

1 **TITLE**

2 Application of a unifying reward-prediction error (RPE)-based framework to explain underlying
3 dynamic dopaminergic activity in timing tasks

4

5 **AUTHORS**

6 Allison E. Hamilos¹ & John A. Assad^{1,2}

7

8 **AFFILIATIONS**

9 ¹Department of Neurobiology, Harvard Medical School, Boston, Massachusetts, 02115, USA.

10 ²Istituto Italiano di Tecnologia, Genova, Italy.

11 Correspondence: A.H. (ahamilos@mit.edu) or J.A.A. (jassad@hms.harvard.edu).

12

13 **SUMMARY**

14 **This manuscript is intended as a theoretical companion to [Hamilos et al., 2020](#)¹, in which we**
15 **examined the role of dopaminergic neurons (DANs) in self-timed movements. In that study,**
16 **we recorded DAN signals in mice trained to initiate a licking movement after a self-timed**
17 **delay following a start-timing cue. DAN signals both before the start-timing cue and during**
18 **the timing interval predicted the timing of movement onset, up to seconds before the**
19 **movement itself. In particular, dopaminergic signals “ramped up” from the time of the cue**
20 **to the time of the movement. On a given trial, the slope of the ramping was predictive of**
21 **when the movement would occur, with steep slope associated with early movement and**
22 **shallow slope with late movement, reminiscent of a ramp-to-threshold process.**

23

24 **Ramping dopaminergic signals were recently proposed in a theoretical framework that**
25 **examined temporal-difference learning under resolved state uncertainty (Mikhael *et al.*,**
26 **2019²; Mikhael & Gershman, 2019³; Gershman, 2014⁴). Here, we show that an adapted**
27 **version of Mikhael *et al.*'s model recapitulates the ramping dopaminergic signaling observed**
28 **in our self-timed movement task. We also applied the model to results reported in a recent**
29 **temporal bisection study, in which mice categorized time intervals as relatively short or long**
30 **compared to a criterion interval (Soares *et al.*, 2016⁵). The model successfully predicted the**
31 **relative amplitude of dynamic DAN signals observed in the bisection task. These combined**
32 **results suggest a common neural mechanism that broadly underlies timing behavior: trial-**
33 **by-trial variation in the rate of the internal “pacemaker,” manifested in DAN signals that**
34 **reflect stretching or compression of the derivative of the subjective value function relative to**
35 **veridical time. In this view, faster pacemaking is associated with relatively high amplitude**
36 **dopaminergic signaling, whereas slower pacemaking is associated with relatively low levels**
37 **of dopaminergic signaling.**

38

39

40 **MAIN TEXT**

41 *Nigrostriatal dopaminergic signaling controls the moment-to-moment decision of when to*

42 *move*

43 Clues from human movement disorders and pharmacological studies have long suggested a
44 connection between the neurotransmitter dopamine and the timing of movement initiation^{3,5-13}. We
45 recently showed that dopaminergic signaling controls the moment-to-moment timing of
46 movements in mice¹. We recorded dopaminergic signals with fiber photometry in mice executing

47 a self-timed movement task, in which animals received juice rewards for withholding movement
48 for a proscribed interval (3.3 s) after a start-timing cue and then initiating movement (a first-lick)
49 within a rewarded time window (3.3-7 s, [Figure 1a](#)). We observed two aspects of dopaminergic
50 signaling that predicted movement timing: 1) pre-trial baseline signaling of nigrostriatal dopamine
51 neurons (DANs), and, 2) slow “ramping” signals that built up over the course of seconds between
52 the start-timing cue and the self-timed movement. Although self-timed movements occurred with
53 variable timing relative to the start-timing cue¹, DAN signaling rose to about the same level at the
54 moment of movement onset, reminiscent of a ramp-to-threshold process ([Figure 1b](#)). DAN signals
55 were not explained by ongoing nuisance movements nor optical artifacts and were best modeled
56 with timing-dependent predictors, including a baseline offset term whose amplitude was
57 proportional to the mouse’s timing on the upcoming trial, as well as a stretch feature that encoded
58 percentages of elapsed time between the cue and self-timed movement¹. DAN ramping activity
59 predicted first-lick time on single trials, independently of trial history, and optogenetic
60 manipulation of DANs bidirectionally shifted movement timing, with activation early-shifting
61 movements versus inhibition late-shifting movements. Together, these results indicate that
62 dopaminergic signaling during self-timing controls the moment of movement onset.
63

64 **Figure 1 | Nigrostriatal**
65 **dopaminergic signaling during a**
66 **self-timed movement task.**

67 Figure adapted from [Hamilos *et al.*, 2020¹](#) and used with permission.

68 **a**, Schematic of self-timed movement
69 task.

70
71 **b**, Top: Average DAN GCaMP6f
72 responses from 12 mice; Bottom:
73 Responses of tdTomato, a non-activity-
74 dependent fluorophore used to control
75 for optical artifacts. The different
76 colored traces correspond to averaged
77 trial responses with different first-lick
78 times (ranging from 1-4 s in 250 ms
79 increments). Traces are plotted up to
80 150 ms before first-lick. Averaged
81 traces are aligned relative to both start-
82 timing cue onset (left of x-axis break)

83 and first-lick (right of x-axis break); the break in the x-axis indicates the change in plot alignment.
84 **c**, Cue-aligned average DAN GCaMP6f signals at lower gain show post-movement RPE-like
85 signals. Movement onset occurs just before the peak response for each curve. Mice were rewarded
86 for first-licks made later than 3.3 s, but were not rewarded for earlier first-licks.

87

88

89 ***A temporal-difference learning model of dynamic dopaminergic signaling***

90 We were interested in understanding the origin of the dynamic dopaminergic signals we observed

91 in our self-timed movement task and how they fit into the context of prior work on the dopamine

92 system. A framework that has explained many disparate experimental results from the

93 dopaminergic system is temporal difference (TD) learning with reward-prediction errors (RPE)^{2,14}.

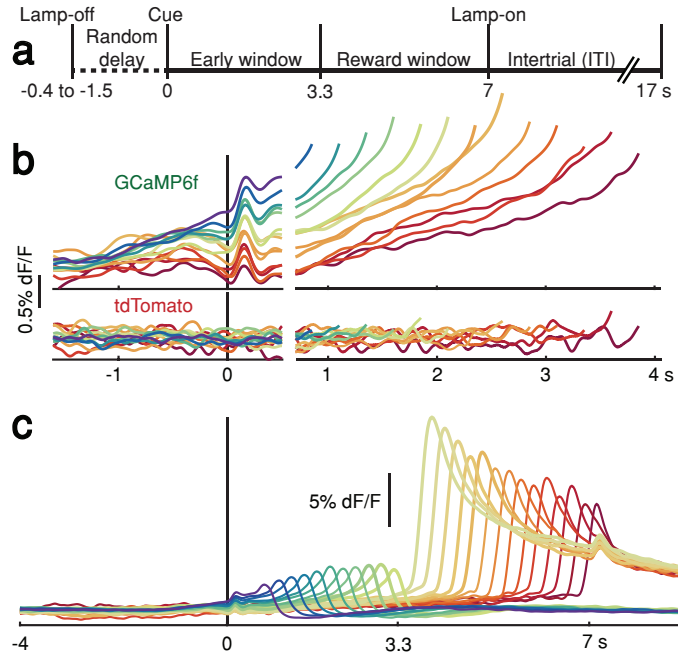
94 In this framework, DAN activity is thought to reflect the moment-to-moment difference in the

95 animal's expectation versus its perception of the value of its current state, where value is defined

96 as the temporally-discounted expectation of total future reward. In classical trace-conditioning

97 paradigms, DANs fire in transient bursts to unexpected rewards and reward-predicting cues,

98 whereas they pause their firing when expected reward is omitted. Indeed, we observed RPE-like



99 signals in the cue-related transient, dips in activity after unrewarded first-licks, and surges in
100 activity following rewarded first-licks (Figure 1c). Persistence of RPE-like signals in well-trained
101 animals has been suggested to arise from the inherent imprecision in neural timing¹⁰, which may
102 reflect the animal's moment-to-moment uncertainty of its current state—i.e., its position in time—
103 and, by extension to our task, uncertainty about its accuracy for a given self-timed lick³. Indeed,
104 positive-going RPE-like signals were strongest for first-licks closest to the reward-boundary (3.3
105 s), presumably when the mouse's "confidence" of reward was lowest, consistent with the greatest
106 RPE occurring when the mice were least certain of success (Figure 1c).

107
108 Whereas RPE-frameworks have explained *transient* bursts and pauses in DAN activity during
109 trace conditioning and other types of learning experiments, DAN activity can also change more
110 slowly^{2,3}. For example, "ramping" signals build up over seconds during goal-directed navigation¹⁵,
111 bandit tasks in which animals must complete multiple goals to receive reward^{16,17}, and tasks with
112 visual cues of proximity to reward¹⁸. It has been suggested that DANs could signal different
113 information via slow changes in activity (e.g., motivation, ongoing value, vigor) compared to fast-
114 timescale activity (e.g., post-hoc RPE signals for learning), and a number of proposals have
115 suggested that DANs multiplex different kinds of information over different timescales and
116 contexts^{17,19}.

117
118 However, recent models have proposed RPE-based explanations that may be able to reconcile
119 these seemingly disparate dopamine signals^{2,3,18}. While these models do not refute the possibility
120 that DANs could encode other types of information (e.g., value, vigor, etc.), they are attractive for
121 their parsimonious explanation of how fast time-scale phenomena and slowly-evolving ramps

122 could arise from the same underlying RPE-based calculation. In short, these models employ
123 principles from TD learning to show how certain shapes of the value function (i.e., the assignment
124 of values to the series of behavioral states comprising a task) can give rise to a *continuously*
125 *changing* RPE, even in well-trained animals^{2,3,18,20}.

126

127 We were interested in whether an RPE-based framework could explain the results found in our
128 self-timed movement task as well as results from other timing tasks⁵. To approach this question,
129 we applied a key feature of TD learning algorithms to determine what an RPE-like signal would
130 look like in different kinds of timing tasks. Specifically, we took advantage of the fact that *RPE is*
131 *proportional to the derivative of the subjective value function under conditions of state*
132 *uncertainty*^{2,3}, as is the case during timing tasks in which the animal must rely on its own internal
133 representation of time to guide behavior².

134

135 Thus, if the value landscape for a given behavioral task is known, and if DAN activity encodes
136 RPE, the RPE-based framework makes predictions about the expected shape of dynamic DAN
137 activity during the task. In a recent study, similar applications of this principle predicted the
138 ramping DAN signals that were observed in virtual reality (VR) tasks in which animals were
139 moved passively through VR spaces, as well as when the animals passively viewed abstract,
140 dynamic visual cues indicating proximity to reward¹⁸, suggesting the ramping in our task could be
141 explained from similar principles.

142

143

144

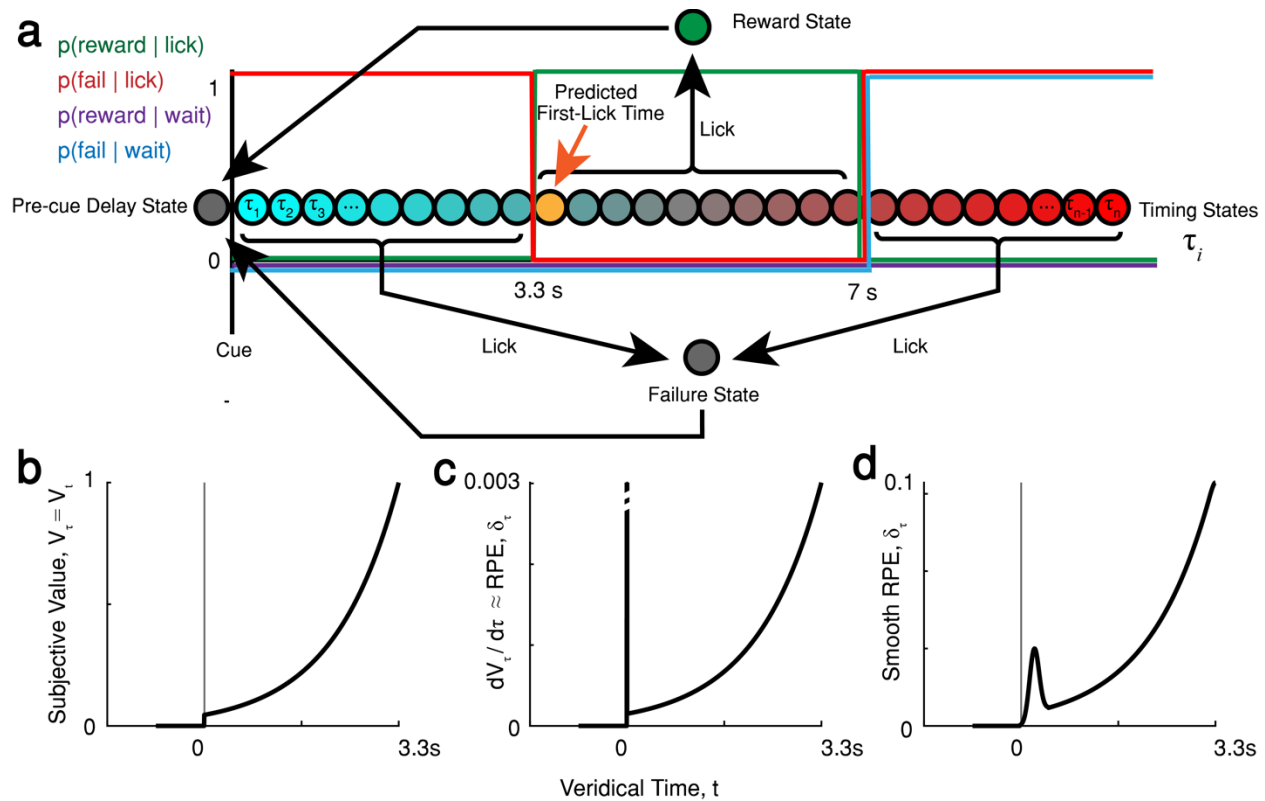
145 ***RPE-predictions for DAN responses during self-timed movement***

146 In a simple TD learning model of self-timed movement, time may be modeled as a continuous set
147 of states through which a Markov agent must traverse to receive reward²¹ (Figure 2a). At each state
148 transition (timestep), the agent must decide whether to move (lick) or to wait based on the
149 probability of transitioning to a reward or failure state. If the agent is an optimal timer, its
150 subjective approximation of its current state, τ , accurately tracks veridical time, t , and it will thus
151 withhold movement until the first moment at which reward will be available in response to licking
152 (3.3 s in our experiment).

153

154 The value landscape of this model can be understood intuitively. When the cue event occurs, a
155 well-trained agent can expect an increased possibility of reward in the next few seconds; thus, at
156 this moment, value increases. However, reward never occurs within the first 3.3 s of the standard
157 timing task we implemented; thus, value at the cue is necessarily lower than value at 3.3 s. In fact,
158 value will constantly increase as time approaches 3.3 s. Thus, as long as the agent withholds licks,
159 the value landscape, V_t , during the first few seconds is a monotonically increasing, convex
160 function⁴ (Figure 2b). If the agent is an optimal timer, the subjective approximation of the value
161 function, \hat{V}_τ , matches the true value function, and $\hat{V}_\tau = V_t$.

162



163

164 **Figure 2 | Value and RPE Landscapes for an optimal timer predict DAN responses**
 165 **during the self-timed movement task.** **a**, State space and probability of state transition for
 166 an optimal timer. Gold-shaded state is the first state from which reward is available, and thus
 167 is when the first-lick is predicted to occur. **b**, Estimated value function \hat{V}_t , where $\hat{V}_t \approx V_t$ for
 168 an optimal timer. An exponential value landscape is shown, consistent with prior literature².
 169 However, any sufficiently convex function could be implemented with the same result²⁻⁴. The
 170 agent is expected to first-lick at the peak of the trajectory. **c**, RPE function for an optimal
 171 timer, estimated as $\delta_\tau \approx \hat{V}'_\tau$, the derivative of the subjective value function. Y-axis scaled
 172 to show ramp. **d**, Predicted DAN GCaMP6f signals for an optimal timer. The RPE function
 173 was smoothed with a gaussian kernel spanning *ca.* 10% of the interval to approximate
 174 GCaMP6f off-dynamics.

175

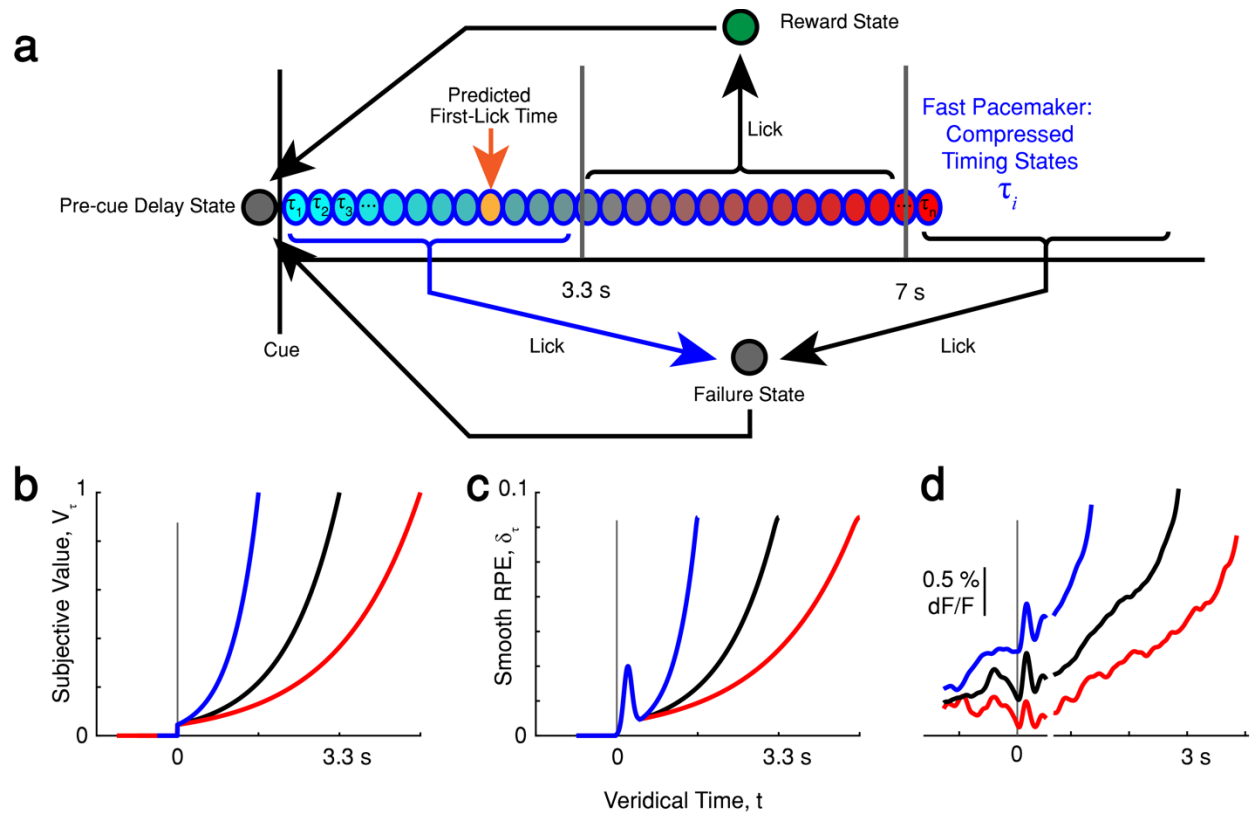
176 However, we assume that, because the timer does not have access to the true state identity, t , it is
 177 never certain of its subjective approximation of its state, τ . Under conditions of state uncertainty,
 178 RPE is approximately the derivative of the subjective value function^{2,18}, $\delta_\tau \approx \hat{V}'_\tau$, where δ_τ is
 179 RPE at subjective time τ , and \hat{V}'_τ is the time-derivative of the subjective value function. Thus, the
 180 shape of the RPE function, δ_τ is also quite simple: a transient increase at the cue followed by a

181 slowly-evolving ramp (Figure 2c). If the RPE function is measured by a calcium indicator such as
182 GCaMP6f, the binding kinetics of the indicator would tend to blur the RPE function, which we
183 approximated by smoothing (Figure 2d).

184

185 The modeled RPE function mirrors the shape of the dynamics observed in DAN signals: a cue-
186 related transient followed by a slow ramp up to the time of first-lick. However, unlike the optimal
187 timer in this model, mice, like humans, exhibit suboptimal timing behavior with variability
188 proportional to the duration of the timed interval¹⁰. It has been proposed that this variability in
189 timing results from imprecision in an internal clock, referred to classically as the internal
190 “pacemaker²²”. When the pacemaker is fast, self-timed movements occur relatively early, whereas
191 when the pacemaker is slow, later movements occur. These changes in the pacemaker rate would
192 correspond to the mouse traversing the set of subjective states, τ , at different rates than the passage
193 of veridical time, t (Figure 3a), resulting in relative *compression* and *stretching*, respectively, in
194 the subjective value function, \hat{V}_τ (Figure 3b), with corresponding compression/stretching of the
195 RPE function (Figure 3c).

196



197

198 **Figure 3 | Compressed and stretched Value and RPE Landscapes for a sub-optimal**
 199 **timer predict dynamic DAN responses during the self-timed movement task, but do not**
 200 **capture baseline offsets.** **a**, Simple state space of self-timed movement task for a suboptimal
 201 timer with a fast pacemaker. The fast pacemaker “compresses” state space^{3,21}, resulting in
 202 traversal of the timing states faster than veridical time. The mouse can only make a decision
 203 based on which state it believes itself in; thus first-lick is expected to occur early (gold-shaded
 204 state). **b**, A compressed subjective value function (\hat{V}_τ , blue) reflects relatively fast traversal
 205 through the value landscape compared with that of veridical time (V_t , black). Conversely,
 206 stretched \hat{V}_τ (red) reflects slow traversal, consistent with a slow pacemaker. The animal is
 207 expected to lick at the peak of the trajectory. **c**, Smoothed estimated RPE function ($\hat{V}'_\tau \approx \delta_\tau$).
 208 Compression/stretching of the value function produces ramping dynamics similar to those
 209 observed in DANs (**d**) and striatal dopamine¹. However, this model alone does not explain the
 210 more tonic baseline offsets that were anti-correlated with upcoming movement time (**d** and
 211 [Figure 1b](#)). **d**, Average DAN GCaMP6f signals (12 mice, 3 timepoints replotted from [Figure](#)
 212 [1b](#), plotted up to 150 ms before first-lick). Break in x-axis as in [Figure 1b](#).
 213

214 Strikingly, as this simple RPE-based model predicts, DAN signals observed during our self-timed
 215 movement task show different ramping dynamics depending on when the animal actually moved
 216 ([Figure 3d](#)), consistent with compression/stretching of the subjective value and RPE functions.

217 When the animal moved relatively early (perhaps corresponding to a fast pacemaker), DAN
218 ramping unfolded with a steeper slope, as if the ramping period were *compressed*. Conversely,
219 when the animal moved late (perhaps corresponding to a slow pacemaker), DAN ramping unfolded
220 with a shallower slope, as if the ramping interval were *stretched*. The idea of
221 compression/stretching of DAN ramps was supported by our encoding model¹, for which we
222 needed to add a timing-dependent “stretch factor” to best capture the variance in GCaMP6f signals
223 during the timed interval. Together, these observations could be explained by DANs encoding an
224 RPE-like signal related to the animal’s “belief” of its position in objective time, τ , as derived from
225 its position along the subjective value trajectory during the timing interval of the task.

226

227 In fact, a recent model described how a timing mechanism instantiated by the nigrostriatal system
228 could lead to (the well-known) variability in self-timed intervals by stretching or compressing of
229 subjective value trajectories³. The model posits that dopamine modulates the pacemaker rate
230 (consistent with pharmacological and lesion studies), with increased dopamine availability (or
231 efficacy) speeding the pacemaker, and decreased dopamine slowing the pacemaker^{6-8,11-13}. In turn,
232 the pacemaker controls the encoding of subjective time, and thus the steepness of the value
233 function with respect to objective, veridical time. It follows that variation in dopamine availability
234 would compress or stretch the value landscape to varying degrees from trial-to-trial. This model is
235 consistent with our findings of variable ramping slope in DANs signals from trial-to-trial. It is also
236 consistent with neural recordings from striatal spiny projection neurons and parietal cortical
237 neurons during similar self-timed movement tasks, for which temporal sequences of striatal and
238 cortical firing during timing were compressed for early movements and stretched for late
239 movements^{23,24}.

240

241 While the RPE-based view of DAN activity captures the dynamic DAN signals we observed, our
242 simple RPE model alone does not capture the *baseline offsets* in DAN signals that were predictive
243 of movement timing even after controlling for previous trial outcome and ongoing nuisance
244 movements¹. More complex RPE-based explanations for these *tonic* offsets in DAN signals could
245 be imagined with further assumptions (e.g., states like the pre-cue delay could also contain timing
246 states that create offsets before the trial begins, etc.), but a parsimonious explanation for how and
247 why these offsets emerge requires further investigation. Mohebi *et al.* recently showed baseline
248 differences in the amount of dopamine in the nucleus accumbens core that were correlated with
249 the recent history of reward rate: higher recent reward rates were related to higher tonic dopamine¹⁷.
250 However, in our task, animals tended to move later toward the end of sessions, resulting in periods
251 of relatively high reward rate when the average tonic baseline signal was *lower* (baseline preceding
252 rewarded trials—by definition, later movements—was systematically lower in our task, [Figure 1b-](#)
253 [c](#)), suggesting a more complex relationship between tonic DAN activity and reward rate in our
254 task. While the origin of offsets in DAN signals remains unclear, these offsets were nonetheless
255 inversely related to the first-lick time, and thus directly related to the (inferred) pacemaker rate,
256 consistent with pharmacological and lesion studies positing a positive correlation between
257 dopamine availability and pacemaker rate^{3,6-8,11-13}.

258

259 Ramping signals in our photometry experiments were measured from a population of DANs. An
260 important future question is whether ramps are also present at the level of individual neurons, or
261 rather represent a progressive recruitment of individual neurons, or some combination of both.
262 Prior studies have reported ramping signals in individual neurons during tasks with visual feedback

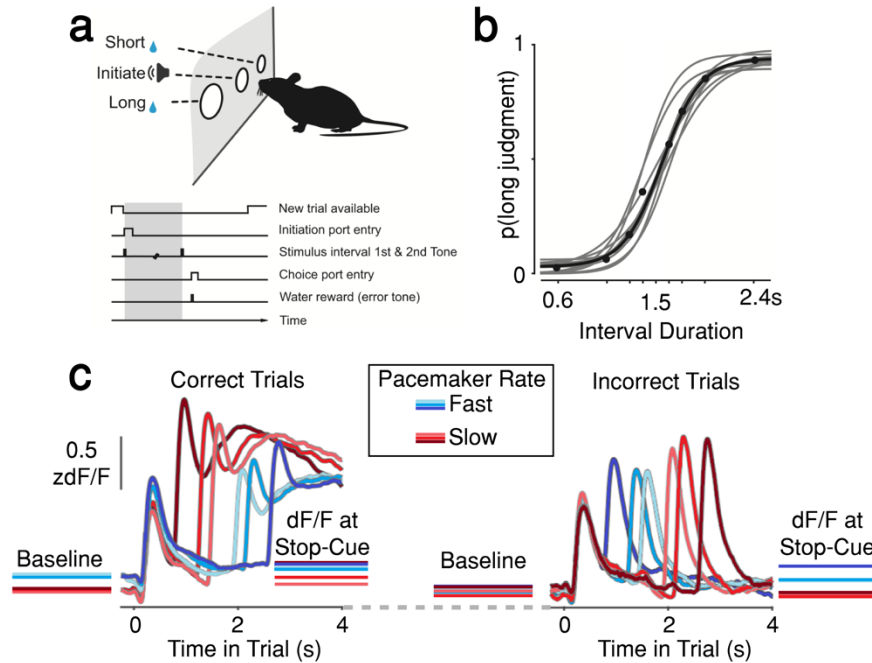
263 of distance to reward¹⁸, whereas others have observed decoupling between DAN firing rates and
264 downstream dopamine release¹⁷, making it unclear whether electrophysiology would be capable
265 of addressing this question. Observation of individual neurons expressing calcium indicators with
266 GRIN-lens equipped endoscopes may be better suited to this question.

267

268 *RPE-based predictions for DAN responses during a temporal bisection task*

269 Whereas DAN signals during our self-timed movement task were consistent with classic
270 observations of the influence of dopamine on the speed of the pacemaker, a recent study employing
271 a different timing task found more complex DAN dynamics during timing. Soares *et al.* recorded
272 SNc DAN GCaMP6f signals with fiber photometry as mice executed a classic temporal bisection
273 perceptual task⁵ (Figure 4a). Trials began when mice entered a nose-poke port and received an
274 auditory start-timing cue. Mice had to remain in the port throughout a variable timing interval,
275 which was terminated with a stop-timing auditory cue. Mice then reported whether the interval
276 was shorter or longer than a criterion time (1.5 s) by choosing a left or right nose-poke port
277 corresponding to a “long” or “short” judgment. Mice were trained to categorize intervals spanning
278 0.6-2.4 s. As expected, trials with more extreme intervals were easier for the mice, whereas trials
279 with intervals closer to the 1.5 s criterion time elicited chance performance (Figure 4b).

280



281

282 **Figure 4 | A temporal bisection task shows relatively high DAN signals during the**
283 **timing interval when the inferred pacemaker rate is relatively fast.** Figures adapted from
284 [Soares et al., 2016](#)⁵ with permission of authors and AAAS. **a**, Task schematic. **b**,
285 Psychometric curve for timing intervals of different duration. Criterion time: 1.5 s. **c**, Start-
286 timing cue-aligned average Snc DAN GCaMP6f signals. Second peak occurs just after the
287 stop-timing cue (intervals: 0.6, 1.05, 1.26, 1.74, 1.95, 2.4 s). Figure recolored to indicate
288 average inferred pacemaker rate. Red: slow; blue: fast. Note: colors intended to indicate
289 category of clock speed, *not* relative pacemaker speed within category. Relative dF/F
290 amplitude during baseline and immediately prior to stop-timing cue shown left and right.
291 dF/F amplitudes during timing are higher when the inferred pacemaker rate is fast. Left:
292 Correct trials. Right: Incorrect trials show the same dF/F relationship with pacemaker rate.
293

294 DANs exhibited complex dynamics during the bisection task, starting with a sharp transient after
295 the start-timing cue and ending with a second transient after the stop-timing cue ([Figure 4c](#)).
296 Between the start-timing and stop-timing cues, DAN signals exhibited a U-shape with increasing
297 time, which was visible for trials with longer intervals but was truncated prematurely for the shorter
298 intervals. The authors focused their analyses on the transient occurring *after the stop-timing* cue.
299 Short judgments (suggesting a slow pacemaker) were accompanied by relatively high-amplitude
300 transients after the stop-cue, whereas long judgments (suggesting a fast pacemaker) showed

301 relatively low-amplitude transients. These results seemed to suggest that relatively *high* DAN
302 activity reflected a *slow* pacemaker, the opposite of what is expected based on the bulk of
303 pharmacological and lesion studies³, as well as the trend we observed during our self-timed
304 movement task.

305
306 This surprising finding could be a unique feature of the bisection task. Unlike self-timed
307 movements, in which animals directly report elapsed time with a movement, the temporal bisection
308 task requires an additional computational step, in which the timed interval must be categorized as
309 “long” or “short.” However, prior pharmacological studies employing the bisection task found
310 results consistent with the classic view that higher dopamine availability is associated with a faster
311 pacemaker^{3,25}—opposite the interpretation of Soares *et al.*, but consistent with the findings of our
312 self-timed movement task.

313
314 The discrepancy between our results and those found by Soares *et al.* could perhaps be traced to
315 differences in the way DAN signals were analyzed. We focused our attention on DAN signals
316 unfolding *during timing* in our self-timed movement task, whereas these signals were not explored
317 by Soares *et al.* We thus asked two questions: 1. What correlations exist between DAN signals and
318 pacemaker rate in the bisection task *before* the timing interval? And, 2. What correlations exist
319 *during* the timing interval itself?

320
321 Before addressing these questions, we note that the relationship between pacemaker and bisection
322 judgment is not as straightforward as in self-timed movement, and thus we recolored [Figure 4c](#) to
323 clarify this, employing the following intuition: For a trial to be correct in the bisection task, on

324 average, the pacemaker must be either accurate or “conservatively inaccurate.” In other words, a
325 correct “short” judgment requires either accurate timing or a *slow* pacemaker (Figure 4c, red
326 curves). Conversely, a correct “long” judgment requires either accurate timing or a *fast* pacemaker
327 (Figure 4c, blue curves).

328

329 When we considered DAN signals *before* the timing interval for correct trials in the Soares *et al.*
330 study (Figure 4c, left), we noticed what appears to be two strata of signal levels. Trials with “long”
331 judgments (fast pacemaker on average) had relatively high baseline signals, whereas trials with
332 “short” judgments (slow pacemaker on average) had lower baseline signals, consistent with the
333 relationship between baseline offsets and pacemaker rate that we observed in our self-timed
334 movement task. As in our task, these baseline offsets remained present during the timing interval,
335 resulting in the same stratification of dF/F signals immediately prior to the stop-timing cue (except
336 for the very shortest interval, 0.6 s, which overlaps decaying GCaMP6f signals related to the start-
337 timing cue, likely causing an artifactual inflation of the signal just prior to the stop-cue due to the
338 off-kinetics of the calcium indicator or kinetics of calcium clearance more generally). Thus, it
339 generally appears that DAN activity was *higher* on trials with fast pacemaker rates, both during
340 and before the interval in which the animal was actually timing. Intriguingly, *incorrect* trials (to
341 the right in Figure 4c) showed a relative convergence of the baseline signals preceding the start-
342 cue, but then signals diverged during the timing interval, resulting in relatively *high* signals just
343 before the stop-cue for incorrect “long” choices (i.e., a fast pacemaker, blue), but relatively low
344 signals just before the stop-cue for incorrect “short” choices (i.e., a slow pacemaker, red). This is
345 consistent with the patterns observed on correct trials. Interpreted thusly, the Soares *et al.* result is
346 consistent both with our results and with classic pharmacological studies relating higher/lower

347 dopamine availability to faster/slower pacemaker rates, respectively. Soares *et al.* presented their
348 subsequent analyses with these baseline differences normalized-out in some way (Figure 3 of
349 Soares *et al.*). It is possible that this “zeroing out” of the baseline offset may have hindered efforts
350 to detect consistent effects during the timing interval due to reordering of the traces.

351

352 Because baseline offsets in the bisection task appear similar to those in our self-timed movement
353 task, we asked whether dynamic DAN signals in the bisection task could similarly be explained
354 by the task’s RPE landscape. In their investigation of the stop-timing cue-related transient, Soares
355 *et al.* showed that its amplitude is well-explained by a combination of temporal surprise and
356 behavioral performance, and we applied these parameters to derive a value landscape consistent
357 with their bisection task.

358

359 The inferred value landscape of the bisection task for an optimal agent was built from a few
360 assumptions (Figure 5a):

361

- 362 1. As in our self-timed movement task, value increases immediately at the start-cue and
363 continues to rise toward the time of expected potential reward delivery.
- 364
- 365 2. Because the longest interval is 2.4 s, the time until potential reward is known to be no more
366 than ~3 s (including the time to report judgment). However, due to temporal uncertainty
367 and the fact that a false start (leaving the port before the stop-timing cue) results in an error
368 and loss of reward, there is a second jump in the value function at the time of the stop-cue

369 when the feedback of the tone reorients the value function and indicates the opportunity to
370 collect reward within a few hundred milliseconds.

371

372 3. Because value is temporally discounted at the start-cue by the possibility of the longest-
373 possible interval, any stop-cue occurring before 2.4 s results in a sudden “teleportation”
374 through the value landscape to the final limb of the task that occurs just before the judgment
375 and ascertainment of trial outcome, similar to the jump in the value function in a recently-
376 reported, virtual reality, spatial teleportation task¹⁸. Thus, assuming the value function
377 trends upwards steadily, the amplitude of RPE-related transients following the stop-cue
378 would *decrease* as the interval duration increases, because the sudden jump in the value
379 function becomes progressively smaller.

380

381 4. To capture aspects related to behavioral performance, we additionally included contours in
382 the value function during the timing interval to reflect the probability of a correct choice
383 for intervals of different lengths. Specifically, a relative minimum in the value function
384 occurs near 1.5 s, when predicted performance is worst. However, a stop-timing tone near
385 the criterion time also results in a smaller jump in the value function because the probability
386 of a correct decision is also lower. Thus, the increase in value at the moment of decision
387 was adjusted by the probability of a correct choice.

388

389 5. As in the simple RPE-model of our self-timed movement task, we modeled changes in
390 pacemaker rate as compression/stretching of the subjective value landscape with respect to
391 veridical time.

392

393 6. The agent traverses timing states during the timing interval, similar to the timing states in
394 the self-timed movement task, but unlike our task, the bisection task does not require the
395 agent to decide when to move. We assume the need to make a timed movement imposes a
396 need for the agent to be relatively certain of its subjective timing state, τ , to make a decision,
397 even though it is uncertain of its true state, t . The bisection task, on the other hand, is more
398 similar to classical conditioning tasks in which the timing interval is not in the agent's
399 control, and thus subjective state uncertainty increases with the distance from the last state-
400 informative cue³. Thus, we took into account temporal blurring of the subjective state
401 function, which would tend to reduce the convexity of the subjective value function and
402 reduce the amplitude of ramping during the timing interval³. However, adding temporal
403 blurring does not substantially change the fit-shape in our simplified model, and versions
404 with or without blurring can reproduce the shape of the dynamic DAN signals.

405

406 Together, we arrived at a model of the RPE landscape for each of the six tested interval durations
407 (Figure 5b,c). Importantly, this simple RPE-based model accurately captures the relative
408 categorical amplitudes of the stop-timing cue-related transients, as follows: If the instantaneous
409 DAN activity at the time of the stop-timing cue is relatively high, this would indicate that the
410 animal is further along in the subjective value trajectory, resulting in 1) a *long* judgment, and 2) a
411 relatively *smaller* RPE transient, because the underlying subjective value was higher at that
412 moment. Conversely, if instantaneous DAN activity is relatively low at the stop-timing cue, this
413 would indicate that the animal is not very far along the subjective value trajectory, leading to 1) a

414 *short* judgement and 2) a relatively *larger* stop-cue-related RPE transient, because the underlying
415 subjective value was relatively low just before the stop-cue.

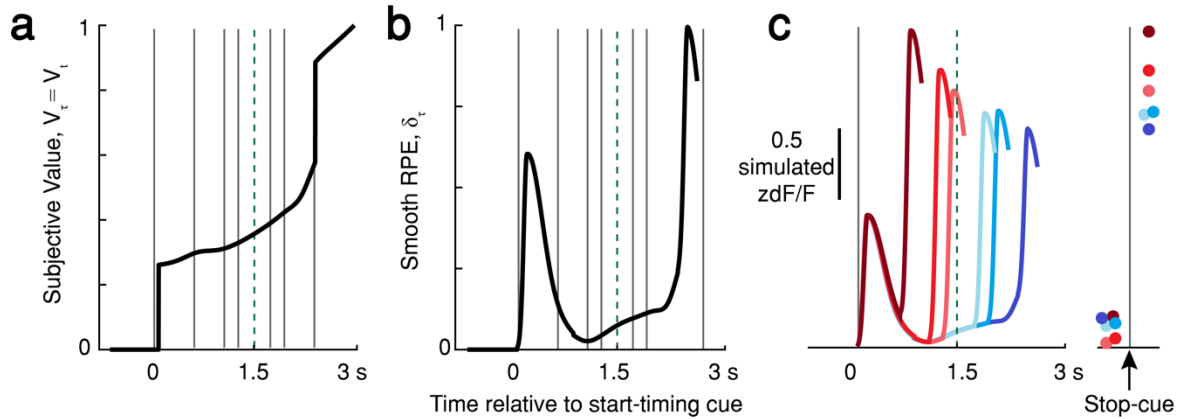
416

417 Now consider a particular (objective) time interval near the criterion time, for which the animal
418 makes a mix of “long” and “short” choices (e.g., 1.74 s; [Figure 4b](#)). Soares *et al.* found that the
419 amplitude of the stop-timing cue-related GCaMP6f transient tended to be bigger when the animal
420 incorrectly made short choices, and this was taken as evidence that elevated DAN activity *slows*
421 the internal clock. However, our model predicts that the size of the stop-cue-related transient will
422 be inversely related to the amplitude of the underlying subjective value at that point, and thus
423 inversely related to elapsed *subjective* time. It thus follows that if subjective time is more advanced
424 on a given trial (i.e., faster pacemaker), the animal would tend to choose the long judgment on that
425 trial, and the stop-timing RPE transient would be *smaller*. Conversely, if subjective time is less
426 advanced on a trial (i.e., slower pacemaker), the animal would tend to choose the short judgment,
427 and the stop-timing RPE transient would be *larger*.

428

429 Our RPE model accurately predicts the results of Soares *et al.*; however, our model holds that
430 elevated DAN activity *speeds* the internal clock, consistent with most pharmacological studies but
431 *opposite* the interpretation of Soares *et al.* Thus, our RPE-based model suggests a parsimonious
432 explanation for DAN activity in both the self-timed movement and temporal bisection paradigms,
433 with (1) relatively high DAN activity corresponding to a fast pacemaker; manifesting in (2)
434 compression of the value landscape; thereby leading to (3) early movements (in the self-timed
435 movement task) or long judgments (in the temporal bisection task).

436



437

438 **Figure 5 | Subjective Value and RPE Landscapes for the temporal bisection task predict**
439 **dynamic DAN responses during the temporal bisection task, but do not capture baseline**
440 **offsets. a**, Estimated value function \hat{V}_t , where $\hat{V}_t \approx V_t$ for an optimal timer on a 2.4 s trial.
441 Grey lines: test interval times. Green dashed line: criterion time (1.5 s). Value increases
442 approaching the time when reward is available, increasing abruptly at the start- and stop-
443 timing cues (0 and 2.4 s). **b**, Smoothed RPE function for an optimal timer, estimated as $\delta_t \approx$
444 \hat{V}'_t , the derivative of the subjective value function. The RPE function was smoothed with
445 an asymmetrical gaussian kernel spanning *ca.* 28% of the interval to approximate GCaMP6f
446 off-dynamics. **c**, Predicted DAN GCaMP6f signals for an optimal timer for the six test
447 interval times. Traces truncated before reward collection for clarity. Colors indicate
448 conservative pacemaker speed for a correct judgment (red: slow, blue: fast). Right: relative
449 simulated dF/F amplitude just before the stop-timing cue and subsequent peak response.
450 *Amplitude just before the stop-timing cue is directly proportional to pacemaker speed; peak*
451 *amplitude after the stop-timing cue is inversely proportional to pacemaker speed.*
452

453 *Limitations of the RPE-based model*

454 The simple RPE-based models presented here explain dynamic DAN signals in both the bisection
455 task and our self-timed movement task, but they do not explain the origin of baseline offsets.
456 Mohebi *et al.*¹⁷ recently-proposed that baseline offsets in ventral striatal dopamine levels could
457 reflect the average recent reward rate, but we found that offset amplitude in DAN signals is at least
458 partially independent of recent trial history during the self-timed movement task. It is possible that
459 baseline variation arises from slow, random fluctuations in DAN activity, but further work is
460 needed to explore the origins of these signals.

461

462 A second issue is the impact of optogenetic DAN activation and suppression on the rate of the
463 pacemaker. In our self-timed movement task, DAN activation promoted early movements,
464 consistent with increasing the pacemaker rate, whereas suppression promoted late movements,
465 consistent with slowing the pacemaker rate¹. However, Soares *et al.* reported an opposite effect for
466 optogenetic manipulation during the bisection task, at least for DAN activation⁵.

467

468 This difference between the tasks could be reconciled by a recent theoretical model proposed by
469 Mikhael and Gershman to explain the behavior of the pacemaker in a wide range of classical
470 conditioning and timing studies³. Their model shows that the pacemaker rate is expected to be
471 updated at the time of reinforcement by a Hebbian-like, bidirectional learning rule. If reward
472 occurs exactly at the expected time, there is no update in the pacemaker rate. However, if
473 reinforcement occurs before the expected time, this is interpreted as feedback that the pacemaker
474 was running too slowly; thus, the update rule increases the pacemaker rate leading to expectation
475 of reward at an earlier time on the next trial. Conversely, if reinforcement occurs after it was
476 expected, this is interpreted as feedback indicating an overly fast pacemaker, resulting in an update
477 that slows the pacemaker rate and creates the expectation of a later reward on the next trial. The
478 same principles apply to ongoing RPE during timing tasks.

479 In our self-timed movement task, we activated or inhibited DANs only *up to* the time of first-lick,
480 which Mikhael and Gershman's model predicts will produce an effect on the pacemaker rate
481 consistent with the sign of the manipulation (activate: increase, inhibit: decrease). However, Soares
482 *et al.* continued optical stimulation *past* the end of the timing interval, until the end of the trial.
483 When Mikhael and Gershman modeled stimulation in the Soares *et al.* task, they found that
484 simulated DAN activation increased the pacemaker rate during the timing interval, but the

485 continuing stimulation after the stop-timing cue rapidly counteracted this effect, resulting in
486 *slowing* of the modeled pacemaker between the stop-cue and the judgment, leading to an effect on
487 pacemaker rate *inconsistent* with the sign of the manipulation, as observed in Soares *et al.* If this
488 model is correct, the effect of stimulation on the animal's judgment in the Soares *et al.* task may
489 have arisen due to continued manipulation of DAN activity *after* the timing interval had ended. A
490 "retrospective" effect of this sort might seem counterintuitive, but such retrospective effects have
491 long been observed in perceptual studies, in which recall of sensory stimuli can be enhanced by
492 additional sensory cues presented shortly after stimulus offset^{26,27}, suggesting that sensory events
493 are "buffered" briefly and can be altered by neural activity occurring between the sensory event
494 and the perceptual decision. It is possible that a similar process could occur in the bisection task if
495 DAN stimulation extends past the timing interval, although this is speculative. More work is
496 needed to reconcile the optogenetic results in the self-timed movement and bisection tasks. To
497 start, it would be informative to repeat the optogenetic experiments in the bisection task with
498 optical stimulation limited to the period of the timed intervals only.

499

500

501

502

503

504

505 **REFERENCES**

506

507 1. Hamilos, A.E. *et al.* Dynamic dopaminergic activity controls the timing of self-timed
508 movement. Preprint at: <https://doi.org/10.1101/2020.05.13.094904> (2020).

509 2. Mikhael, J. G., Kim, H. R., Uchida, N., & Gershman, S. J. Ramping and State
510 Uncertainty in the Dopamine Signal. Preprint at

511 <https://www.biorxiv.org/content/10.1101/805366v1> (2019).

512 3. Mikhael, J. G. & Gershman, S. J. Adapting the flow of time with dopamine. *J.*
513 *Neurophysiol.*, **121**, 1748–1760 (2019).

514 4. Gershman, S. J. Dopamine Ramps Are a Consequence of Reward Prediction
515 Errors. *Neural Comp.* **26**, 467–471 (2014).

516 5. Soares, S., Atallah, B. & Paton, J. Midbrain dopamine neurons control judgment of
517 time. *Science* **354**, 1273–1277 (2016)

518 6. Dews, P. B. & Morse, W. H. Some observations on an operant in human subjects and
519 its modification by dextro amphetamine. *J. Exp. Anal. Behav.* **1**, 359–364 (1958).

520 7. Schuster, C. & Zimmerman, J. Timing behavior during prolonged treatment with dl-
521 amphetamine. *J. Exp. Anal. Behav.* **4**, 327–330 (1961).

522 8. Meck, W. H. Affinity for the dopamine D2 receptor predicts neuroleptic potency in
523 decreasing the speed of an internal clock. *Pharmacol. Biochem. Behav.* **25**, 1185–
524 1189 (1986).

525 9. Malapani, C., Rakitin, B. C., Levy, R., Meck, W. H., Deweer, B., Dubois, B., &
526 Gibbon, J. (1998). Coupled Temporal Memories in Parkinson’s Disease: A
527 Dopamine-Related Dysfunction. *J. Cog. Neuro.*, 3(May), 316–331.

- 528 10. Raitin, B. C., Penney, T. B., Gibbon, J., Hinton, S. C., & Meek, W. H. Scalar
529 Expectancy Theory and Peak-Interval Timing in Humans. *Journal of Experimental*
530 *Psychology: Animal Behavior Processes* **24**, 15–33 (1998).
- 531 11. Lutig, C. & Meck, W. H. Chronic treatment with haloperidol induces deficits in
532 working memory and feedback effects of interval timing. *Brain Cogn.* **58**, 9–16
533 (2005).
- 534 12. Meck, W. H. Neuroanatomical localization of an internal clock: A functional link
535 between mesolimbic, nigrostriatal, and mesocortical dopaminergic systems. *Brain*
536 *Res.* **1109**, 93–107 (2006).
- 537 13. Merchant, H., Harrington, D. L. & Meck, W. H. Neural Basis of the Perception and
538 Estimation of Time. *Annu. Rev. Neurosci* **36**, 313–36 (2013).
- 539 14. Schultz, W., Dayan, P. & Montague, P. R. A neural substrate of prediction and
540 reward. *Science* **275**, 1593–1599 (1997).
- 541 15. Howe, M. W., Tierney, P. L., Sandberg, S. G., Phillips, P. E. M. & Graybiel, A. M.
542 Prolonged dopamine signalling in striatum signals proximity and value of distant
543 rewards. *Nature* **500**, 575–579 (2013).
- 544 16. Hamid, A. A. et al. Mesolimbic dopamine signals the value of work. *Nat. Neurosci.*
545 **19**, 117–126 (2016).
- 546 17. Mohebi, A et al. Dissociable dopamine dynamics for learning and motivation. *Nature*
547 **570**, 65–70 (2019).
- 548 18. Kim, H. R. et al. A unified framework for dopamine signals across timescales.
549 Preprint at <https://www.biorxiv.org/content/10.1101/803437v1> (2019).

- 550 19. Engelhard, B et al. Specialized coding of sensory, motor and cognitive variables in
551 VTA dopamine neurons. *Nature* **570**, 509–513 (2019).
- 552 20. Starkweather, C. K., Babayan, B. M., Uchida, N., & Gershman, S. J. Dopamine
553 reward prediction errors reflect hidden-state inference across time. *Nat. Neuro.* **20**,
554 581–589 (2017).
- 555 21. Sutton, R. S., & Barto, A. G. *Reinforcement Learning*. MIT Press, Cambridge,
556 (2000).
- 557 22. Gibbon, J., Malapani, C., Dale, C. L., & Gallistel, C. Toward a neurobiology of
558 temporal cognition: advances and challenges. *Curr. Opin. Neurobiol.* **7**, 170–184
559 (1997).
- 560 23. Mello, G. B. M., Soares, S. & Paton, J. J. A scalable population code for time in the
561 striatum. *Curr. Biol.* **25**, 1113–1122 (2015).
- 562 24. Maimon, G. & Assad, J. A. A cognitive signal for the proactive timing of action in
563 macaque LIP. *Nat. Neuro.* **9**, 948–955 (2006).
- 564 25. Morgan, L., Killeen, P.R., Fetterman, J.G. Changing rates of reinforcement perturbs
565 the flow of time. *Behav. Processes* **30**, 259–271 (1993).
- 566 26. Gegenfurtner, K. R., & Sperling, G. Information Transfer in Iconic Memory
567 Experiments. *J. Exp. Psych: Human Perception and Performance* **19**, 845–866 (1993).
- 568 27. Herrington, T. M., & Assad, J. A. Neural activity in the middle temporal area and
569 lateral intraparietal area during endogenously cued shifts of attention. *J. Neurosci.* **29**,
570 14160–14176 (2009).

571

572

573 **ACKNOWLEDGEMENTS**

574 We thank J.G. Mikhael and S.J. Gershman for discussions on temporal difference learning models
575 and analytical methods; The work was supported by NIH grants UF-NS109177 and U19-
576 NS113201, and NIH core grant EY-12196. A.E.H. was supported by a Harvard Lefler Predoctoral
577 Fellowship and a Harvard Quan Predoctoral Fellowship.

578

579 **AUTHOR CONTRIBUTIONS**

580 A.E.H. and J.A.A conceived the project. A.E.H. performed all experiments and implemented the
581 computational models. A.E.H. and J.A.A. analyzed the data and wrote the paper.

582

583 **DECLARATION OF INTERESTS**

584 The authors have no relevant interests to declare.

585

586 **CODE AVAILABILITY**

587 All custom behavioral software and analysis tools are available

588 at <https://github.com/harvardschoolofmouse>.

589

590 **DATA AVAILABILITY**

591 The data that support the findings of this study are available from the corresponding author upon
592 reasonable request.