We overlay the mean (+/- standard deviation) of the opportunity cost
448 sequence, $\mathcal{C}_t^{\text{del}}$ (shaded area in fig. 4; averaged over all trials produced by applying the two
449 fitted PGD models on the data sequence and conditioning the resulting average within-
450 trial deliberation cost on context). To facilitate our qualitative comparison, we convert
451 cost to spikes/step simply by adjusting the y-axis of the deliberation cost. The observed
452 urgency signals then lie within the uncertainty of the context-conditioned deliberation cost
453 signals computed from the fitted PGD models. There are multiple features of the qualitative
454 correspondence exhibited in fig. 6b: (1) the linear rise in time; (2) the same slope across
455 both fast and slow conditions; and (3) the baseline offset between conditions, where the fast
456 condition is offset to higher values than the slow condition. Such features would remain
457 descriptive in the absence of a theory. With the theory we have presented here, however,
458 each has their respective explanations via the interpretation of urgency as the opportunity
459 cost of deliberation: (1) the subject uses a constant cost per token jump, (2) this cost rate
460 refers to moment-to-moment decisions, irrespective of context, that is reflective of the use
461 of the context-agnostic stationary reward, and (3) trial-aware planning over contexts leads
462 to an opportunity cost baseline offset with a sign given by the reward rate deviation $\rho_\alpha - \rho$
463 with respect to the stationary average, ρ .

464 Up to now, the computational and neural basis for urgency has remained largely un-
465 explored in normative approaches, which also typically say little about adaptation effects
466 (see [34] for a notable exception). In summary, we exploited the adaptation across context
467 switches to learn the model and explained earlier responses in high reward rate contexts
468 as the result of a higher opportunity cost of deliberation. While this qualitative effect is
469 expected, we go beyond existing work by quantitatively predicting the average dependence

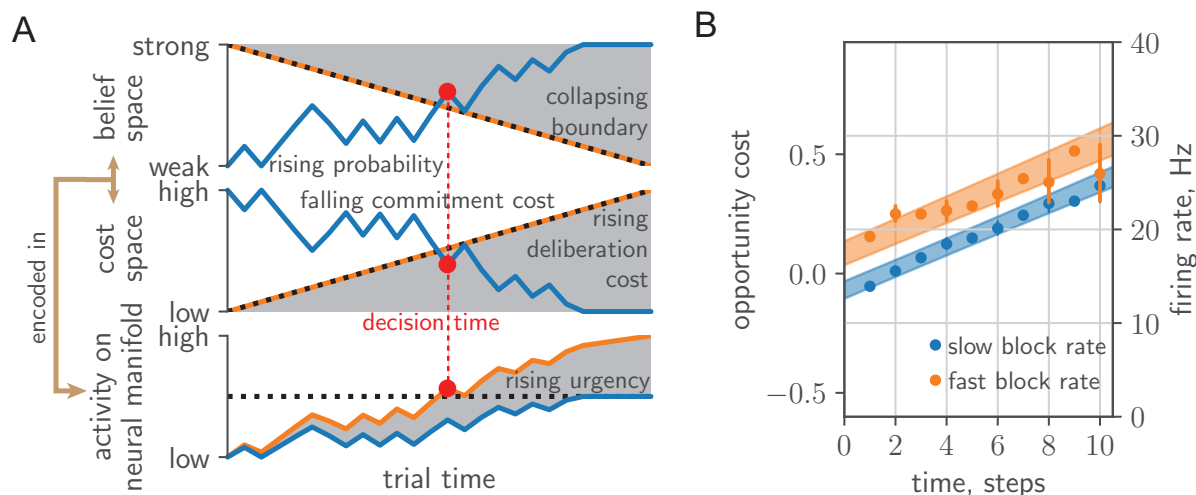


Figure 6. *Comparing neural urgency and collapsing decision boundaries.* (a) Top: Rising belief (blue) meets collapsing decision boundary (black dashed) in belief space. Middle: Falling commitment cost (blue) meets rising deliberation cost (black-dashed) in cost space. Bottom: Belief/commitment cost is encoded (blue) into a low-dimensional neural manifold, with the addition of an urgency signal (orange) (*c.f.* fig.8 in [7]). The decision (red circle) is taken when the sum passes a fixed threshold (black-dashed). (b) Deliberation cost maps onto the urgency signal extracted from zero-evidence conditioned cell-averaged firing rate in PMd (200ms time steps).

470 on both time and state (fig. 5b-e) as well as the qualitative form of urgency signal (fig. 6b).
 471 Taken together, the data is thus consistent with our interpretation that neural activity un-
 472 derlying context-conditioned decisions is gated by opportunity costs reflective of a trial-aware
 473 timescale hierarchy computed using performance estimation on multiple timescales.

474

DISCUSSION

475 We introduced PGD, a heuristic decision-making algorithm for continuing tasks that
 476 gates deliberation based on performance. We constructed a foraging example for which
 477 PGD is the optimal strategy with respect to the average-adjusted value function of average-
 478 reward reinforcement learning (AR-RL). While this will not be true in general, PGD does
 479 strike a balance between strategy complexity and return. The PGD decision rule does not
 480 depend on task specifics and exploits the stationarity of the environment statistics while
 481 simultaneously hedging against longer term non-stationarity in reward context. It does so
 482 by splitting the problem into two fundamental components—learning the statistics of the
 483 environment in order to compute the opportunity cost of commitment, and tracking one’s
 484 own performance in that environment with which to compute the opportunity cost of de-
 485 liberation. This splitting is not only crucial to making efficient use of the opportunity cost
 486 of time in non-stationary settings. Building on the field’s current understanding of how the
 487 cortico-basal ganglia system supports higher-level decision-making [35], we propose that the
 488 cost of deliberation arises from performance estimated on multiple, behaviourally-relevant
 489 timescales and is broadcast to multiple, lower-level decision-making areas to gate the speed
 490 of their respective evidence-driven attractor dynamics. Incorporating this cost into existing

491 models of such dynamics [32, 36, 37] is an interesting direction for future work. Consistent
492 with this picture, PGD’s explanatory power was borne out at both the behavioural and
493 neural levels for the tokens task data we analyzed. In particular, a deliberation cost con-
494 structed from trial-aware planning was supported independently by both these data sources.
495 We used behavioural data to fit and validate the theory, and neural recordings to provide
496 evidence of one of the neural correlates it proposes: the temporal profile of neural urgency.

497 *Scientific and clinical implications* In our proposal, we have linked two important and
498 related, but often disconnected fields: the systems neuroscience of the neural dynamics of
499 decision-making and the cognitive neuroscience of opportunity cost and reward sensitivity.
500 The view that tonic dopamine encodes average reward is two decades old [3]. However, the
501 existence of a reward representation decomposed by timescale has received increasing em-
502 pirical support only in recent years, from cognitive results [38–40] to a recent unified view of
503 how dopamine encodes reward prediction errors using multiple discount factors [41, 42] and
504 of dopamine as encoding both value and uncertainty [43]. Dopamine’s effect on time per-
505 ception has been proposed [44] and has empirical support [45], but the mechanism by which
506 its putative effect on decision speed is implicated in the neural dynamics of the decision-
507 making areas driving motor responses was unknown. Our theory fills this explanatory gap
508 by considering dynamic evidence tasks and parametrizing urgency using a multiple-timescale
509 representation of performance. One candidate for the latter’s neural implementation is in
510 the complex spatio-temporal filtering of dopamine via release-driven tissue diffusion and
511 integration via DR1 and DR2 binding kinetics [46]. Subsequent neural filtering and compu-
512 tation by striatal network activity could also play a role [47]. The study of spatiotemporal
513 filtering of dopamine is increasingly accessible experimentally [48, 49] and provides an excit-
514 ing direction for multiscale analysis of behaviour. Our proposal that urgency is the means
515 by which the neural representation of reward ultimately affects neural dynamics in decision-
516 making areas frames a timely research question on which these experimental methods could
517 shed light.

518 We applied PGD to decisions playing out in PMd, a decision-making area relevant to arm
519 movements. PGD appears to be relevant to other kinds of decisions, however. For instance,
520 a large body of work has studied decisions through recordings in lateral intraparietal cortex
521 in random dot motion tasks whose environment is formally similar to that of the tokens task.
522 One seminal study identified an urgency signal with the same properties as those exhibited
523 by the tokens task: a linear rise at early trial times that is independent of trial evidence
524 and an offset with sign given by the reward rate deviation of the current context, here two
525 and four-choice trials [17]. While decision boundaries obtained using AR-RL are evidence-
526 independent, these models require tailored cost functions that are fit to those experiments
527 in a procedure that assumes optimality *a priori* [7]. Here, we offer an alternative explanation
528 that behaviour is in fact suboptimal, with the decision boundary determined directly by the
529 estimated opportunity cost only. PGD decision boundaries are thus independent of evidence
530 by construction. In contrast to the tokens task, however, context in these random dot task
531 experiments was sampled randomly and thus its dynamics lacked temporal correlation [17].
532 In this case, a natural hypothesis from our approach is that a pair of performance filters,
533 one for each context, tracks the reward history in two parallel streams. In this case, our
534 theory would predict that the ratio of slopes of urgency across the two contexts reflects
535 the ratio of context-conditioned reward rates. An estimation procedure described in the
536 [Methods](#) for this data [17] agrees to within 20% error, providing support for the hypothesis
537 that PGD underlies non-human primate behaviour on this widely-studied task. Within the

538 context of the drift-diffusion models typically used to understand neural activity for that
539 task, PGD provides a principled mechanism that implements collapsing decision boundary.
540 PGD is thus easily incorporated into such models and testing the generality of our theory
541 using tailored experiments in this setting is an important next step.

542 Urgency may play a role in both decision and action processes, potentially providing a
543 transdiagnostic indicator of a wide range of cognitive and motor impairments in Parkinson’s
544 disease and depression [50]. Our theory offers a means to ground these diverse results in
545 neural dynamics by formulating opportunity cost estimation as the underlying causal factor
546 linking vigor impairments (e.g. in Parkinson’s disease) and dysregulated dopamine signalling
547 in the reward system [50–52]. We provide a concrete proposal for a signal filtering system
548 that extracts a context-sensitive opportunity cost from a reward prediction error sequence
549 putatively encoded by dopamine. Neural recordings of basal ganglia provide a means to
550 identify the neural substrate for this system.

551 *Commitment cost estimation* Beyond the estimation of the opportunity cost of deliber-
552 ation, we assumed that the agent had a precise estimate of the expected reward, which it
553 used to compute the within-trial commitment cost. For the tokens task, a recorded signal in
554 dorsal lateral prefrontal cortex of non-human primates correlates strongly with belief [20],
555 equivalent to the expected reward for binary rewards). How this quantity is computed by
556 neural systems is not currently known. However, for a general class of tasks, a generic,
557 neurally plausible means to learn the expected reward is via distributional value codes [43].
558 For example, the Laplace code is a distributional value representation that uses an ensemble
559 of units over a range of temporal discount factors and reward sensitivities [53]. The authors
560 show that expected reward is linearly decodeable from this representation.

561 *Experimental predictions* A feature of our decision-making theory is that it is highly
562 vulnerable to falsification. First, with regards to behaviour via the shape of the action
563 policy using our survival probability representation (*c.f.* fig. 5b-e,g,h), PGD varies markedly
564 with reward structure and thus provides a wealth of predictions for how observed behaviour
565 should be altered by it. For example, a salient feature of the standard tokens task is its
566 reflection symmetry in the tokens difference, N_t . We can break this symmetry for which the
567 theory predicts a distinctly asymmetric shape (fig. S10; for details see Methods). Our theory
568 is also prescriptive for neural activity via the temporal profile of neural urgency. The slope of
569 $\mathcal{C}_t^{\text{del}}$ remained fixed across blocks for relatively short block lengths used in the data analyzed
570 here. In the opposite limit, $T_{\text{block}}/\tau_{\text{long}} \gg 1$, $\rho_k^{\tau_{\text{long}}}$ approaches ρ_α except when undergoing
571 large, transient excursions after context switches. Thus, the deliberation cost is given by
572 the first component in eq. (5) most of the time, with the context specific reward rate as the
573 slope. One simple prediction is that the slope of urgency should exhibit increasing variation
574 as the duration of the blocks increases.

575 *Reinforcement learning theory* We suggest how to generalize average-adjusted value
576 functions to context-varying opportunity cost of time in a way that reduces to AR-RL
577 when context is fixed or not tracked. This adds a continuing task perspective to episodic
578 AR-RL, in line with recent work in machine learning, which is arguably the more appropriate
579 reinforcement learning setting for many decision-making experiments in neuroscience. The
580 epistemic perspective entailed in the estimation of these costs parallels a recent epistemic
581 interpretation of the discount-reward formulation as encoding knowledge about the volatility
582 of the environment [54].

583 Our work also suggests a new class of reinforcement learning algorithms between model-
584 based and model-free: only parts of the algorithm need adjustment upon task structure

585 variation. This is reminiscent of how the effects of complex state dynamics are decoupled
586 from reward when using a successor representation [55], but tailored for the average-reward
587 rather than the discount-reward formulation. We have left analysis of the algorithmic com-
588 plexity of PGD to future work, but expect performance improvements, as with successor
589 representations, in settings where decoupling the learning of environment statistics from the
590 learning of reward structure is beneficial.

591 *Comparison with humans* In the space of strategies, PGD lies in a regime between fully
592 exploiting assumed task knowledge (average-case optimal) and assumption-free adaptation
593 (worst-case optimal). Highly incentivized human behaviour is likely to be more structured
594 than PGD because of access to more sophisticated learning. While some humans land on
595 the optimal one-and-done policy in the fast condition when playing the tokens task [56],
596 most do not. The human brain likely has all the components needed to implement PGD.
597 Nevertheless, the situations in which we actually exploit PGD, if any, are as yet unclear. In
598 particular, how PGD and AR-RL relate to existing behavioural models tailored to explain
599 relative-value, context-dependent decision-making in humans [4], such as scale and shift
600 adaptation[57], is an open question. Whether or not PGD is built into our decision-making,
601 the question remains if PGD is optimal with respect to some bounded rational objective.
602 In spite of the many issues with the latter approach [58], using it to further understand the
603 computational advantages of PGD is an interesting direction for future work.

604 Despite our putative access to sophisticated computation, humans still exhibit measurable
605 bias in how we incorporate past experience [59]. One simple example is the win-stay/lose-
606 shift strategy, a more rudimentary kind of performance-gated decision-making than PGD,
607 which explains how humans approach the rock-paper-scissors game [60]. In that work,
608 numerical experiments demonstrated that this strategy outperforms at a population level the
609 optimal Nash equilibrium for this game, demonstrating that the use of such seemingly sub-
610 optimal strategies can confer a surprising evolutionary advantage. This example supports
611 the claim that relatively simple and nimble strategies such as PGD make for attractive
612 candidates when acknowledging that a combination of knowledge and resource limitations
613 over task, development, and evolutionary timescales have shaped decision-making in non-
614 stationary environments.

615

METHODS

616 Code for simulations and main figure generation (written in Python 3) is publicly acces-
617 sible as a online repository: https://github.com/mptouzel/dyn_opp_cost/.

618

Patch leaving task

619 We devised a mathematically tractable patch leaving task for which PGD learning is
620 optimal with respect to the average-adjusted value function. Here the value is simply the
621 return from the patch. This value function is related, but not equivalent to the marginal
622 value of optimal foraging, for which the decision rule is $\mathcal{C}_t^{\text{del}} > r_{\text{max}} - \mathcal{C}_t^{\text{com}} = \bar{r}_t$ [5]. This
623 choice of task allowed us to compare PGD's convergence properties relative to conventional
624 AR-RL algorithms that make use of value functions. In contrast to PGD, the latter requires
625 exploration. For a comparison generous to the AR-RL algorithm, we allowed it to circumvent
626 exploration by estimating the value function from off-policy decisions obtained from the

627 PGD algorithm using the same learning rate. We then compared them to PGD using their
628 on-policy, patched-averaged reward. This made for a comparison based solely between the
629 parameters of the respective models. If we did not allow for this, the AR-RL algorithms
630 would have to find good learning signals by exploring. In any form, this exploration would
631 lead them converge substantially slower. This setting thus provides a lower bound on the
632 convergence times of the AR-RL algorithm.

633 In this task, the subject randomly samples (with replacement) d patches, each of a dis-
634 tinct, fixed, and renewable richness defined by the maximum return conferred. These maxi-
635 mum returns are sampled before the task from a richness distribution, $p(r_{\max})$, with $r_{\max} > 0$
636 and are fixed throughout the experiment. The trials of the task are temporally extended
637 periods during which the subject consumes the current patch. After a time t in a patch,
638 the return is defined $r(t) = r_{\max}(1 - (\lambda t)^{-1})$. This patch return profile, $1 - (\lambda t)^{-1}$, is shared
639 across all patches and saturates in time with rate λ , a parameter of the environment that
640 sets the reference timescale. The return diverges negatively for vanishing patch leaving times
641 for mathematical convenience, but also evokes situations where leaving a patch soon after
642 arriving is prohibitively costly (e.g. when transit times are long). A stationary policy is then
643 a leaving time, t_s , for each of d patches, where the s -subscript indexes the patch. Given any
644 policy, the stationary reward rate for uniformly random sampling of patches is then defined
645 as

$$\rho = \frac{\sum_{s=1}^d r_s(t_s)}{\sum_{s=1}^d t_s}. \quad (6)$$

646 We designed this task to (1) emphasize the speed-return trade-off typical in many delibera-
647 tion tasks, and (2) have a tractable solution with which to compare convergence properties
648 of PGD and AR-RL value function learning algorithms.

649 A natural optimal policy is the one that maximizes the average-adjusted trial return,
650 $Q(r, t) = r - \rho t$. Given the return profile we have chosen, the corresponding optimal decision
651 time, t_s^* , in the s th patch obtained by maximizing $r - \rho t$ is $t_s^* = \sqrt{r_{\max,s}/(\lambda\rho)}$, which scales
652 inversely with the reward rate so that decision times are earlier for larger reward rates,
653 because consumption (or more generally deliberation) at larger reward rates costs more. We
654 chose this return profile such that stationary PGD learning gives exactly the same decision
655 times: the condition $C_t^{\text{del}} = C_t^{\text{com}}$ for patch s here takes the form $\rho t_s = r_{\max,s}/(\lambda t_s)$. Thus,
656 they share the same optimal reward rate, ρ^* . Using t_s^* for each patch in eq. (6) gives a
657 self-consistency equation for ρ with solution $\rho^* = \lambda\mu_1^2/4\mu_{1/2}^2$, where $\mu_n = \langle r_{\max}^n \rangle_{p(r_{\max})}$ (we
658 have assumed d is large here to remove dependence on s). Described so far in continuous
659 time, the value function was implemented in discrete time such that the action space is
660 a finite set of decision times selected using the greedy policy, $t^* = \operatorname{argmax}_t \hat{Q}(r, t)$, where
661 $\hat{Q}(r, t)$ is the estimated trial return. As a result, there is a finite lower bound on the
662 performance gap, i.e. the relative error, $\epsilon = (\rho^* - \rho)/\rho^* > 0$ for the AR-RL algorithm.
663 Approaching this bound, convergence time for both PGD and AR-RL learning is limited
664 by the integration time τ of the estimate $\hat{\rho}_k^r$ (c.f. eq. (8)) of ρ . We note that PGD learns
665 faster in all parameter combinations tested. To demonstrate the insensitivity of PGD to the
666 state space representation, at 5×10^5 time steps into the experiment we shuffled the labels
667 of the states. PGD is unaffected, while the value function-based AR-RL algorithm is forced
668 to relearn and in fact does so slower than in the initial learning phase, due to the much
669 larger distance between two random samples, than between the initial values (chosen near
670 the mean) and the target sample.

671

Filtering performance history

672 For unit steps of discrete time, the step-wise update of the performance estimate, $\hat{\rho}_t^\tau$, is

$$\hat{\rho}_t^\tau = (1 - \beta)\hat{\rho}_{t-1}^\tau + \beta R_t, \quad (7)$$

673 with $\beta = 1/(1 + \tau)$ called the learning rate, and τ the characteristic width of the exponential
 674 window of the corresponding continuous time filter over which the history is averaged. We
 675 add τ as a superscript when denoting the estimate to indicate this. Exceptionally, here t
 676 indexes absolute time rather than trial time. Note that a continuous-time formulation of
 677 the update is possible via an event-based map given the decision times in which the reward
 678 event sequence is given as a sum of delta functions. In either case, to leading order in β ,
 679 $\hat{\rho}_t^\tau \approx \beta \sum_i^t R_i$, i.e. the filter sums past rewards. Thus, when $\tau \sim \mathcal{O}(t) \gg 1$, $\beta \sim \mathcal{O}(1/t) \ll 1$
 680 and so $\hat{\rho}_t^\tau \approx \beta \sum_i^t R_i \rightarrow \rho$ when t is large.

681 The rewards in this task are sparse: $R_t = 0$ except when a trial ends and the binary
 682 trial reward R_k (1 or 0) is received. A cumulative update of eq. (7) that smooths the
 683 reward uniformly over the trial duration and is applied once at the end of each trial is
 684 thus more computationally efficient. Resolving a geometric series leads to the cumulative
 685 update [8, 28]

$$\hat{\rho}_k^\tau = (1 - \beta)^{T_k} \hat{\rho}_{k-1}^\tau + (1 - (1 - \beta)^{T_k}) \rho_k^{\text{trial}}, \quad (8)$$

686 where the smoothed reward, $\rho_k^{\text{trial}} = R_k/T_k$, can be interpreted as a trial-specific reward rate.
 687 The initial estimate, $\hat{\rho}_0^\tau$, is set to 0. Exceptionally, $\hat{\rho}_1^\tau = R_1/T_1$, after which eq. (8) is used.
 688 Using the first finite sample as the first finite estimate is both more natural and robust than
 689 having to adapt from zero. We will reuse this filter for different τ and denote the filtered
 690 estimate from its application with a τ -superscript, $\hat{\rho}_k^\tau$. For example, the precision of $\hat{\rho}_k^{\tau_{\text{long}}}$
 691 as an estimate of a stationary reward rate ρ is set by how many samples it averages over,
 692 which is determined by the effective length of its memory given by τ_{long} . Since we assume
 693 the subject has learned the expected reward, \bar{r}_t , we use it instead of R_k when computing
 694 ρ_k^{trial} .

695

Tokens task: a random walk formulation

696 The tokens task is a continuing task of episodes (here trials), which can be formulated
 697 using the token difference, N_t . Each trial effectively presents to the agent a realization
 698 of a finite-length, unbiased random walk, $\mathbf{N}_{t_{\text{max}}} = (N_0, \dots, N_{t_{\text{max}}})$ with $N_t = \{-t, \dots, t\}$
 699 and $N_0 = 0$. We express time in units of these steps. The agent observes the walk and
 700 reports its prediction of the sign of the final state, $\text{sign}(N_{t_{\text{max}}}) = \pm 1$ (t_{max} is odd to exclude
 701 the case it has no sign). The time at which the agent reports is called the decision time,
 702 $t^{\text{dec}} \in \{0, 1, \dots, t_{\text{max}}\}$. For a greedy policy, $\text{sign}(N_t)$ can be used as the prediction (and
 703 the reporting action selected randomly if $N_{t^{\text{dec}}} = 0$). The decision-making task then only
 704 involves choosing when to decide. In this case, the subject receives reward $R = \Theta(N_{t_{\text{max}}} N_{t^{\text{dec}}})$
 705 at the end of the random walk, i.e. a unit reward for a correct prediction, otherwise nothing
 706 (Θ is the Heaviside function: $\Theta(x) = 1$ if $x > 0$, zero otherwise).

707 An explicit action space beyond decision time is not necessary for the case of greedy
 708 actions. It can nevertheless be specified for illustration in an Markov decision process (MDP)
 709 formulation: the agent waits ($a_t = 0$ for $t < t^{\text{dec}}$) until it reports its prediction, $a_{t^{\text{dec}}} = \pm$,
 710 after which actions are disabled and the prediction is stored in an auxiliary state variable

711 used to determine the reward at the end of the trial. A MDP formulation for a general class
 712 of perceptual decision-making tasks, including the tokens and random dots task, is given in
 713 [Methods](#)).

714 Perfect accuracy in this task is possible if the agent reports at t_{\max} since $R = \Theta(N_{t_{\max}}^2) =$
 715 1. The task was designed to study reward rate maximizing policies. In particular, the task
 716 has additional structure that allows for controlling what this optimal policy is through the
 717 incentive to decide early, α , incorporated into the trial duration for deciding at time t in the
 718 trial,

$$T(t) = t + (1 - \alpha)(t_{\max} - t) + T_{\text{ITI}}. \quad (9)$$

719 Here, a dead time between episodes is added via the inter-trial interval, T_{ITI} , to make
 720 suboptimal the strategy of predicting randomly at the trial’s beginning. We emphasize that
 721 it is through the trial duration that α serves as a task parameter controlling the strength
 722 of the incentive to decide early. When α is fixed, we denote the corresponding optimal
 723 stationary reward rate, ρ_α , obtained from the reward rate maximizing policy. This policy
 724 shifts from deciding late to deciding early as α is varied from 0 to 1 (*c.f.* [fig. S9f,g](#)).

725 We consider a version of the task where α is variable across two episode types, a slow
 726 ($\alpha = 1/4$) and fast ($\alpha = 3/4$) type. The agent is aware that the across-trial α dynamics
 727 are responsive (maybe even adversarial), whereas the within-trial random walk dynamics
 728 (controlled by the positive jump probability, here $p = 1/2$) can be assumed fixed (see the
 729 next section for how p factors into the expression for the expected reward, \bar{r}_t).

730 Expected trial reward for the tokens task

731 We derived and used an exact expression for the expected reward in a trial of the tokens
 732 task. We derive that expression here as well as a simple approximation. The state sequence
 733 is formulated as a t_{\max} -length sequence of random binary variables, $\mathbf{S}_{t_{\max}} = (S_1, \dots, S_{t_{\max}})$,
 734 $S_t = \pm 1$, $i = 1, 2, \dots, t_{\max}$. Consider a simple case in which each is an independent and
 735 identically distributed Bernoulli sample, $P(s) = p^{\frac{1+s}{2}}(1-p)^{\frac{1-s}{2}}$, for jump probability $p \geq 1/2$.
 736 The distribution of $\mathbf{S}_{t_{\max}}$ is then

$$P(\mathbf{s}_{t_{\max}}) = \prod_{i=1}^{t_{\max}} P(s_i). \quad (10)$$

737 We will use this distribution to compute expectations of quantities over this space of trajec-
 738 tories, namely the sign of $N_t = \sum_{i=1}^t S_i$, for some $0 \leq t \leq t_{\max}$ and in particular the sign of
 739 the final state, $\xi := \text{sgn}(N_{t_{\max}}) \in \{+, -\}$ given $N_t = n$. Note that N_t is even if t is even and
 740 same with odd values. We remove the case of no sign in $N_{t_{\max}}$ by choosing t_{\max} to be odd.

741 First, consider predicting $\text{sgn}(N_t)$ with no prior information. The token difference, $-t \leq$
 742 $N_t \leq t$, appears directly in $P(\mathbf{s}_{t_{\max}})$. Marginalizing (here just integrating out) the additional
 743 degrees of freedom leads to a binomial distribution in the number of S_i for $i \leq t$ for which
 744 $S_i = +1$, $N_t^+ = \sum_{i=1}^t \Theta(s_i) = (t + N_t)/2$,

$$P(N_t^+ = n) = \binom{t}{n} p^n (1-p)^{t-n}, \quad (11)$$

745 with $n \in \{0, \dots, t\}$ and $N_t = 2N_t^+ - t$. Thus, the probability that $N_t > 0$, i.e. $N_t^+ > t/2$, is

$$P(N_t > 0) = \sum_{n=0}^t \binom{t}{n} p^n (1-p)^{t-n} \Theta(n - t/2). \quad (12)$$

746 Now consider predicting $\xi = \text{sgn}(N_{t_{\max}})$, given the observation $N_t = n$. Define $t' = t_{\max} - t$
747 as the remaining time steps to the predicted time and $N_{t'} = \sum_{i=t+1}^{t_{\max}} s_i$, i.e. the total count
748 in the remaining part of the realization. Then the probability of $\xi = +$ conditioned on the
749 state $N_t = n$, denoted $p_{n,t}$, is defined in the same way as $P(N_t > 0)$,

$$p_{n,t}^+ := P(\xi = + | N_t = n) = \sum_{n'=0}^{t'} \binom{t'}{n'} p^{n'} (1-p)^{t'-n'} \Theta(n' - (t' - n)/2). \quad (13)$$

750 where $N_{t'}^+ = n'$ is the number of positive jumps in the remaining $t' = t_{\max} - t$ steps and we
751 have used $N_{t_{\max}} = N_t + N_{t'} = N_{t'}^+ - (t' - N_t)/2$. The $\Theta(n' - (t' - n)/2)$ factor effectively changes
752 the lower bound of the sum to $\max\{0, \lceil (t' - n)/2 \rceil\}$, where $\lceil \cdot \rceil$ rounds up. If $\lceil (t' - n)/2 \rceil \leq 0$
753 then $p_{n,t}^+ = 1$ since the sum is over the domain of the distribution, which is normalized.
754 Otherwise, the lower bound is $\lceil (t' - n)/2 \rceil$, and the probability of $\xi = +1$ is

$$p_{n,t}^+ = \sum_{n'=\lceil (t'-n)/2 \rceil}^{t'} \binom{t'}{n'} p^{n'} (1-p)^{t'-n'}. \quad (14)$$

755 For odd t_{\max} , the probability that $\xi = -$ is denoted $p_{n,t}^- = 1 - p_{n,t}^+$. For the symmetric case,
756 $p = 1/2$,

$$p_{n,t}^+ = \frac{1}{2^{t'}} \sum_{n'=\lceil (t'-n)/2 \rceil}^{t'} \binom{t'}{n'}, \quad (15)$$

757 when $\lceil (t' - n)/2 \rceil > 0$ and 1 otherwise. This expression is equivalent to equation 5 in [16],
758 which was instead expressed using $N_{t'}^-$.

759 The space of trajectories, i.e. of $\mathbf{s}_{t_{\max}}$, maps to a space of trajectories for $p_{n,t}^+$ defined on
760 an evolving lattice in belief space. The expected reward in this case is,

$$\bar{r}_t := \langle r | N_t = n \rangle = \mathbb{E}[\Theta(N_{t_{\max}} N_t) | N_t = n] \quad (16)$$

$$= \max\{p_{n,t}^+, 1 - p_{n,t}^+\} \quad (17)$$

$$= b_t, \quad (18)$$

761 where the belief of correct report $b_t := \max\{p_{n,t}^+, 1 - p_{n,t}^+\}$. The commitment cost $\mathcal{C}_t^{\text{com}} =$
762 $r_{\max} - \bar{r}_t$, then also evolves on a lattice (see fig. 3(b)). More generally, $\bar{r}_t = \Delta r b_t + r_{\text{incorrect}}$
763 for Δr the difference of correct r_{correct} (here 1) and incorrect $r_{\text{incorrect}}$ (here 0) rewards. Since
764 $r_{\max} = r_{\text{correct}}$, we have $\mathcal{C}_t^{\text{com}} = \Delta r (1 - b_t)$. For $p = 1/2$ and $\Delta r = 1$, $\mathcal{C}_{t=0}^{\text{com}} = 1/2$.

765 The shape of $p_{n,t}^+$ is roughly sigmoidal, admitting the approximation,

$$p_{n,t}^+ \approx \frac{1}{1 + \exp[-(at + b)n]} \quad (19)$$

766 where fitting constants a and b depend on t_{\max} . For $t_{\max} = 15$, $a = 0.03725$ and $b = 0.3557$.
767 We demonstrate the quality of this approximation in fig. S5. Approximation error is worse
768 at t near t_{\max} . More than 95% of decisions times in the data we analyzed occur before
769 12 time steps, where the approximation error in probability is less than 0.05. A similar
770 approximation without time dependence was presented in [16]. We nevertheless used the
771 exact expression eq. (15) in all calculations.

772

PGD implementation and fitting to relaxation after context switches

773 We identified the times of the context switches in the data and their type (slow-to-fast
774 and fast-to-slow). Taking a fixed number of trials before and after each event, we averaged
775 the decision times over the events to create two sequences of average decision times around
776 context switches (the result is shown in [fig. 4a,b](#)). We used a uniformly weighted squared-
777 error objective, minimized with the standard (Nelder-Mead) simplex routine in python's
778 scientific computing library's optimization package.

779

Survival probabilities over the action policy

780 Behavioural analyses typically focus on response time distributions. From the perspective
781 of reinforcement learning, this is insufficient to fully characterize the behaviour of an agent.
782 Instead, the full behaviour is given by the action policy. In this setting, a natural represen-
783 tation of the policy is the probability to report as a function of both the decision time *and*
784 the environmental state (see [fig. 5](#)). These are computed from the histograms of $(N_{t^{\text{dec}}}, t^{\text{dec}})$,
785 over trials. However, the histograms themselves do not reflect the preference of the agent
786 to decide at a particular state and time because they are biased by the different frequencies
787 with which the set of trajectories visit each state and time combination. While there are
788 obviously the same number of trajectories at early and late times, they distribute over many
789 more states at later times and so each state at later times is visited less on average than states
790 at earlier times. We can remove this bias by transforming the data ensemble to the ensemble
791 of two random variables: the state conditioned on time $(N_t|t)$, and the event that $t = t^{\text{dec}}$.
792 Conditioning this ensemble on the state gives $P(t = t^{\text{dec}}|N_t, t) = p(N_t, t = t^{\text{dec}}|t)/p(N_t|t)$. To
793 reduce estimator variance, we focus on the corresponding survival function, $P(t < t^{\text{dec}}|N_t, t)$.
794 So, $P(t < t^{\text{dec}}|N_t, t) = 1$ when $t = 0$ and decays to 0 as t and $|N_t|$ increase. Unlike the
795 unconditioned histograms, these survival probabilities vary much more smoothly over state
796 and time. This justifies the use of the interpolated representations displayed in [fig. 5b-e](#).
797 Note that to simplify the analysis, we have binned decision times by the 200 ms time step
798 between token jumps. This is justified by the small deviations from uniformity of decision
799 times modulo the time step shown in [fig. S11](#).

800

Episodic decision-making and dynamic programming solutions of value iteration

801 We generalize the mathematical notation and description of an existing AR-RL formu-
802 lation and dynamic programming solution of the random dots task [7], a binary perceptual
803 evidence accumulation task extensively studied in neuroscience. To align notation with
804 convention in reinforcement learning theory, exceptionally here s denotes the belief state
805 variable, ie. a representation of the task state sufficient to make the decision (e.g. the to-
806 kens difference, N_t , in the case of the tokens task). We connect this extended formulation to
807 account for a dynamic deliberation cost. We write it in discrete time, though the continuous
808 time version is equally tractable.

809 The problem is defined by a recursive optimality equation for the value function $V(s|t)$
810 in which the highest of the action values, $Q(s, a|t)$, is selected. We formalize the non-
811 stationarity within episodes by conditioning on the trial time, t , where $t = 0$ is the trial start
812 time. $Q(s, a|t)$ is the action-value function of average-reward reinforcement learning [11], i.e.

813 the expected sum of future reward deviations from the average when selecting action a when
 814 in state s , at possible decision time t within a trial, and then following a given action policy
 815 π thereafter. The action set for these binary decision tasks consists of *report left* ($-$), *report*
 816 *right* ($+$), and *wait*. When *wait* is selected, time increments and beliefs are updated with
 817 new evidence. We use a decision-time conditioned, expected trial reward function, $r(s, a|t)$
 818 with $a = \pm$, that denotes the reward expected to be received at the end of the trial after
 819 having reported \pm in state s at time t during the trial. Note that $r(s, a|t)$ can be defined
 820 in terms of a conventional reward function $r(s, a)$ if the reported action, decision time, and
 821 current time are stored as an auxiliary state variable so they can be used to determine the
 822 non-zero reward entries at the end of the trial.

823 The average-reward formulation of $Q(s, a|t)$ naturally narrows the problem onto deter-
 824 mining decisions within only a single episode of the task. To see this, we pull out the
 825 contribution of the current trial,

$$Q(s, a|t) = \mathbb{E}^\pi \left[\sum_{t'=t}^T R_{t'} - \rho \mid S_t = s, A_t = a \right] + V(s|T+1) \quad (20)$$

826 where T is the (possibly stochastic) trial end time and $V(s|T+1)$ is the state value at the
 827 start of the following trial, which does not depend on s_t and a_t for independently sampled
 828 trials. Following conventional reinforcement learning notation, the expectation \mathbb{E}^π is over
 829 all randomness conditioned on following the policy, π , which itself could be stochastic [11].
 830 When trials are identically and independently sampled, the state at the trial start is the
 831 same for all trials and denoted s_0 with value V_0 . Thus, the value at the start of the trial
 832 $V(s|t=0) = V(s|T+1) = V_0$ equals that at the start of the next trial and so, by construction,
 833 the expected trial return (total trial rewards minus trial costs) must vanish (we will show
 834 this explicitly below). Note that the value shift invariance of eq. (20) can be fixed so that
 835 $V_0 = 0$.

836 The *optimality equation* for $V(s|t)$ arises from a greedy action policy over $Q(s, a|t)$: it
 837 selects the action of the largest of $Q(s, -|t)$, $Q(s, +|t)$, and $Q(s, wait|t)$. The value expression
 838 for the wait-action is incremental, and so depends on the value at the next time step. In
 839 contrast, expression for the two reporting actions integrates over the remainder of the trial
 840 since no further decision is made and so depends on the value at the start of the following
 841 trial. The resulting optimality equation for the value function $V(s|t)$ is then

$$\begin{aligned} V(s|t) &= \max_a Q(s, a|t) , \\ Q(s, \pm|t) &= r(s, \pm|t) - \sum_{t'=t+1}^T c_{t'} + V(s|t = T+1) , \\ Q(s, wait|t) &= -c_t + \mathbb{E}_{s_{t+1}|s} [V(s_{t+1}|t+1)] , \\ V(s|t=0) &= V(s|t = T+1) . \end{aligned} \quad (21)$$

842 Here, $t = 0, 1, \dots, t_{\max}$ within the current trial and $t = T+1, T+2 \dots$ in the following
 843 trial, with t_{\max} the latest possible decision time in a trial, and $T = T(t)$ the decision-time
 844 dependent trial duration. For inter-trial interval T_{ITI} , T satisfies $T_{\text{ITI}} \leq T \leq t_{\max} + T_{\text{ITI}}$.
 845 c_t is the cost rate at time t . The second term in $Q(s, wait|t)$ uses the notation $\mathbb{E}_{x|y}[z]$, i.e.
 846 the expectation of z with respect to $p(x|y)$. The last line in eq. (21) is the self-consistency
 847 criterion imposed by the AR-RL formulation, which demands that the expected value at

848 the beginning of the trial be the expected value at the beginning of the following trial. The
 849 greedy policy then gives a single decision time for each state trajectory as the first time when
 850 $Q(s, -|t) > Q(s, wait|t)$ or $Q(s, +|t) > Q(s, wait|t)$, with the reporting action determined
 851 by which of $Q(s, -|t)$ and $Q(s, +|t)$ is larger. For given c_t , dynamic programming provides
 852 a solution to eq. (21) [7] by recursively solving for $V(s|t)$ by back-iterating in time from the
 853 end of the trial. For most relevant tasks, to never report is always sub-optimal, so the value
 854 at $t = t_{\max}$ is set by the best of the two reporting (\pm) actions, which do not have a recursive
 855 dependence on the value and so can seed the recursion.

856 We now interpret this general formulation in terms of opportunity costs. For the choice
 857 of a static opportunity cost rate of time, $c_t = \rho$. This is the AR-RL case. As in [7], a
 858 constant auxiliary deliberation cost rate, c , incurred only up to decision time can be added,
 859 $c_t = \rho + c\Theta(t^{\text{dec}} - t)$. Of course, ρ is unknown *a priori*. For this solution method, its value
 860 can be found by exploiting the self-consistency constraint, $V(s|t = 0) = V(s|t = T + 1)$. This
 861 dependence can be seen formally by taking the action value eq. (20), choosing a according
 862 to π to obtain the state value, $V(s|t)$, and evaluating it for $t = 0$,

$$V(s|t = 0) = \mathbb{E}_{t^{\text{dec}}} \left[\sum_{t=0}^T R_t - \rho \right] + V(s|t = T + 1) \quad (22)$$

$$= \mathbb{E}_{t^{\text{dec}}} [r(t^{\text{dec}}) - \rho T(t^{\text{dec}})] + V(s|t = T + 1) \quad (23)$$

$$= \bar{R} - \rho \bar{T} + V(s|t = T + 1) . \quad (24)$$

863 Here, $\bar{R} = \mathbb{E}_{t^{\text{dec}}} [r(t^{\text{dec}})]$ and $\bar{T} = \mathbb{E}_{t^{\text{dec}}} [T(t^{\text{dec}})]$ denotes the expectations over the trial en-
 864 semble that, when given the state sequence, transforms to an average over t^{dec} , the trial deci-
 865 sion time, defined as when $V(s|t)$ achieves its maximum on the state sequence, $(s_0, \dots, s_{t_{\max}})$.
 866 The expected trial reward function, $r(t) := \max_{a \in \{-, +\}} r(s, a|t)$ is the expected trial reward
 867 for deciding at t . Imposing the self-consistency constraint on eq. (24) recovers the definition
 868 $\rho = \bar{R}/\bar{T}$.

869 Asymmetric switching cost model

870 Here, we present the model component that accounts for the asymmetric relaxation
 871 timescales after context switches. The basic assumption is that tracking a signal at a higher
 872 temporal resolution should be more cognitively costly, so that adapting from faster to slower
 873 environments should happen more quickly than the reverse, so as to not pay this cost un-
 874 necessarily. We now develop this idea formally (see fig. S4).

875 Let T_{track} and T_{sys} be the timescale of tracking and of the tracked system, respectively.
 876 One way to interpret the mismatch ratio, $T_{\text{sys}}/T_{\text{track}}$, is via an attentional cost rate, q .
 877 This rate should decay with T_{track} : the slower the timescale of tracking, the lower the
 878 cognitive cost. For simplicity, we set $q = 1/T_{\text{track}}$ (fig. S4a). Integrating this cost rate over a
 879 characteristic time of the system is then the tracking cost, $Q = qT_{\text{sys}} = T_{\text{sys}}/T_{\text{track}}$, which is
 880 also the mismatch ratio. We propose that Q enters the algorithm via a scale factor on the
 881 integration time of the reward filter for $\hat{\rho}_k^{\tau_{\text{context}}}$, τ_{context} . We redefine τ_{context} as

$$\tau_{\text{context}} \leftarrow \frac{\tau_{\text{context}}}{1 + Q^\nu} , \quad (25)$$

882 where ν is a sensitivity parameter that captures the strength of the nonlinear sensitivity of
 883 the speed up (for $\nu > 1$) or slow down (for $\nu < 1$) in adaptation with the tracking cost,

884 Q (fig. S4a shows how this timescale varies over Q for three values of ν). A natural choice
 885 for T_{sys} is T_k , the trial duration. For T_{track} , we introduce the filtered estimate of the trial
 886 duration, $\hat{T}_k^{\tau_{\text{context}}}$ (computed using the same simple low-pass filter *c.f.* eq. (8)). Thus, the
 887 tracking timescale adapts to the system timescale. As a result of how τ_{context} is lowered by Q
 888 for $\nu > 1$, this adaptation is faster in the fast-to-slow transition relative to the slow-to-fast
 889 transition.

890 Prediction for asymmetric rewards

891 Given a payoff matrix, $\mathbf{R} = (r_{s,a})$, where $r_{s,a}$ is the reward for reporting $a \in \{-, +\}$ in the
 892 trial realization leading to s , here the sign of $N_{t_{\text{max}}}$, and the probability that the rightward
 893 choice is correct, $p_{n,t}^+$, the expected reward for the two reporting actions in a trial is given
 894 by the matrix equation

$$\langle r|a = +, n, t \rangle \quad \langle r|a = -, n, t \rangle = \begin{bmatrix} p_{n,t}^+ & 1 - p_{n,t}^+ \end{bmatrix} \begin{bmatrix} r_{++} & r_{+-} \\ r_{-+} & r_{--} \end{bmatrix}.$$

895 Here, the corresponding reported choice is $a^* = \operatorname{argmax}_{a \in \{-, +\}} \langle r|a, n, t \rangle$. In this paper and
 896 in all existing tokens tasks, \mathbf{R} was the identity matrix. In this case, and for all cases where
 897 \mathbf{R} is a symmetric matrix, $\mathbf{R} = \mathbf{R}^\top$, an equivalent decision rule is to decide based on the sign
 898 of N_t . When \mathbf{R} is not symmetric, however, this is no longer a valid substitute. Asymmetry
 899 can be introduced through the actions and the states.

900 Using an additional parameter γ , we introduce asymmetry via a bias for $+$ actions that
 901 leaves the total reward unchanged by replacing the payoff matrix with

$$\mathbf{R}_{\text{asym}} = \begin{bmatrix} r_{++}(1 + \gamma) & r_{+-}(1 - \gamma) \\ r_{-+}(1 + \gamma) & r_{--}(1 - \gamma) \end{bmatrix},$$

902 The result for $\gamma = -0.6, 0$, and 0.6 is shown in fig. S10. For $\gamma > 0$ the decision boundary for
 903 $a = +$ shifts up proportional to γ . For $\gamma < 0$ the decision boundary for $a = -$ shifts down
 904 proportional to $-\gamma$. The explanation is that the components are set and exchange where
 905 the decision is exchanged, $N_t = 0$ for the symmetric case. This changes to $N_t \propto \pm\gamma$ for the
 906 asymmetric $\gamma \neq 0$ case.

907 Comparing reward rates and slopes of urgency

908 Reference [17] parametrize urgency with the saturation value, u_∞ , and the half-maximum,
 909 $\tau_{1/2}$. The initial slope is given by their ratio. We used the context-conditioned values
 910 published in Table 1 in [17] for the $n = 70$ (no 90° control) dataset. The context-conditioned
 911 reward rates, ρ_α , are computed as the accuracy $\langle R \rangle|_\alpha$ divided by the average trial time, $\langle T \rangle|_\alpha$
 912 for choice number $\alpha \in \{2, 4\}$ as context. We computed $\langle R \rangle|_{\alpha=2} = 0.71$ and $\langle R \rangle|_{\alpha=4} = 0.49$.
 913 The trial time is the sum of the response time, the added time penalty if incorrect, and the
 914 inter-trial interval. We computed the response times $t_{\text{response}, \alpha=2} = 0.527$ and $t_{\text{response}, \alpha=4} =$
 915 0.725 . While the dataset contains the response times, it does not have the latter two. The
 916 time penalty was on the order of 1 second, as was the time penalty [61]. Under those
 917 estimates, the reward rates are $\rho_{\alpha=2} = 0.40$ and $\rho_{\alpha=4} = 0.22$. The ratio between slopes is
 918 1.8 and the ratio of reward rates was 2.3 giving an error of about 20%.

919

ACKNOWLEDGMENTS

920 We would like to acknowledge helpful discussions with Jan Drugowitsch, Becket Ebitz,
921 and Paul Masset, and to Anne Churchland for sharing data from [17]. MPT acknowledges
922 support from IVADO via their postdoctoral fellowship award. PC acknowledges support
923 from NSERC Discovery Grant (RGPIN-2016-05245). GL acknowledges support from FRQS
924 Research Scholar Award, Junior 1 (LAJGU0401-253188), NSERC Discovery Grant (RGPIN-
925 2018-04821), and the Canada CIFAR AI Chair program.

-
- 926 [1] David I Green, “Pain-Cost and Opportunity-Cost,” *The Quarterly Journal of Economics* **8**,
927 **218–229** (1894).
- 928 [2] Vektor Dewanto, George Dunn, Ali Eshragh, Marcus Gallagher, and Fred Roosta, “Average-
929 reward model-free reinforcement learning: a systematic review and literature mapping,”
930 [arXiv:2010.08920 \[cs.LG\]](https://arxiv.org/abs/2010.08920).
- 931 [3] Nathaniel D Daw and David S Touretzky, “Long-term reward prediction in TD models of the
932 dopamine system,” *Neural computation* **14**, 2567–2583 (2002).
- 933 [4] Lindsay E Hunter and Nathaniel D Daw, “Context-sensitive valuation and learning,” *Current*
934 *Opinion in Behavioral Sciences* **41**, 122–127 (2021).
- 935 [5] Nils Kolling and Thomas Akam, “(Reinforcement?) Learning to forage optimally,” *Current*
936 *Opinion in Neurobiology* **46**, 162–169 (2017).
- 937 [6] Yael Niv, Nathaniel D Daw, and Peter Dayan, “How fast to work : Response vigor , motivation
938 and tonic dopamine,” in *Neural Information Processing Systems* (2005).
- 939 [7] Jan Drugowitsch, Rubén Moreno-Bote, Anne K Churchland, Michael N Shadlen, and Alexan-
940 dre Pouget, “The Cost of Accumulating Evidence in Perceptual Decision Making,” *The Journal*
941 *of Neuroscience* **32**, 3612 LP – 3628 (2012).
- 942 [8] A Ross Otto and Nathaniel D Daw, “The opportunity cost of time modulates cognitive effort,”
943 *Neuropsychologia* **123**, 92–105 (2019).
- 944 [9] A Ross Otto and Eliana Vassena, “It’s all relative: Reward-induced cognitive control mod-
945 ulation depends on context.” *Journal of Experimental Psychology: General* **150**, 306–313
946 (2021).
- 947 [10] Germain Lefebvre, Aurélien Nioche, Sacha Bourgeois-gironde, and Stefano Palminteri, “Con-
948 trasting temporal difference and opportunity cost reinforcement learning in an empirical
949 money-emergence paradigm,” *Proceedings of the National Academy of Sciences* **115**, E11446
950 LP – E11454 (2018).
- 951 [11] Richard S Sutton and Andrew G Barto, *Reinforcement learning: An introduction, 2nd ed.*,
952 Adaptive computation and machine learning. (The MIT Press, Cambridge, MA, US, 2018)
953 pp. xxii, 526–xxii, 526.
- 954 [12] Khimya Khetarpal, Matthew Riemer, Irina Rish, and Doina Precup, “Towards Continual
955 Reinforcement Learning: A Review and Perspectives,” [arXiv:2012.13490 \[cs.LG\]](https://arxiv.org/abs/2012.13490).
- 956 [13] Roger Ratcliff, “A theory of memory retrieval.” *Psychological Review* **85**, 59–108 (1978).
- 957 [14] Gaurav Malhotra, David S Leslie, Casimir J H Ludwig, and Rafal Bogacz, “Time-varying
958 decision boundaries : insights from optimality analysis,” *Psychon Bull Rev* **25**, 971–996 (2018).
- 959 [15] Jochen Ditterich, “Evidence for time-variant decision making,” *European Journal of Neuro-*
960 *science* **24**, 3628–3641 (2006).

- 961 [16] Paul Cisek, Geneviève Aude Puskas, and Stephany El-Murr, “Decisions in Changing Con-
962 ditions: The Urgency-Gating Model,” *The Journal of Neuroscience* **29**, 11560 LP – 11571
963 (2009).
- 964 [17] Anne K Churchland, Roozbeh Kiani, and Michael N Shadlen, “Decision-making with multiple
965 alternatives,” *Nature Neuroscience* **11**, 693–702 (2008).
- 966 [18] David Thura and Paul Cisek, “Deliberation and Commitment in the Premotor and Primary
967 Motor Cortex during Dynamic Decision Making,” *Neuron* **81**, 1401–1416 (2014).
- 968 [19] David Thura, Ignasi Cos, Jessica Trung, and Paul Cisek, “Context-Dependent Urgency Influ-
969 ences Speed–Accuracy Trade-Offs in Decision-Making and Movement Execution,” *The Journal*
970 *of Neuroscience* **34**, 16442 LP – 16454 (2014).
- 971 [20] David Thura, Jean-François Cabana, Albert Feghaly, and Paul Cisek, “Unified neural dy-
972 namics of decisions and actions in the cerebral cortex and basal ganglia,” *bioRxiv* (2020),
973 [10.1101/2020.10.22.350280](https://doi.org/10.1101/2020.10.22.350280).
- 974 [21] David Thura and Paul Cisek, “The Basal Ganglia Do Not Select Reach Targets but Control
975 the Urgency of Commitment,” *Neuron* **95**, 1160–1170.e5 (2017).
- 976 [22] Peter Janssen and Michael N Shadlen, “A representation of the hazard rate of elapsed time
977 in macaque area LIP,” *Nature Neuroscience* **8**, 234–241 (2005).
- 978 [23] Satohiro Tajima, Jan Drugowitsch, and Alexandre Pouget, “Optimal policy for value-based
979 decision-making,” *Nature Communications* **7**, 12400 (2016).
- 980 [24] Anton Schwartz, “A Reinforcement Learning Method for Maximizing Undiscounted Rewards,”
981 in *International Conference on Machine Learning*, Vol. 0 (1993).
- 982 [25] Yael Niv, Nathaniel D Daw, Daphna Joel, and Peter Dayan, “Tonic dopamine: opportunity
983 costs and the control of response vigor,” *Psychopharmacology* **191**, 507–520 (2007).
- 984 [26] Sara M Constantino and Nathaniel D Daw, “Learning the opportunity cost of time in a patch-
985 foraging task,” *Cogn Affect Behav Neurosci.* **15**, 837 (2015).
- 986 [27] Benjamin Y Hayden and Yael Niv, “The case against economic values in the orbitofrontal
987 cortex (or anywhere else in the brain),” *PsyArXiv* [10.31234/osf.io/7hgup](https://doi.org/10.31234/osf.io/7hgup).
- 988 [28] Nathaniel D Daw, “Advanced Reinforcement Learning,” in *Neuroeconomics*, edited by Paul W
989 Glimcher and Ernst B T Neuroeconomics (Second Edition) Fehr (Academic Press, San Diego,
990 2014) 2nd ed., Chap. 16, pp. 299–320.
- 991 [29] These were primarily as a result of the experimenter adapting to fluctuations in motivation
992 of the subject. D. Thura. Personal communication.
- 993 [30] Edward Vul, Noah Goodman, Thomas L Griffiths, and Joshua B Tenenbaum, “One and
994 Done? Optimal Decisions From Very Few Samples,” *Cognitive Science* **38**, 599–637 (2014).
- 995 [31] Single subject behavioural data shared by Thomas Thierry.
- 996 [32] Surya Ganguli, James W Bisley, Jamie D Roitman, Michael N Shadlen, Michael E Goldberg,
997 and Kenneth D Miller, “One-Dimensional Dynamics of Attention and Decision Making in
998 LIP,” *Neuron* **58**, 15–25 (2008).
- 999 [33] David Thura and Paul Cisek, “Modulation of Premotor and Primary Motor Cortical Activity
1000 during Volitional Adjustments of Speed-Accuracy Trade-Offs,” *The Journal of Neuroscience*
1001 **36**, 938 – 956 (2016).
- 1002 [34] Kiyohito Iigaya, Yashar Ahmadian, Leo P Sugrue, Greg S Corrado, Yonatan Loewenstein,
1003 William T Newsome, and Stefano Fusi, “Deviation from the matching law reflects an optimal
1004 strategy involving learning over multiple timescales,” *Nature Communications* **10**, 1466 (2019).
- 1005 [35] Long Ding and Joshua I. Gold, “The Basal Ganglia’s Contributions to Perceptual Decision
1006 Making,” *Neuron* **79**, 640–649 (2013).

- 1007 [36] Kong-Fatt Wong and Xiao-Jing Wang, “A Recurrent Network Mechanism of Time Integration
1008 in Perceptual Decisions,” *The Journal of Neuroscience* **26**, 1314 – 1328 (2006).
- 1009 [37] Alex Roxin and Anders Ledberg, “Neurobiological Models of Two-Choice Decision Making
1010 Can Be Reduced to a One-Dimensional Nonlinear Diffusion Equation,” *PLOS Computational
1011 Biology* **4**, e1000046 (2008).
- 1012 [38] David Meder, Nils Kolling, Lennart Verhagen, Marco K Wittmann, Jacqueline Scholl, Kristof-
1013 fer H Madsen, Oliver J Hulme, Timothy E J Behrens, and Matthew F S Rushworth, “Simul-
1014 taneous representation of a spectrum of dynamically changing value estimates during decision
1015 making,” *Nature Communications* **8** (2017), 10.1038/s41467-017-02169-w.
- 1016 [39] “Predictive Representations in Hippocampal and Prefrontal Hierarchies,” .
- 1017 [40] Jan Zimmermann, Paul W Glimcher, and Kenway Louie, “Multiple timescales of normalized
1018 value coding underlie adaptive choice behavior,” *Nature Communications* **9**, 3206 (2018).
- 1019 [41] HyungGoo R Kim, Athar N Malik, John G Mikhael, Pol Bech, Iku Tsutsui-Kimura, Fangmiao
1020 Sun, Yajun Zhang, Yulong Li, Mitsuko Watabe-Uchida, Samuel J Gershman, and Naoshige
1021 Uchida, “A Unified Framework for Dopamine Signals across Timescales,” *Cell* **183**, 1600–
1022 1616.e25 (2020).
- 1023 [42] Paul Masset, Athar N. Malik, HyungGoo R. Kim, Pol Bech, and Naoshige Uchida, “A
1024 diversity of discounting horizons explains ramping diversity in dopaminergic neurons,” in
1025 *COSYNE Abstracts* (2021).
- 1026 [43] Angela J Langdon and Nathaniel D Daw, “Beyond the Average View of Dopamine,” *Trends
1027 in Cognitive Sciences* **24**, 499–501 (2020).
- 1028 [44] John G Mikhael and Samuel J Gershman, “Adapting the flow of time with dopamine,” *Journal
1029 of Neurophysiology* **121**, 1748–1760 (2019).
- 1030 [45] Ido Toren, Kristoffer C Aberg, and Rony Paz, “Prediction errors bidirectionally bias time
1031 perception,” *Nature Neuroscience* **23**, 1198–1202 (2020).
- 1032 [46] Lars Hunger, X Arvind Kumar, and X Robert Schmidt, “Abundance Compensates Kinet-
1033 ics : Similar Effect of Dopamine Signals on D1 and D2 Receptor Populations,” *Journal of
1034 Neuroscience* **40**, 2868–2881.
- 1035 [47] Julia Cox and Ilana B Witten, “Striatal circuits for reward learning and decision-making,”
1036 *Nature Reviews Neuroscience* **20**, 482–494 (2019).
- 1037 [48] Helen N Schwerdt, Hideki Shimazu, Ken-ichi Amemori, Satoko Amemori, Patrick L Tierney,
1038 Daniel J Gibson, Simon Hong, Tomoko Yoshida, Robert Langer, Michael J Cima, and Ann M
1039 Graybiel, “Long-term dopamine neurochemical monitoring in primates,” *Proceedings of the
1040 National Academy of Sciences* **114**, 13260 LP – 13265 (2017).
- 1041 [49] Tommaso Patriarchi, Jounhong Ryan Cho, Katharina Merten, Mark W Howe, Aaron Marley,
1042 Wei-hong Xiong, Robert W Folk, Gerard Joey Broussard, Ruqiang Liang, Min Jee Jang,
1043 Haining Zhong, Daniel Dombeck, Mark Von Zastrow, Axel Nimmerjahn, Viviana Gradinaru,
1044 John T Williams, and Lin Tian, “Ultrafast neuronal imaging of dopamine dynamics with
1045 designed genetically encoded sensors,” *Science* **4422** (2018), 10.1126/science.aat4422.
- 1046 [50] Matthew A Carland, David Thura, and Paul Cisek, “The Urge to Decide and Act: Implica-
1047 tions for Brain Function and Dysfunction,” *The Neuroscientist* **25**, 491–511 (2019).
- 1048 [51] Samuel J Gershman and Naoshige Uchida, “Believing in dopamine,” *Nature Reviews Neuro-
1049 science* **20**, 703–714 (2019).
- 1050 [52] Andrew Westbrook and Todd S Braver, “Dopamine Does Double Duty in Motivating Cognitive
1051 Effort,” *Neuron* **91**, 708 (2016).

- 1052 [53] Pablo Tano, Peter Dayan, and Alexandre Pouget, “A Local Temporal Difference Code for
1053 Distributional Reinforcement Learning,” in *Advances in Neural Information Processing Sys-*
1054 *tems*, Vol. 33, edited by H Larochelle, M Ranzato, R Hadsell, M F Balcan, and H Lin (Curran
1055 Associates, Inc., 2020) pp. 13662–13673.
- 1056 [54] William Fedus, Carles Gelada, Yoshua Bengio, Marc G Bellemare, and Hugo Larochelle,
1057 “Hyperbolic Discounting and Learning over Multiple Horizons,” [arXiv:1902.06865 \[stat.ML\]](https://arxiv.org/abs/1902.06865).
- 1058 [55] I Momennejad, E M Russek, J H Cheong, M M Botvinick, N D Daw, and S J Gershman,
1059 “The successor representation in human reinforcement learning,” *Nature Human Behaviour*
1060 **1**, 680–692 (2017).
- 1061 [56] Personal communication, Thomas Thierry.
- 1062 [57] Stefano Palminteri and Maël Lebreton, “Context-dependent outcome encoding in human re-
1063 inforcement learning,” *Current Opinion in Behavioral Sciences* **41**, 144–151 (2021).
- 1064 [58] Ernest S Davis and Gary F Marcus, “Computational limits don’t fully explain human cognitive
1065 limitations,” *Behavioral and Brain Sciences* **43**, e7 (2020).
- 1066 [59] Arman Abrahamyan, Laura Luz Silva, Steven C Dakin, Matteo Carandini, and Justin L Gard-
1067 ner, “Adaptable history biases in human perceptual decisions,” *Proceedings of the National*
1068 *Academy of Sciences* **113**, E3548 LP – E3557 (2016).
- 1069 [60] Zhijian Wang, Bin Xu, and Hai-Jun Zhou, “Social cycling and conditional responses in the
1070 Rock-Paper-Scissors game,” *Scientific Reports* **4**, 5830 (2014).
- 1071 [61] A. Churchland. Personal communication.

1072

Supplemental Materials

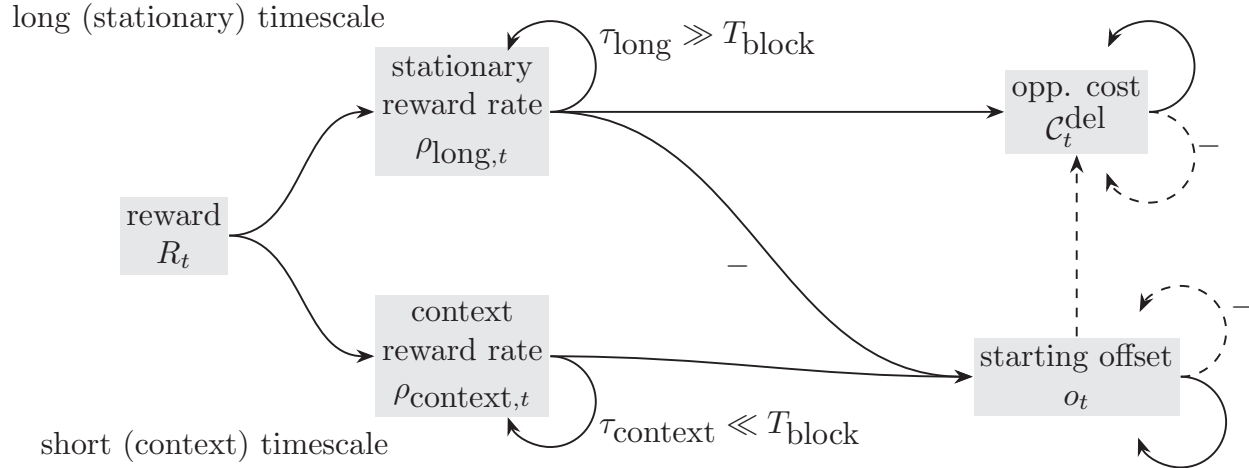


Figure S1. *Reward filtering scheme for online computation of within-trial opportunity cost.* With t denoting absolute time, the reward sequence, R_t , is integrated on both a stationary (τ_{long}) and context (τ_{context}) filtering timescale to produce estimates of the stationary and context-specific reward rates, respectively. These are large and small, respectively, relative to the average context switching timescale, T_{block} . The estimate of the context-specific offset, o_t is computed by time-integrating the difference of these two estimates. In this filtering, when a trial terminates, the effective operation is that C_t^{del} is set to o_t , and the latter is zeroed. Thus, the opportunity cost starts at this offset and then integrates ρ_{long} , $C_{t,k}^{\text{del}} = o_{T_{k-1},k-1} + \rho_{\text{long},k-1}t$, where $o_{T_{k-1},k-1} = (\rho_{\text{context},k-1} - \rho_{\text{long},k-1})T_{k-1}$. Notes on the computational graph: Arrows pass the value at each time step (dashed arrows only pass the value when a trial terminates). Links annotated with ‘-’ multiply the passed quantity by -1 .

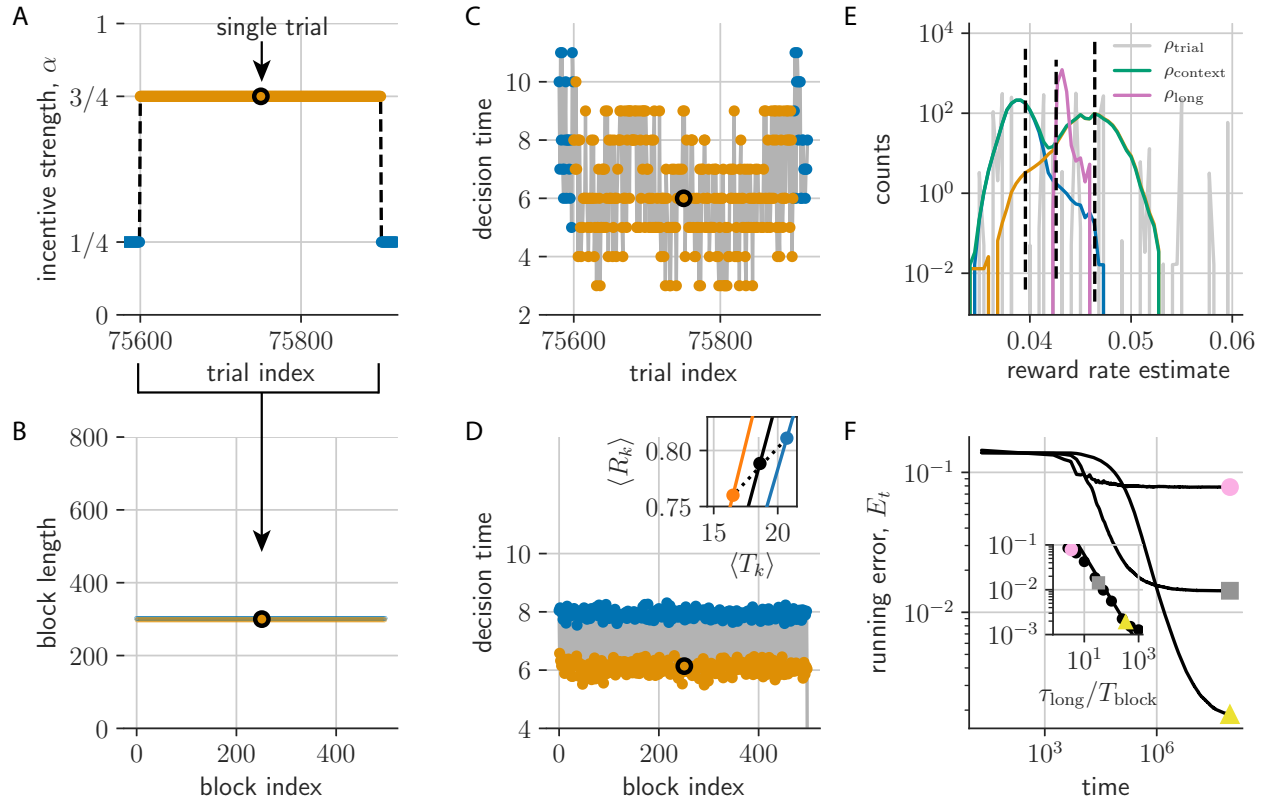


Figure S2. *PGD agent plays the tokens task with periodic α -dynamics.* (a) Trials are grouped into alternating trial blocks of constant α (fast (orange) and slow (blue) conditions). (b) Here, trial block durations are constant over the experiment. (c) Decision times over the trials from (a) distribute widely, but relax after context switches. (d) Block-averaged decision times remain stationary. Inset shows the context-conditioned trial-averaged reward $\langle R_k \rangle$ and trial duration $\langle T_k \rangle$ (orange and blue dots; black is unconditioned average; $\langle \cdot \rangle$ denotes the trial ensemble average). Lines pass through the origin (slope given by the respective reward rate). (e) Distribution of estimates have lower variance than the trial reward rates, ρ^{trial} (gray). The conditioned averages of $\hat{\rho}_k^{\text{context}}$ shown as blue and orange. (f) The relative error in estimating ρ , $E_t = \frac{1}{t} \sum_k^t |\hat{\rho}_k^{\tau_{\text{long}}} - \rho| / \rho$, for $\tau_{\text{long}} = 10^3$ (circle), 10^4 (square), 10^5 (triangle). Inset shows that $E_{T_{\text{exp}}} \propto (\tau_{\text{long}}/T_{\text{block}})^{-1}$ over a grid of τ_{long} and T_{block} as expected (black line).

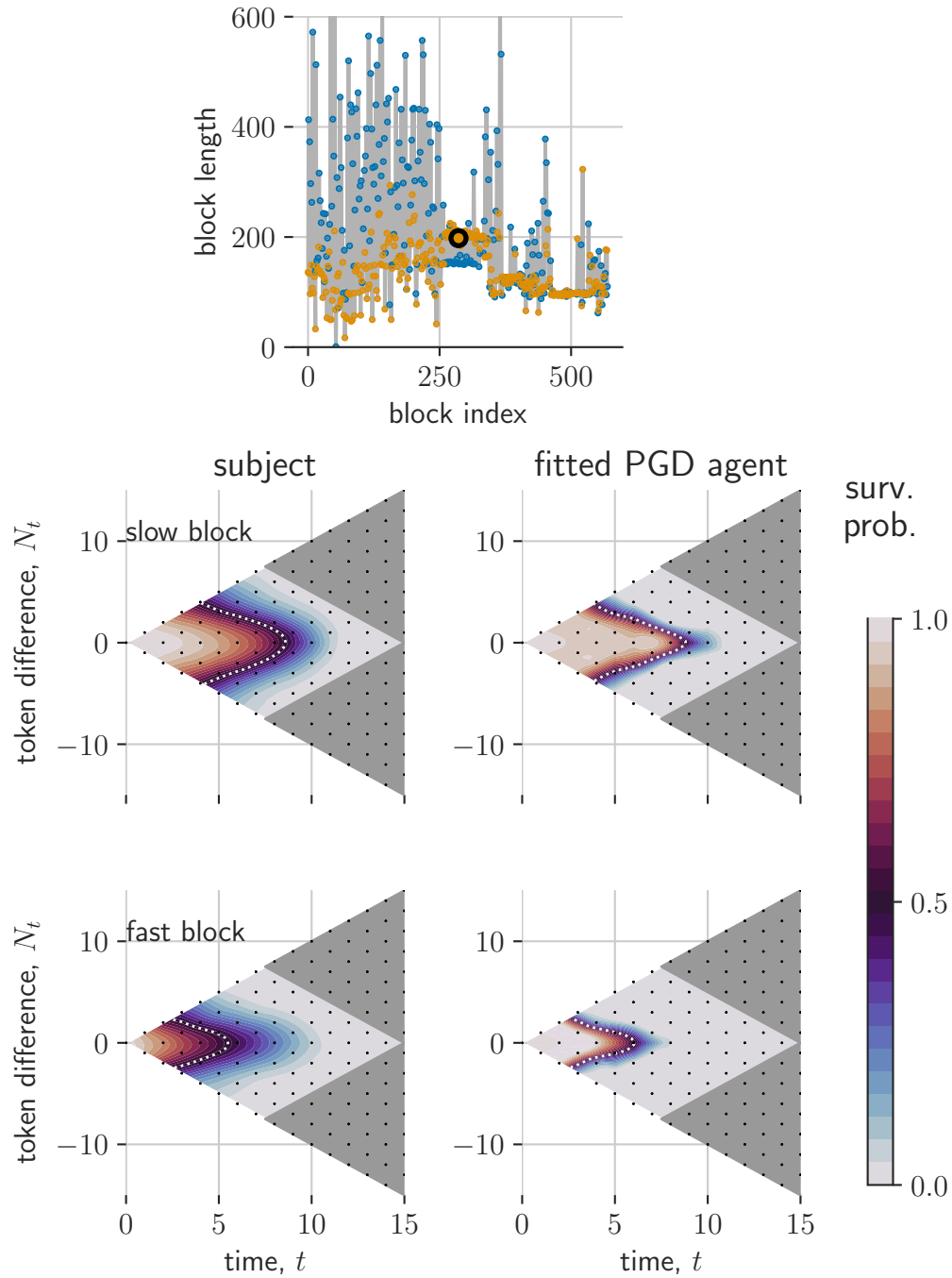


Figure S3. Comparison of PGD and NHP in non-stationary α dynamics from [19]: Subject 2. Same as fig. 5.

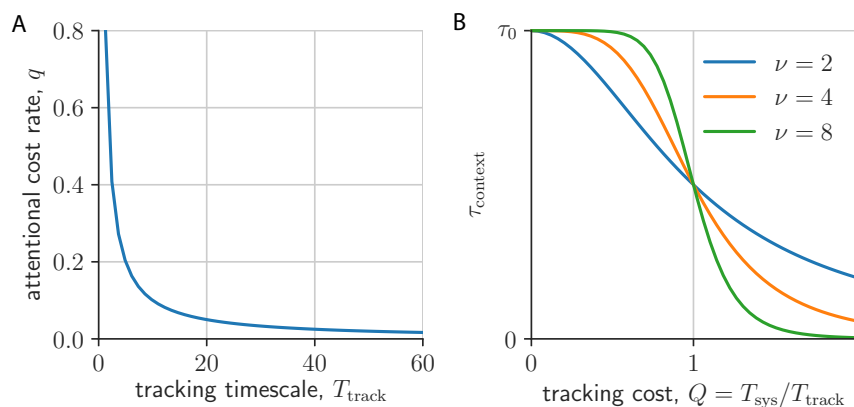


Figure S4. *Asymmetric switching cost model*. (a) Attentional cost rate, q , is set to be inversely proportional to tracking timescale, T_{track} . (b) Filtering timescale τ_{context} is scaled down with tracking cost, $Q = T_{\text{sys}}/T_{\text{track}}$ from a base timescale, here denoted τ_0 (shown for three values of sensitivity $\nu = 2, 4, 8$).

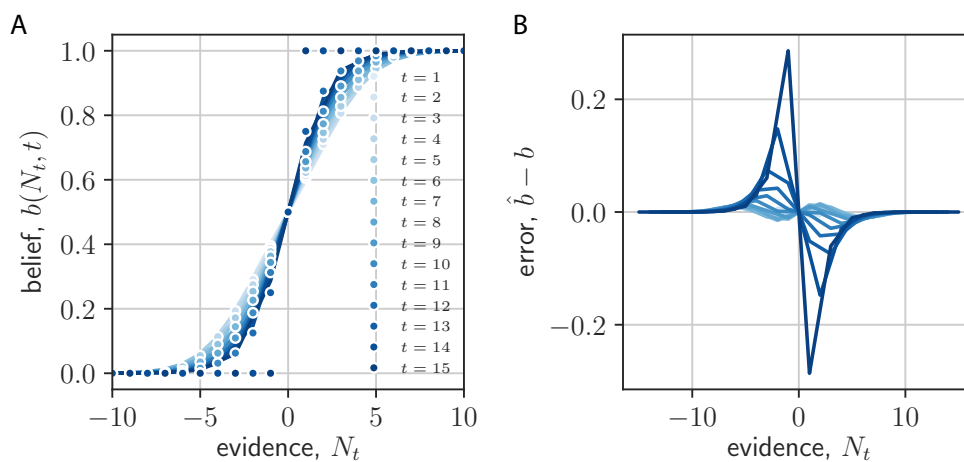


Figure S5. *sigmoidal approximation to expected reward*. (a) the approximation explained in [Methods: State-conditioned expected trial reward](#), for different dec,[p]ision times. (b) The error in the approximation for different decision times.

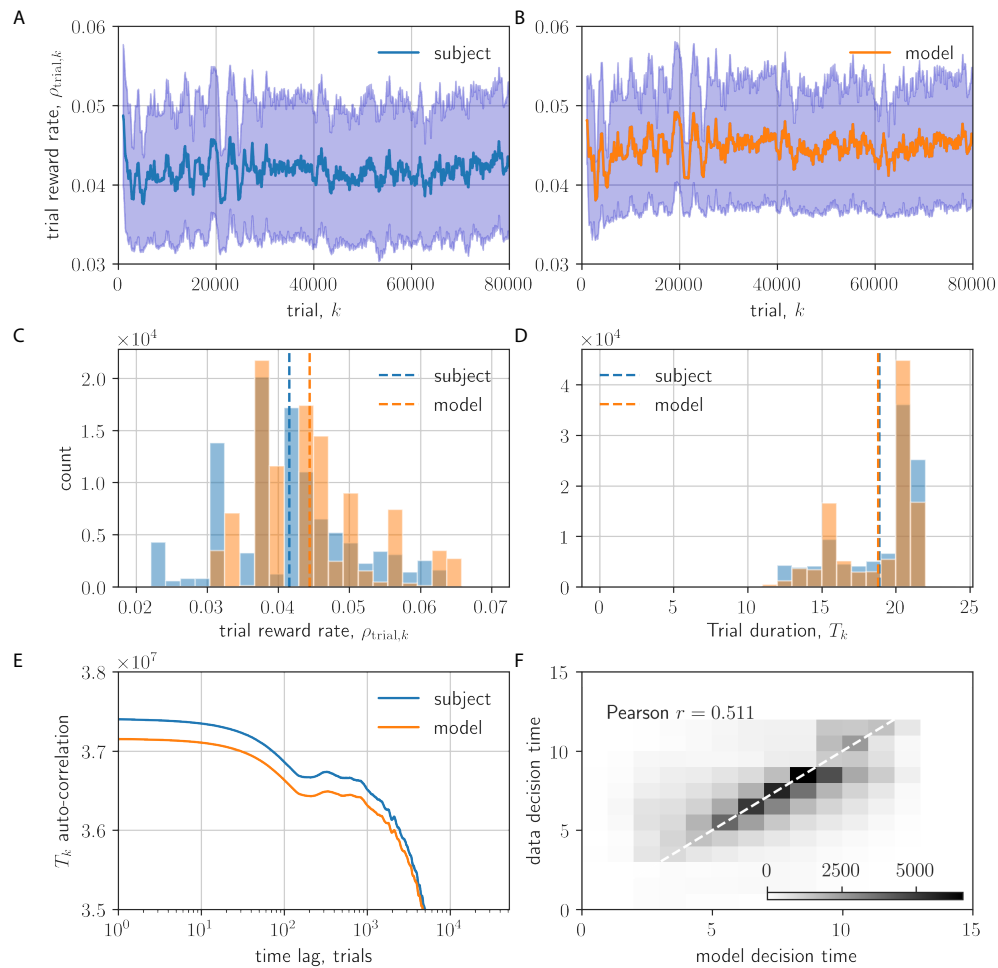


Figure S6. *Model validation on behavioural statistics from [19].* (a,b) Running average (last 1000 trial) of trial reward rate ρ_k^{trial} . (c,d) Histograms of trial reward rate, ρ_k^{trial} (c) and trial duration, T_k (d). (e) Auto-correlation function of trial duration. (f) Data vs. model decision time (gray-scale is count; white dashed line is perfect correlation; actual Pearson correlation is shown)

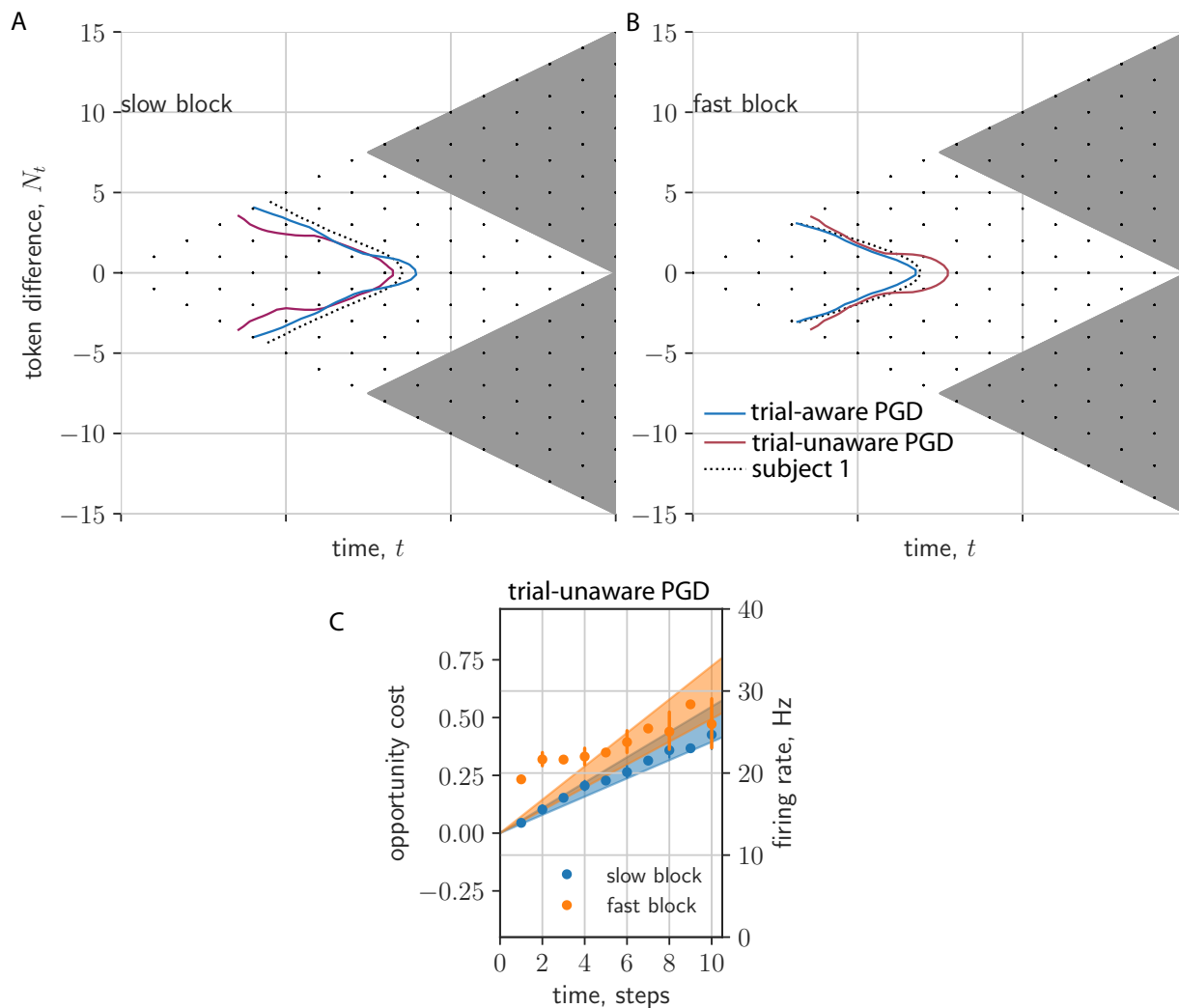


Figure S7. *Comparison of trial-aware and trial-unaware results.* (a,b) 1/2-Survival probability contours for subject 1 (dashed), trial-aware PGD (blue), and trial-unaware PGD (red) for slow (a) and fast (b) context-conditioned data. (c) Opportunity cost for trial-unaware PGD (compare with [fig. 2b](#)). Opportunity cost range adjusted here such that data within standard error of trial-unaware PGD model prediction for slow block (blue).

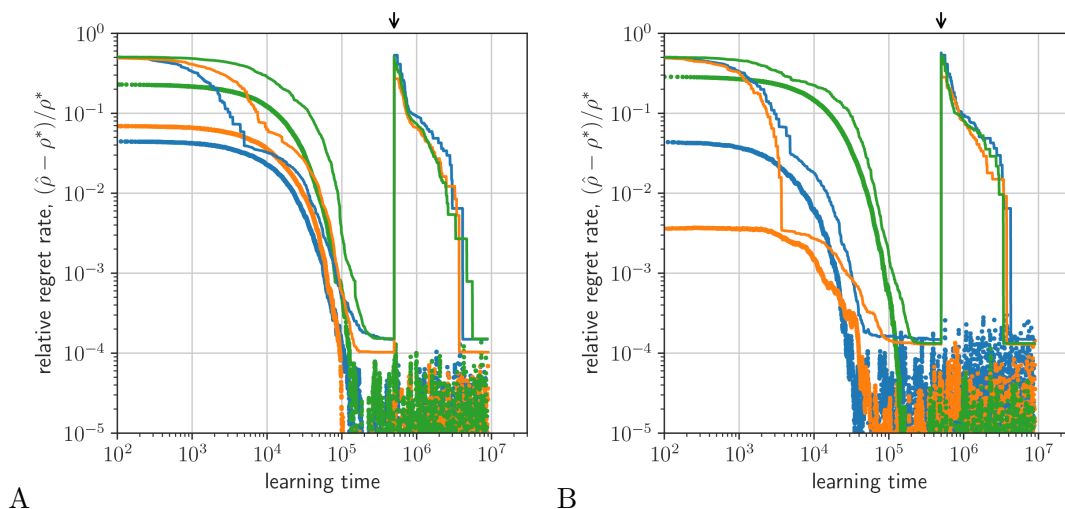


Figure S8. *Comparison of PGD and AR-RL learning on a patch leaving task.* Performance is defined as relative regret rate, $(\hat{\rho} - \rho^*)/\rho^*$ (PGD (dots); AR-RL (lines)). (a) Performance over different sizes of the state vector ($d = 100$ (blue), 200 (orange), 300 (green)). (b) Performance over different learning rates (parametrized by integration time constant, $\tau = 1 \times 10^4$ (blue), 2×10^4 (orange), 3×10^4 (green)). (parameters: $\lambda = 1/5$; r_{\max} sampled uniformly on $[0, 1]$). A random state label permutation is made at the time indicated by the black arrow. Values were initialized at -1 .

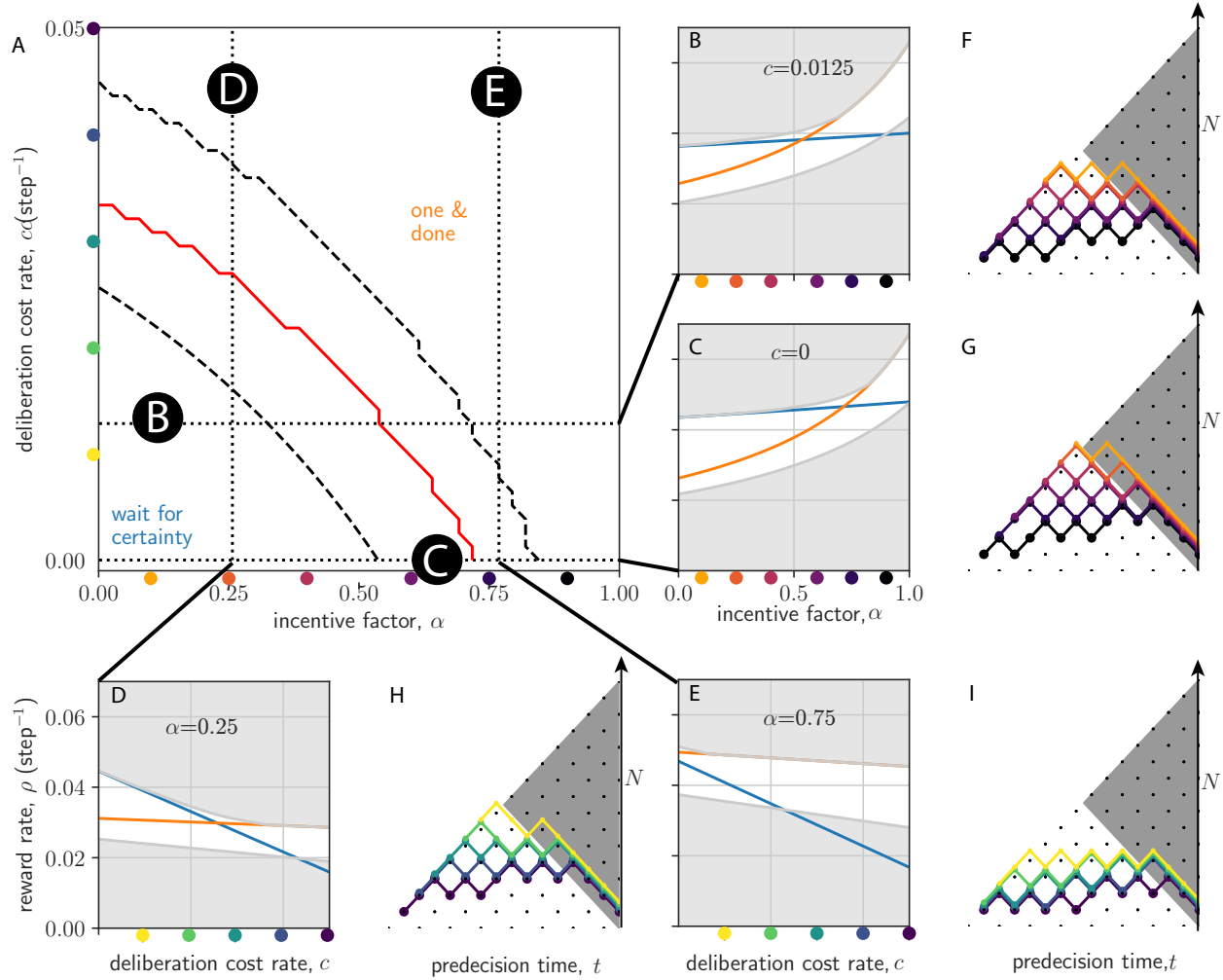


Figure S9. *Reward rate optimal strategies in (α, c) plane.* (a) The reward-rate maximizing policy interpolates from the wait-for-certainty strategy at weak incentive (low α) and low deliberation cost (low c), to the one-and-done strategy at strong incentive (high α) and high deliberation cost (high c). Dashed lines bound a transition regime between the two extreme strategies. Red line denotes where they have equal performance. (b-e) Slices of the (α, c) -plane. Shown are the reward rate as a function of α (b,c) and c (d,e) (wait-for-certainty strategy is shown in blue; one-and-done strategy is shown in orange). N is the magnitude of the token difference

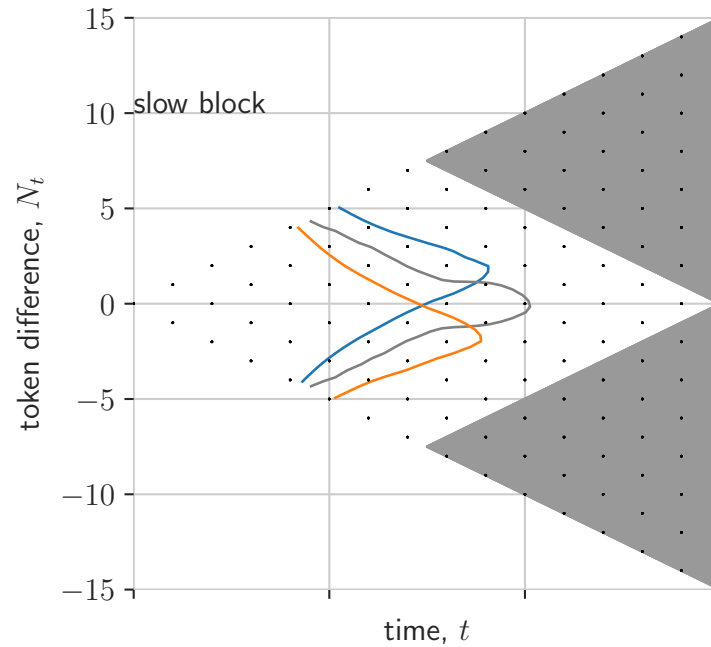


Figure S10. *Asymmetric action rewards skew survival probability.* Here, we plot the half-maximum of the PGD survival probability for three values of the action reward bias, $\gamma = -0.6, 0, 0.6$ (blue, black and orange, respectively). Other model parameters same as in fitted model.

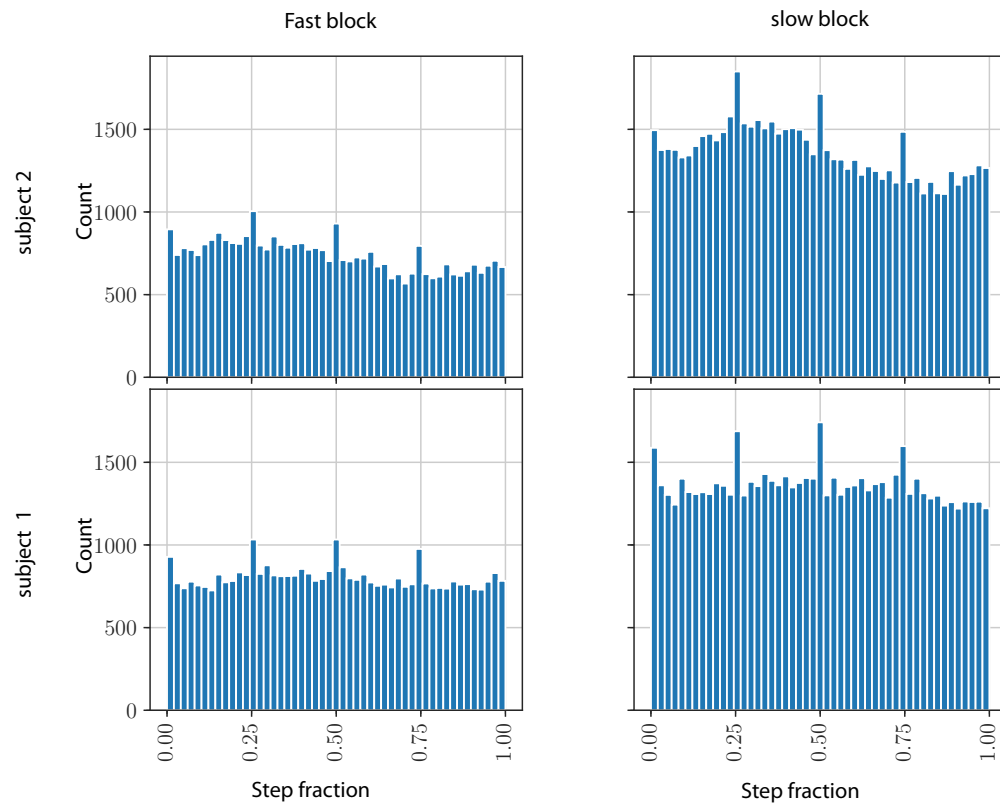


Figure S11. *Decision times relative to token jumps.* Here, we plot the histograms of decision times using their position between token jumps, the step fraction. The data is separated by α and monkey.