

1 **Distinct value computations support rapid** 2 **sequential decisions**

3 Andrew Mah¹, Shannon S. Schiereck¹, Veronica Bossio^{1†},
 Christine M. Constantinople^{1*}

¹ Center for Neural Science, New York University; New York, NY 10003.

† Present address: Zuckerman Institute, Columbia University; New York, NY 10027.

*Corresponding author. E-mail: constantinople@nyu.edu.

4 **The value of the environment determines animals’ motivational states and sets**
5 **expectations for error-based learning¹⁻³. How are values computed? Rein-**
6 **forcement learning systems can store or “cache” values of states or actions**
7 **that are learned from experience, or they can compute values using a model of**
8 **the environment to simulate possible futures³. These value computations have**
9 **distinct trade-offs, and a central question is how neural systems decide which**
10 **computations to use or whether/how to combine them⁴⁻⁸. Here we show that**
11 **rats use distinct value computations for sequential decisions within single tri-**
12 **als. We used high-throughput training to collect statistically powerful datasets**
13 **from 291 rats performing a temporal wagering task with hidden reward states.**
14 **Rats adjusted how quickly they initiated trials and how long they waited for re-**
15 **wards across states, balancing effort and time costs against expected rewards.**
16 **Statistical modeling revealed that animals computed the value of the environ-**
17 **ment differently when initiating trials versus when deciding how long to wait**

18 **for rewards, even though these decisions were only seconds apart. Moreover,**
19 **value estimates interacted via a dynamic learning rate. Our results reveal how**
20 **distinct value computations interact on rapid timescales, and demonstrate the**
21 **power of using high-throughput training to understand rich, cognitive behav-**
22 **iors.**

24 **Introduction**

25 There are many ways to compute value. Reinforcement learning provides a powerful frame-
26 work for describing how animals or agents learn the value of different states and actions from
27 experience and use those value estimates to guide behavior³. The value of the environment, or
28 how much reward it is expected to yield, is important for motivation and sets expectations for
29 reinforcement learning¹⁻³.

30 There are many reinforcement learning methods for computing value that differ in their
31 implementation, computational demands, and flexibility^{3,6,8,9}. For instance, some algorithms
32 use a model of the world to flexibly estimate the value of states or actions by mental simu-
33 lation or planning. Other algorithms cache values from direct experience, without an explicit
34 model of the environment. These different reinforcement learning methods, which remarkably
35 are thought to be supported by distinct neural circuits^{10,11}, have trade-offs between flexibility
36 and computational efficiency^{6,8,9}. They also represent two ends of a continuum⁴⁻⁸. A central
37 question in neuroscience and psychology is determining how values are computed in animals
38 including humans¹². Moreover, neurobiologically-inspired value computations will likely lead
39 to advances in next generation artificial intelligence¹³.

40 However, it is difficult to determine the value computations that subjects use, especially over

41 behaviorally relevant timescales of seconds. In standard two-alternative forced choice tasks, the
42 behavioral read-out is a binary choice, and the underlying values driving choice are obscure.
43 State-of-the-art methods for revealing how values are computed use regression models that pool
44 data over entire behavioral sessions¹⁴, or pre-determined subsets of trials¹⁵, thereby obscuring
45 moment-by-moment changes in value computations. Therefore, whether or how multiple value
46 computations interact on rapid timescales in the same subject is unclear.

47 **Results**

48 **Rats' deliberative and motivational decisions are sensitive to the value of** 49 **the environment.**

50 We developed a temporal wagering task for rats, in which they were offered one of several
51 water rewards on each trial, the volume of which (5, 10, 20, 40, 80 μ L) was indicated by a tone
52 (Fig. 1a). The reward was assigned randomly to one of two ports, indicated by an LED. The rat
53 could wait for an unpredictable delay to obtain the reward, or at any time could terminate the
54 trial by poking in the other port (“opt-out”). Wait times were defined as how long rats waited
55 before opting out. Trial initiation times were defined as the time from opting-out or consuming
56 reward to initiating a new trial. Reward delays were drawn from an exponential distribution, and
57 on 15-25 percent of trials, rewards were withheld to force rats to opt-out, providing a continuous
58 behavioral readout of subjective value (Fig. 1b)¹⁶⁻¹⁸. We used a high-throughput facility to train
59 291 rats using computerized, semi-automated procedures. The facility generated statistically
60 powerful datasets (median = 33,493 behavioral trials, 71 sessions).

61 The task contained latent structure: rats experienced blocks of 40 completed trials (hidden
62 states) in which they were presented with low (5, 10, or 20 μ L) or high (20, 40, or 80 μ L) re-
63 wards¹⁷. These were interleaved with “mixed” blocks which offered all rewards (Fig. 1c). 20 μ L
64 was present in all blocks, so comparing behavior on trials offering this reward revealed contex-

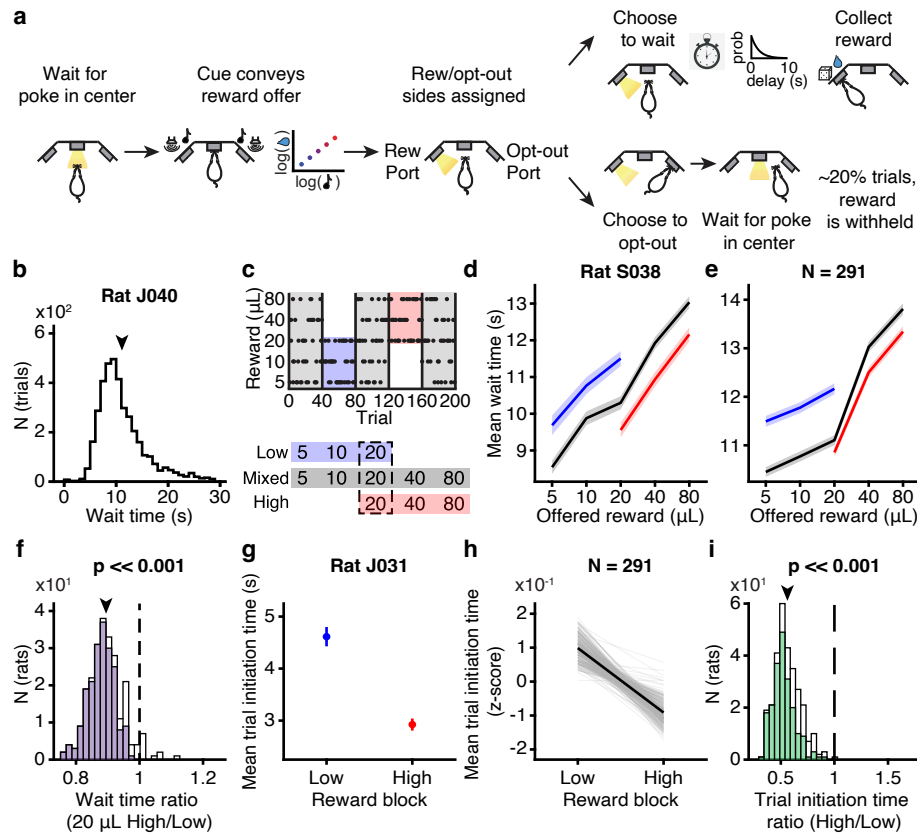


Figure 1: Wait time and trial initiation time were modulated by the value of the environment. **a.** Schematic of behavioral paradigm. **b.** Distribution of wait times for one rat. **c.** Block structure of task. **d-e.** Average wait time on catch trials by reward in each block for (d) one rat and (e) averaged across rats. **f.** Wait time ratio (average wait time for 20 μL in high block/low block) across all rats. Filled boxes indicated rats with $p < 0.05$, Wilcoxon rank-sum test. Population average, $p << 0.001$, Wilcoxon signed-rank test, $N = 291$. **g-h.** Average trial initiation times in high and low blocks for (g) one rat and (h) all rats. **i.** Trial initiation time ratio (average initiation time in high block/low block) across all rats. Filled boxes indicated rats with $p < 0.05$, Wilcoxon rank-sum test. Population average, $p << 0.001$, Wilcoxon signed-rank test, $N = 291$.

65 tual effects (i.e., effects of hidden states). The hidden states differed in their average reward
66 and therefore in their opportunity costs, or what the rat might miss out on by continuing to wait.
67 According to foraging theories, the opportunity cost is the long-run average reward, or the value
68 of the environment¹⁹. In accordance with these theories^{19,20}, rats adjusted how long they were
69 willing to wait for rewards in each block, and on average waited ~ 10 percent less time for 20 μ L
70 in high blocks, when the opportunity cost was high, compared to in low blocks ($p \ll 0.001$,
71 Wilcoxon signed-rank test, $N = 291$; Fig. 1d-f). These are strong contextual effects compared
72 to previous studies^{17,21}.

73 Animals make more vigorous actions when those actions are expected to yield larger or
74 more valuable rewards^{2,22–25}. Therefore, we analyzed how quickly rats initiated trials, as this
75 might also reflect the perceived value of the environment. Indeed, trial initiation times were
76 modulated by blocks in a similar pattern as the wait times, with rats initiating trials more quickly
77 in high compared to low blocks ($p \ll 0.001$, Wilcoxon signed-rank test, $N = 291$; Fig. 1g-i;
78 Extended Data Fig. 1). Previous work suggests that this pattern optimally balances the energetic
79 costs of vigor against the benefits of harvesting reward in environments with different reward
80 rates^{2,25,26}. Therefore, both the trial initiation times, which reflect motivation, and the wait
81 times, which reflect deliberating between waiting and opting-out, were modulated by the value
82 of the environment.

83 Notably, while we used all behavioral trials for analyses of initiation times in this study,
84 sensitivity to the reward blocks was largely driven by initiation times following unrewarded
85 trials (Extended Data Fig. 2, Methods), which accounted for more variance in initiation times.
86 This is consistent with previous studies showing that response outcomes can gate behavioral
87 strategies^{27,28}. There were no major differences in wait times following rewarded or unrewarded
88 trials (Extended Data Fig. 3). To make comparisons between trial initiation and wait times with
89 as much statistical power as possible and the fewest assumptions, we used all behavioral trials

90 for subsequent analyses in this study.

91 **Trial initiation and wait times exhibited distinct temporal dynamics.**

92 Surprisingly, wait and trial initiation times exhibited dramatically different dynamics at
93 block transitions. In mixed blocks, the wait times following high and low blocks converged
94 to a common value, regardless of the previous block type, suggesting the use of a fixed esti-
95 mate of environmental value in mixed blocks (Fig. 2a). Trial initiation times, however, showed
96 longer timescale effects such that initiation times in mixed blocks strongly depended on the
97 previous block identity (Fig. 2b; Extended Data Fig. 4). These longer timescale dynamics,
98 which are reminiscent of incentive contrast effects²⁹, were also evident in the transitions from
99 mixed blocks into high/low blocks for trial initiation times, but not wait times (Extended Data
100 Fig. 5), indicating that trial initiation and wait times utilize distinct estimates of the value of the
101 environment.

102 To better characterize their temporal dynamics, we regressed the trial initiation and wait
103 times against rewards offered on previous trials. We included current rewards as regressors
104 in the wait time model, and restricted this analysis to mixed blocks only. Examination of the
105 regression coefficients revealed qualitatively different dynamics, in which the wait times were
106 explained by the reward offered on the current trial, but the trial initiation times reflected an
107 exponentially weighted effect of previous rewards, consistent with a model-free temporal dif-
108 ference learning rule (Fig. 2c,d). We fit exponential curves to the previous trial coefficients
109 for each rat, and found that the distributions of exponential decay time constant parameters (τ)
110 were significantly different for the trial initiation and wait times ($p \ll 0.01$, Wilcoxon sign-
111 rank test, $N = 291$; Fig. 2e). Moreover, τ parameters were not correlated across models ($r =$
112 0.08 , $p = 0.18$, Pearson linear correlation, $N = 291$, Fig. 2f).

113 To leverage individual variability across rats, we compared rats with fast and slow temporal

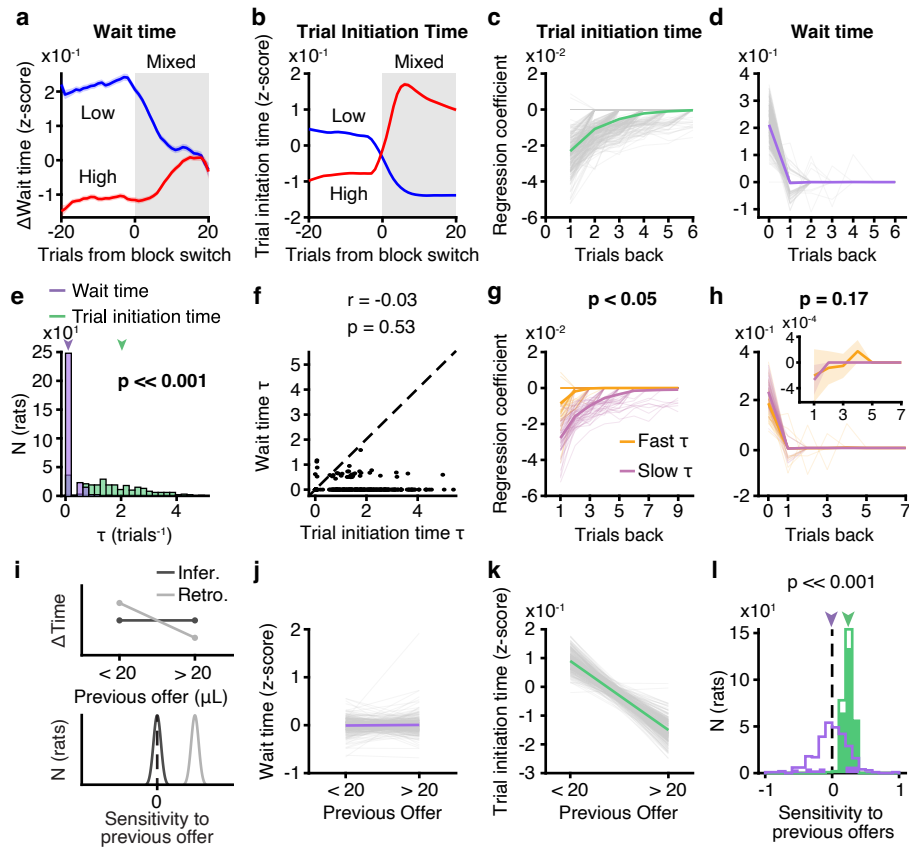


Figure 2: Wait and trial initiation times use distinct estimates of the value of the environment. **a-b.** Mean change in wait times (a) and trial initiation times (b) from low or high blocks to mixed blocks, $N = 291$. Data are mean \pm S.E.M. Data were smoothed with a moving window of 10 trials. **c-d.** Regression coefficients for (c) trial initiation time and (d) wait time. **e-f.** Time constants, τ , of exponential decay parameters fit to previous trial coefficients for wait time (purple) and trial initiation time (green) were (e) significantly different, $p << 0.001$, Wilcoxon sign-rank test, $N = 291$, and (f) uncorrelated, $r = -0.03$, $p = 0.53$, Pearson linear correlation, $N = 291$. **g-h.** Fast or slow initiation time τ (<20th or >80th percentile) meaningfully divided rats based on their initiation time regression coefficients (g; $p << 0.01$, one-tailed permutation test, $N = 116$), but not wait time coefficients (h; $p = 0.1$, one-tailed permutation test, $N = 116$). Inset shows previous trial coefficients for wait times with adjusted y-axis limits. **i.** Predictions for sensitivity to previous offers (behavior conditioned on previous offer <20 μ L - >20 μ L) for fixed (light) versus sequentially-updated (dark) estimates of environmental value, consistent with inferential and retrospective strategies, respectively. **j.** Wait time on 20 μ L catch trials in mixed blocks conditioned on previous reward offer. ($p < 0.05$ for 38/291 rats, Wilcoxon rank-sum test). **k.** Trial initiation time in mixed blocks conditioned on previous reward offer. ($p < 0.05$ for 256/291 rats, Wilcoxon rank-sum test). **l.** Sensitivity to previous offers for wait time (purple) and trial initiation time (green). $p << 0.001$, Wilcoxon sign-rank test, $N = 291$. Colored bars are individual rats with $p < 0.05$, Wilcoxon rank-sum test.

114 integration for trial initiation times (τ from exponential fit to regression coefficients < 20 th or
115 > 80 th percentiles). There were differences in temporal integration for trial initiation times,
116 but not wait times, for these groups (Fig. 2g-h, trial initiation time $p \ll 0.001$, wait time $p =$
117 0.5 , permutation test, $N = 116$). Collectively, these data suggest that within a block, wait times
118 use a fixed estimate of the value of the environment, whereas trial initiation times are sensitive
119 to previous rewards (Fig. 2c,d). Indeed, for almost all rats (87%), wait times for 20 μ L offers
120 in mixed blocks were not significantly different if they were preceded by rewards that were
121 smaller or larger than 20 μ L ($p > 0.05$, Wilcoxon rank-sum test, $N = 253/291$). However, for
122 89% of rats, trial initiation times were significantly modulated by previous rewards, suggesting
123 fixed and incrementally updated estimates of the value of the environment, respectively ($p <$
124 0.05 , Wilcoxon rank-sum test, $N = 256/291$, Fig. 2i-l).

125 One factor that could in principle influence initiation times is satiety. However, satiety
126 effects were small, and modestly apparent as a gradual increase in trial initiation times over the
127 session. To control for these effects, we regressed out initiation times against trial number in
128 order to detrend the slow changes over the course of a session. However, there was no qualitative
129 change in any of our results, including the dynamics at block transitions, if we did not detrend
130 (Extended Data Fig. 6). Because the trials are self-initiated, we suspect that when rats are sated
131 they choose not to initiate trials, thereby minimizing the effects of satiety on behavior, at least
132 compared to other factors that contribute more to the variance in initiation times (e.g., reward
133 history, Fig. 2c).

134 **Computational modeling reveals distinct value computations for sequential** 135 **decisions.**

136 Our data suggest that rats' sequential decisions (when to initiate trials and how long to wait
137 for rewards) reflect different value computations. We developed behavioral models for wait and

138 trial initiation times, inspired by foraging theories¹⁹. The wait time model implemented a trial
139 value function that scaled with the offered reward and decayed to reflect reward probability over
140 time¹⁶. The model's predicted wait time was when the value function fell below the value of
141 the environment (opportunity cost) on each trial (Fig. 3a). This model captured key features of
142 the rats' behavior, including the monotonic relationship between wait time and reward offer in
143 mixed blocks (Fig. 3b) as well as a graded dependence of wait times on the catch probability
144 (Fig. 3c). Different versions of the model estimated the value of the environment using different
145 computations.

146 Analysis of rats' trial initiation times suggests that they estimate the value of the environ-
147 ment as a running average of rewards (Fig. 2c)^{2,17,30}. We refer to this computation as retrospec-
148 tive, as it reflects past experience³¹. Alternatively, rats' wait times reflected the use of discrete
149 estimates of block value (Fig. 2a,d,j). Therefore, rats might infer the current block³¹⁻³⁵, and
150 use fixed estimates of block value based on that inference. We refer to this computation as
151 inferential, since it requires hidden state inference.

152 The inferential model selected the most likely block using Bayes' Rule with a prior that
153 incorporated reward history and knowledge of the block transition structure. This model reca-
154 pitulated the rats' wait times converging to a common value in mixed blocks (Fig. 3d-e). Across
155 rats, the model also captured that wait times for 20 μ L in mixed blocks were not sensitive to pre-
156 vious rewards (Fig 2j; Fig. 3f). This reflects the model's use of a fixed estimate of the value of
157 the environment in each block.

158 In the retrospective case, the value of the environment was estimated as a recency-weighted
159 average of offered rewards according to a temporal-difference learning rule (Fig. 3g). A static
160 learning rate was unable to capture the rats' behavior (Extended Data Fig. 7). Previous work has
161 shown that animals adjust their learning rates depending on the volatility in the environment,
162 since it is advantageous to learn faster in dynamic environments³⁶⁻³⁸. Therefore, our model

163 scaled the learning rate by the trial-by-trial change in the inferential model's beliefs about the
164 hidden state (derivative of the posterior, see Methods).

165 The retrospective model captured several key features of rats' trial initiation times, which
166 we modeled as inversely proportional to the value of the environment² (Fig. 2g-i). First, with
167 a sufficiently small learning rate (<0.1 , Fig. 2g), the model integrated reward history on long
168 timescales such that trial initiation times in mixed blocks depended on the previous block iden-
169 tity. Second, the dynamic learning rate captured the rapid behavioral dynamics at block transi-
170 tions. Finally, integrating over previous trials captured the dependence of trial initiation times
171 on previous rewards in mixed blocks (Fig. 2k, Fig. 3i). We explored versions of the dynamic
172 learning rate that did not reflect inference, including using the unsigned reward prediction error
173 or a running average of reward prediction errors³⁸. However, these models could not capture
174 both short and long timescale dynamics at block transitions (Extended Data Fig. 7). This sug-
175 gests that trial initiation times reflect a retrospective computation that is influenced by subjective
176 belief distributions^{36,37}. In other words, while trial initiation and wait times reflect distinct value
177 computations, those computations interact when states are uncertain via a dynamic learning
178 rate.

179 We fit the retrospective and inferential models to rats' wait times. By several model com-
180 parison metrics, wait times were better fit by the inferential model that used hidden state infer-
181 ence to select block-specific estimates of the value of the environment ($p \ll 0.001$, Wilcoxon
182 signed-rank test, $N = 291$; Fig. 3j, Extended Data Fig. 8), consistent with that model repro-
183 ducing the wait time dynamics (Fig. 2a,3d). We also used the model to identify trials in mixed
184 blocks where the rats were likely to make mistaken inferences. The rats' wait times reflected
185 these mistaken inferences, further indicating that their wait times were well-described by the
186 inferential model ($p \ll 0.001$, Wilcoxon signed-rank test comparing wait times for 20 μ L in
187 misinferred high vs. low blocks, $N = 291$; Extended Data Fig. 9).

188 We also developed a “belief state” model that estimated the value of the environment as the
189 sum of block-specific values weighted by their posterior probabilities. The inferential and belief
190 state models make qualitatively similar predictions about the average wait times. In fact, when
191 the posterior beliefs are stable, which is often the case, the belief state and inferential models
192 are identical, and model comparison did not favor one model over the other (Extended Data Fig.
193 8).

194 To leverage individual differences, we turned to the inferential model of wait times. We
195 added a parameter, λ , that controlled the extent to which the model used an optimal prior, $\lambda =$
196 1, versus an uninformative prior, $\lambda = 0$ (Fig. 3k; Extended Data Fig. 10). We divided the rats
197 into groups with low or high values of λ ($\lambda < 20$ th or > 80 th percentiles; Extended Data Fig.
198 11), and compared the parameters of logistic functions fit to the average wait time dynamics for
199 these groups. Rats with optimal and poor inference exhibited significantly different dynamics
200 at transitions from mixed into low or high blocks, indicated by different inverse temperature
201 parameters (mix to low/high, $p < 0.05$, one-tailed permutation test, $N = 180$ Fig. 3l). There was
202 no difference in the dynamics of trial initiation times for those same groups of rats (mixed to
203 low: $p = 0.3$, mixed to high: $p = 0.2$, one-tailed permutation test, $N = 180$; Fig. 3l). Therefore,
204 individual differences in trial initiation (Fig. 2g,h) and wait times (Fig. 3l) are dissociable.

205 **Block sensitivity for wait times requires structure learning.**

206 Structure learning is the process of learning the hidden structure of environments, including
207 latent states and transition probabilities between them³⁹. If wait and trial initiation times dif-
208 ferentially required knowledge of latent task structure, they should exhibit different dynamics
209 over training. In the final stage of training, when rats were introduced to the hidden states, their
210 wait times for 20 μ L gradually became sensitive to the reward block (Fig. 4a). We observed
211 a gradual increase in the magnitude of reward and block regression coefficients that mirrored

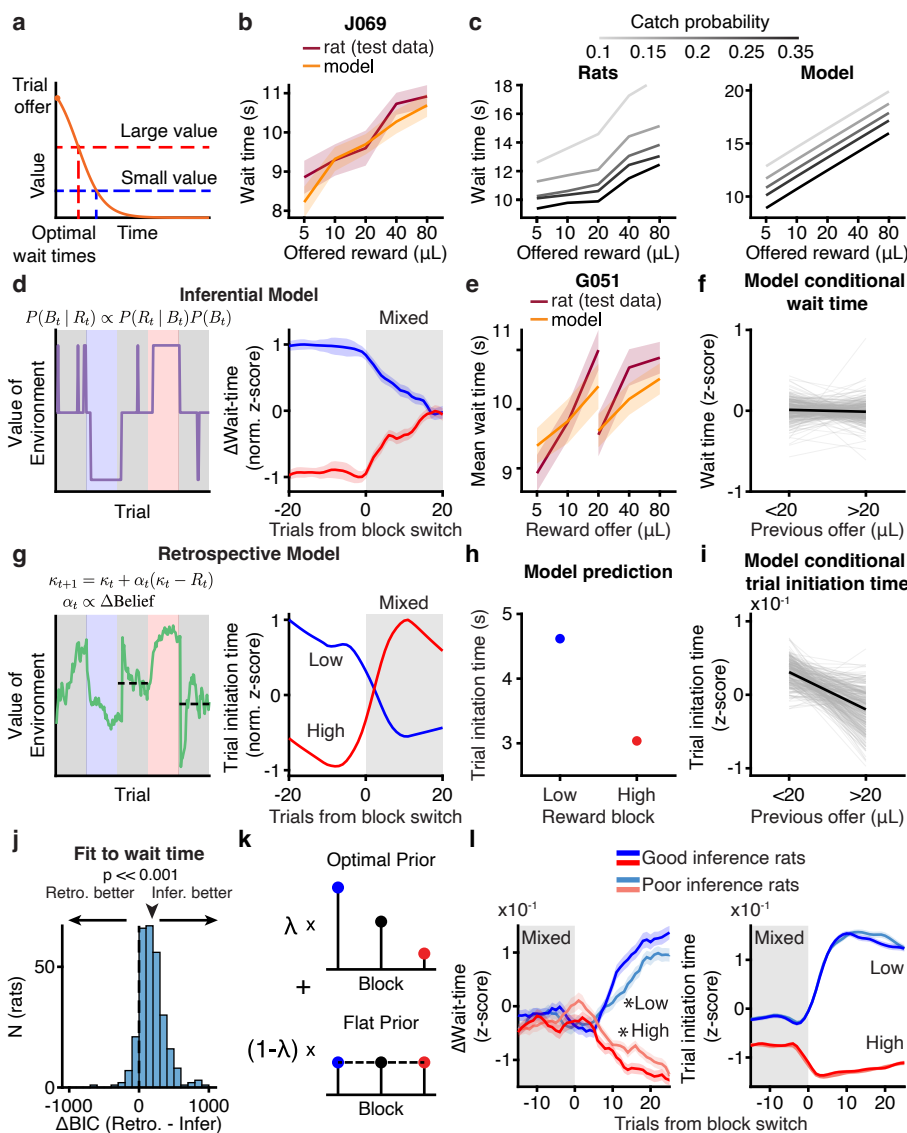


Figure 3: Computational modeling reveals distinct value computations for wait time and trial initiation

a. Model schematic. **b.** Example wait time model performance for mixed blocks only in held-out test data. **c.** Rat population (left; sample sizes in methods) and model (right) wait time data in mixed blocks as a function of catch probability. **d.** Example opportunity cost and wait time dynamics from inferential model. **e.** Inferential model fit to rats can capture wait time behavior in held-out test data. **f.** Inferential model captures conditional wait time trend across rats ($N = 291$). **g.** Example opportunity cost and wait time dynamics from retrospective model. **h.** Retrospective model can qualitatively capture trial initiation time behavior. **i.** Retrospective model captures conditional trial initiation time trend across rats ($N = 291$). **j.** Model comparison using ΔBIC prefers inferential model compared to retrospective model when fit to wait time data ($p \ll 0.001$, Wilcoxon Signed-rank test, $N = 291$).

Figure 3 cont.: **Computational modeling reveals distinct value computations for wait time and trial initiation k.** Schematic for sub-optimal inference model **l**. Transitions from mixed to low (blue) or high (red) blocks for wait time (left) or trial initiation time (right) separated by quality of inference ($\lambda < 20$ th or > 80 th percentile). * $p < 0.05$, one-tailed non-parametric shuffle test comparing logistic fit parameters, $N = 116$. Data are mean \pm S.E.M.

212 the behavioral sensitivity to hidden states (Fig. 4b). In contrast, trial initiation times exhibited
213 block sensitivity on the first session in the final training stage (Fig. 4a). This sensitivity was
214 comparable early and late in training, consistent with animals using previous rewards to a simi-
215 lar extent at these timepoints (Fig. 4c). These data suggest that block sensitivity for wait times,
216 but not trial initiation times, required learned knowledge of hidden task states, and that these
217 decisions reflected computations with distinct learning dynamics.

218 The modest increase in trial initiation time block sensitivity over training is consistent with
219 the gradual use of a dynamic learning rate that reflected learned knowledge of the blocks. A
220 hallmark of the dynamic learning rate was the “overshoot” after transitions from high to mixed
221 blocks (difference between maximum trial initiation time after transitioning and the trial initia-
222 tion time 20 trials post-transition; Fig. 2b). The overshoot became more prominent with training
223 (Fig. 4d), on a similar timescale as block sensitivity for wait times (Fig. 4e), suggesting a shared
224 mechanism.

225 **Reducing state uncertainty did not change trial initiation times.**

226 Why would animals use a retrospective computation at trial initiation, but rely on an inferen-
227 tial computation as rats deliberated just 1-2 seconds later? In non-human primates, the decision
228 to initiate trials can also reflect retrospectively computed values that differ from the values gov-
229 erning the subsequent choice^{40,41}. One possibility is that motivation and approach behavior rely
230 on neural circuits that do not support inference¹¹. Another possibility is that actions more distal
231 to rewards are more likely to be retrospective, because there are more steps required to men-

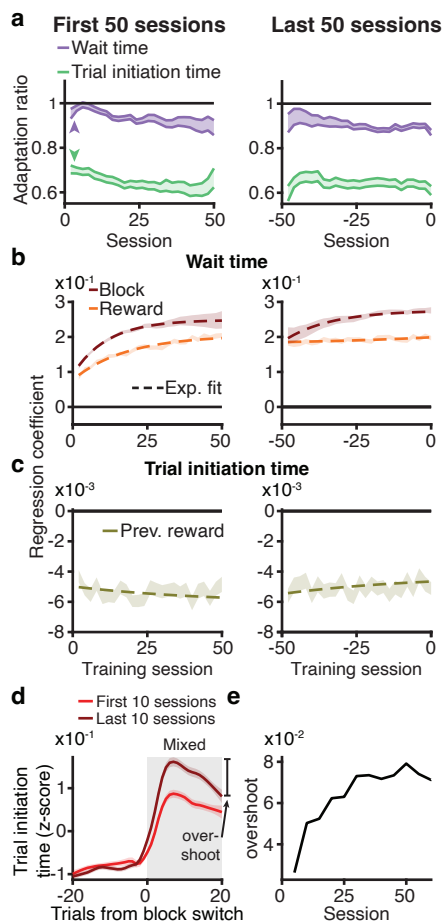


Figure 4: Block sensitivity for wait times requires structure learning. **a.** Wait time adaptation ratio (average wait time for 20 μL in high/low blocks) evolved over training, while trial initiation time ratio (average in high/low blocks) was below 1 on first session. **b.** Linear regression coefficients for block and reward gradually evolved over training for wait time. **c.** Linear regression coefficient for previous reward was relatively stable across training for trial initiation time. **d.** Overshoot in trial initiation time (difference between maximum z-scored trial initiation time and trial initiation time at trial 20 post-transition) was more prominent after structure learning. **e.** Overshoot in trial initiation time dynamics evolved on a similar timescale as block sensitivity for wait times.

232 tally simulate outcomes for forward-looking strategies like planning^{42,43}. According to either
233 hypothesis, the decision of when to initiate a trial is inherently retrospective.

234 Theoretical work in reinforcement learning has suggested that the brain should select the
235 strategy that is the fastest and most accurate when taking into account uncertainty^{8,9}. Therefore,
236 perhaps trial initiation times are retrospective because the rats' subjective beliefs about the
237 inferred state have more uncertainty before they hear the reward offer. Model simulations of a
238 Bayes' optimal observer did show that the reward offer reduced the uncertainty of subjective
239 beliefs about the hidden state (comparing variance of prior to variance of posterior, $p \ll 0.001$,
240 Wilcoxon sign-rank test).

241 To test this hypothesis, we modified the task so that some rats heard the reward cue before
242 they initiated the trial, when the center light turned on; they heard the tone again at trial ini-
243 tiation, as in the standard task (Fig. 5a). Their trial initiation times became sensitive to the
244 offered reward (Fig. 5b). However, trial initiation times for 20 μ L in mixed blocks were still
245 modulated by the previous reward, consistent with the use of incrementally updated estimates
246 of the value of the environment within a block ($p < 0.05$ for 13/16 rats; Fig. 5c). Moreover, how
247 quickly they initiated trials in mixed blocks continued to depend on the previous block identity
248 (Fig. 5d). These data indicate that there may be something inherently retrospective about the
249 motivational decision to initiate a trial.

250 Discussion

251 We used high-throughput training to collect statistically powerful datasets and leverage in-
252 dividual variability across hundreds of animals. Consistent with previous work, rats adjusted
253 their behavior as we varied the richness of the environment in a way consistent with forag-
254 ing theories^{19,30,44-46}, and behavioral economic theories of reference dependence^{47,48}. Notably,
255 we found that animals used multiple, parallel computations to estimate the richness of the en-

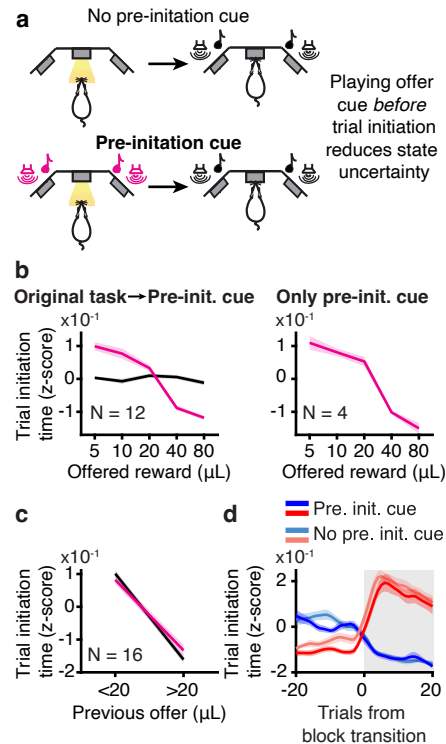


Figure 5: Value computations for motivation do not depend on state uncertainty. **a.** Schematic of pre-initiation cue experiment. **b.** Trial initiation time varied as a function of offered volume for rats that trained on the original task before transitioning to pre-initiation cue task (left) and for rats that trained exclusively on the pre-initiation cue task (right). **c.** Trial initiation times were still sensitive to previous reward (behavior on trials offering 20 μL conditioned on the previous reward offer) after training on the pre-initiation cue task. ($p < 0.05$ for 13/16 rats, Wilcoxon Rank-sum test, $N = 16$). **d.** Trial initiation times in mixed blocks depended on previous block type in pre-initiation cue task.

256 vironment, and rapidly switched between these computations on single trials, indicating that
 257 value computations vary on fine timescales (seconds). Our data are consistent with evidence
 258 for multiple decision-making systems that rely on distinct neural circuits^{10,12,49,50}. While ani-
 259 mals' decisions of how long to wait for rewards relied on hidden state inference, the decision of
 260 when to initiate the trial was governed by a retrospective computation that calculated the value
 261 of the environment as the running average of rewards. Reducing state uncertainty before the
 262 trial did not change the value computations governing trial initiation times, suggesting that this

263 decision may be inherently retrospective, although influenced by subjective belief distributions
264 via a dynamic learning rate.

265 Recent work in psychology, machine learning, and neuroscience has characterized how par-
266 allel value computations might be combined^{4-8,40,51-54}. For instance, in multi-step decision
267 tasks, interaction effects in regression models are thought to reflect the use of combined ret-
268 rospective and inferential value estimates^{14,15}, and hybrid strategies for computing values have
269 been approximated as a weighted average of retrospective and inference-based values⁴⁰. Our
270 findings add to this body of work. Instead of simply combining or averaging values that were
271 computed in different ways, rats seemed to coordinate their dynamics: changes in subjective
272 beliefs about inferred states acted as a gain on retrospective value learning rates. Moreover, we
273 tested the prevailing hypothesis about arbitration between these parallel value computations,
274 namely, that agents should use the value estimate with the lowest uncertainty^{8,9}. We reduced
275 state uncertainty by playing the reward cue before rats initiated trials. However, trial initiation
276 times still reflected retrospective value computations (Fig. 5c-d). We hypothesize that different
277 neural circuits mediate these rapid sequential decisions (starting the trial versus deciding how
278 long to wait), and that these circuits support or favor distinct value computations due to their
279 connectivity and other neurobiological constraints.

280 Alternatively, previous work has suggested that actions more distal to rewards are more
281 likely to be retrospective, because there are more steps required to mentally simulate outcomes
282 for forward-looking strategies like planning^{42,43}. Therefore, one potential reason that trial initi-
283 ation times were retrospective is because they were more distal to rewards. However, in multi-
284 step decision-making tasks (i.e., the two-step task), the first action, which is diagnostic of how
285 value is computed, generally reflects computations that use a model of the world to flexibly esti-
286 mate values^{14,55,56}. Compared to the two-step task, the first action in our task is a similar number
287 of states away from the terminal reward state, but the temporal delays are longer. Therefore, it

288 is possible that temporal proximity to reward may determine how values are computed.

289 It may be counterintuitive that the retrospective computation produced faster dynamics at
290 block transitions than hidden state inference (Fig. 2a,b). Two features of the models explain this
291 observation. First, the inferential model selects the block with the maximum posterior proba-
292 bility. This argmax operation nonlinearly thresholds whether changes in the posterior produce
293 changes in the inferred state. In contrast, the retrospective model’s estimate of the value of the
294 environment is directly influenced by graded, “subthreshold” changes in the posterior via the
295 dynamic learning rate. Subthreshold changes in the posterior necessarily precede changes that
296 cross threshold for inferring a state change. Second, the inferential model’s prior is recursive:
297 the posterior on one trial becomes the prior on the next trial. This means that the prior accu-
298 mulates information over trials to infer state changes, instead of making them instantaneously.
299 Indeed, individual differences in the informativeness of rats’ priors predicted the dynamics of
300 their inferred state changes (Fig. 3l).

301 The contextual effects we observed likely reflect efficient coding of value^{17,57–59}. According
302 to the efficient coding hypothesis, to represent stimuli efficiently, neurons should be tuned to
303 stimulus distributions that animals are most likely to encounter in the world⁶⁰. Recent stud-
304 ies have shown that biases in value-based decision-making, including the contextual effects
305 observed here, reflect efficient value coding^{17,57,58}. Previous studies examined how neurons
306 “adapted” to reward or stimulus distributions over blocks of trials or sessions, implying grad-
307 ual, experience-dependent adjustments in behavioral sensitivity and neural tuning^{17,61,62}. Our
308 findings suggest that if animals have learned the reward or stimulus distributions associated
309 with a particular state, they can condition their subjective value representations on that inferred
310 state, perhaps via discrete, state-dependent adjustments in neural sensitivity⁶³. A major future
311 question is how multi-regional neural circuits represent belief distributions for hidden state in-
312 ference, and condition rapid adjustments in efficient neural representations of value on inferred

313 states.

314 **Methods**

315 **Subjects**

316 A total of 291 Long-evans rats (184 male, 107 female) between the ages of 6 and 24 months
317 were used for this study (*Rattus norvegicus*). The Long-evans cohort also included ADORA2A-
318 Cre ($N=10$), ChAT-Cre ($N=2$), DRD1-Cre ($N=3$), and TH-Cre ($N=12$). Animal use procedures
319 were approved by the New York University Animal Welfare Committee (UAWC #2021-1120)
320 and carried out in accordance with National Institutes of Health standards.

321 Rats were pair housed when possible, but were occasionally single housed (e.g. if fighting
322 occurred between cagemates). Animals were water restricted to motivate them to perform be-
323 havioral trials. From Monday to Friday, they obtained water during behavioral training sessions,
324 which were typically 90 minutes per day, and a subsequent ad libitum period of 20 minutes.
325 Following training on Friday until mid-day Sunday, they received ad libitum water. Rats were
326 weighed daily.

327 **Behavioral training**

328 Rats were trained in a high-throughput behavioral facility in the Constantinople lab using
329 a computerized training protocol. They were trained in custom operant training boxes with
330 three nose ports. Each nose port was 3-D printed, and the face was protected with an epoxied
331 stainless steel washer (McMaster-Carr #92141A056). All ports contained a visible light emit-
332 ting diode (LED; Digikey #160-1850-ND), and an infrared LED (Digikey #365-1042-ND) and
333 infrared photodetector (Digikey #365-1615-ND) that enabled detection of when a rat broke the
334 infrared beam with its nose. Additionally, the side ports contained stainless steel lick tubes
335 (McMaster-Carr #8988K35, cut to 1.5mm) that delivered water via solenoid valves (Lee Com-
336 pany #LHDA1231115H). There was a speaker mounted above each side port that enabled de-

337 livery of stereo sounds (Bohlender Graebener). The behavioral task was instantiated as a finite
338 state machine on an Arduino-based behavioral system with a Matlab interface (Bpod State Ma-
339 chine r2, Sanworks), and sounds were delivered using a low-latency analog output module
340 (Analog Output Module 4ch, Sanworks) and stereo amplifier (Lepai LP-2020TI).

341 Research technicians loaded rats in and out of the training rigs in each session, but the train-
342 ing itself was computer automated. All rig computers automatically pulled version-controlled
343 software from a git repository and wrote behavioral data to a MySQL (MariaDB) database
344 hosted on a synology server. Rig computers automatically loaded each rat's training settings
345 file from the previous session, and following training, wrote a new settings file to the server
346 for the subsequent day of training. Rig computers automatically loaded files for specific rats
347 based on a schedule on the MySQL database. Human intervention was possible but generally
348 unnecessary.

349 **Sound Calibration**

350 We calibrated sounds using a hand-held Precision Sound Level Meter with a 1/2" micro-
351 phone (Bruel & Kjaer, Type 2250). The microphone was calibrated with a sound level calibrator
352 (Bruel & Kjaer, Type 4230). Tones of different frequencies (1, 2, 4, 8, 16kHz) were presented
353 for 10 seconds each; these tones were selected because they are in the trough of the behavioral
354 audiogram for rats⁶⁴. They are also on a logarithmic scale and thus should be equally discrim-
355 inable to the animals. We adjusted the auditory gain in software for each frequency stimulus to
356 match the sound pressure level to 70dB in the rig, measured when the microphone was proximal
357 to the center poke.

358 **Task Logic**

359 LED illumination from the center port indicated that the animal could initiate a trial by
360 poking its nose in that port - upon trial initiation the center LED turned off. While in the center

361 port, rats needed to maintain center fixation for a duration drawn uniformly from [0.8, 1.2]
362 seconds. During the fixation period, a tone played from both speakers, the frequency of which
363 indicated the volume of the offered water reward for that trial [1, 2, 4, 8, 16kHz, indicating
364 5, 10, 20, 40, 80 μ L rewards]. Following the fixation period, one of the two side LEDs was
365 illuminated, indicating that the reward might be delivered at that port; the side was randomly
366 chosen on each trial. This event (side LED ON) also initiated a variable and unpredictable delay
367 period, which was randomly drawn from an exponential distribution with mean = 2.5 seconds.
368 The reward port LED remained illuminated for the duration of the delay period, and rats were
369 not required to maintain fixation during this period, although they tended to fixate in the reward
370 port. When reward was available, the reward port LED turned off, and rats could collect the
371 offered reward by nose poking in that port. The rat could also choose to terminate the trial
372 (opt-out) at any time by nose poking in the opposite, un-illuminated side port, after which a
373 new trial would immediately begin. On a proportion of trials (15-25%), the delay period would
374 only end if the rat opted out (catch trials). If rats did not opt-out within 100s on catch trials, the
375 trial would terminate.

376 The trials were self-paced: after receiving their reward or opting out, rats were free to
377 initiate another trial immediately. However, if rats terminated center fixation prematurely, they
378 were penalized with a white noise sound and a time out penalty (typically 2 seconds, although
379 adjusted to individual animals). Following premature fixation breaks, the rats received the same
380 offered reward, in order to disincentivize premature terminations for small volume offers.

381 We introduced semi-observable, hidden-states in the task by including uncued blocks of
382 trials with varying reward statistics¹⁷: high and low blocks, which offered the highest three
383 or lowest three rewards, respectively, and were interspersed with mixed blocks, which offered
384 all volumes. There was a hierarchical structure to the blocks, such that high and low blocks
385 alternated after mixed blocks (e.g., mixed-high-mixed-low, or mixed-low-mixed-high). The first

386 block of each session was a mixed block. Blocks transitioned after 40 successfully completed
387 trials. Because rats prematurely broke fixation on a subset of trials, in practice, block durations
388 were variable.

389 **Criteria for including behavioral data**

390 In this task, the rats were required to reveal their subjective value of different reward of-
391 fers. To determine when rats were sufficiently trained to understand the mapping between the
392 auditory cues and water rewards, we evaluated their wait time on catch trials as a function of
393 offered rewards. For each training session, we first removed wait times that were greater than
394 two standard deviations above the mean wait time on catch trials in order to remove potential
395 lapses in attention during the delay period (this threshold was only applied to single sessions
396 to determine whether to include them). Next, we regressed wait time against offered reward
397 and included sessions with significantly positive slopes that immediately preceded at least one
398 other session with a positive slope as well. Once performance surpassed this threshold, it was
399 typically stable across months. Occasional days with poor performance, which often reflected
400 hardware malfunctions or other anomalies, were excluded from analysis. We emphasize that the
401 criteria for including sessions in analysis did not evaluate rats' sensitivity to the reward blocks.
402 Additionally, we excluded trial initiation times above the 99th percentile of the rat's cumulative
403 trial initiation time distribution pooled over sessions.

404 **Shaping**

405 The shaping procedure was divided into 8 stages. For stage 1, rats learned to maintain a
406 nose poke in the center port, after which a 20 μL reward volume was delivered from a random
407 illuminated side port with no delay. Initially, rats needed to maintain a 5 ms center poke. The
408 center poke time was incremented by 1 ms following each successful trial until the center poke
409 time reached 1 s, after which the rat moved to stage 2.

410 Stages 2-5 progressively introduced the full set of reward volumes and corresponding au-
411 ditory cues. Rats continued to receive deterministic rewards with no delay after maintaining a
412 1 second center poke. Each stage added one additional reward that could be selected on each
413 trial- stage 2 added 40 μL , stage 3 added 5 μL , stage 4 added 80 μL , and stage 5 added 10 μL .
414 Each stage progressed after 400 successfully completed trials. All subsequent stages used all 5
415 reward volumes.

416 Stage 6 introduced variable center poke times, uniformly drawn from [0.8-1.2] s. Addition-
417 ally, stage 6 introduced deterministic reward delays. Initially, rewards were delivered after a 0.1
418 s delay, which was incremented by 2 ms after each successful trial. After the rat reached delays
419 between 0.5 and 0.8 s, the reward delay was incremented by 5 ms following successful trials.
420 Delays between 0.8 and 1 s were incremented by 10 ms, and delays between 1 and 1.5 s were
421 incremented by 25 ms. Rats progressed to stage 7 after reaching a reward delay of 1.5 s.

422 In stage 7, rats experienced variable delays, drawn from an exponential distribution with
423 mean of 2.5 seconds. Additionally, we introduced catch trials (see above), with a catch proba-
424 bility of 15%. Stage 7 terminated after 250 successfully completed trials.

425 Finally, stage 8 introduced the block structure (see above). We additionally increased the
426 catch probably for the first 1000 trials to 35%, to encourage the rats to learn that they could
427 opt-out of the trial. After 1000 completed trials, the catch probability was reduced to 15-20%.
428 All data in this paper was from training stage 8.

429 **Training for male and female rats**

430 We collected data from both male and female rats (160 male, 114 female). Male and female
431 rats were trained in identical behavioral rigs with the same shaping procedure described above.
432 Early cohorts of female rats experienced the same reward set as the males. However, female
433 rats are smaller, and they consumed less water and performed substantially fewer trials than

Stage	Center poke time	5 μ L	10 μ L	20 μ L	40 μ L	80 μ L	Reward delay	Reward probability	Blocks
1	Increment to 1s			X			0	1	
2	1s			X	X		0	1	
3	1s	X		X	X		0	1	
4	1s	X		X	X	X	0	1	
5	1s	X	X	X	X	X	0	1	
6	Variable (0.8-1.2s)	X	X	X	X	X	Increment to 1.5s	1	
7	Variable (0.8-1.2s)	X	X	X	X	X	Variable (from exponential)	0.85	
8	Variable (0.8-1.2s)	X	X	X	X	X	Variable (from exponential)	0.65-0.85	X

434 the males. Therefore, to obtain sufficient behavioral trials from them, reward offers for female
435 rats were slightly reduced while maintaining the logarithmic spacing: [4, 8, 16, 32, 64 μ L]. For
436 behavioral analysis, reward volumes were treated as equivalent to the corresponding volume
437 for the male rats (e.g., 16 μ L trials for female rats were treated the same as 20 μ L trials for
438 male rats). The auditory tones were identical to those used for male rats. We did not observe
439 any significant differences between the male and female rats, in terms of contextual effects
440 (Extended Data Fig. 12), or behavioral dynamics at block transitions (data not shown).

441 We tracked most female rats' stages in the estrus cycle using vaginal cytology, with vaginal
442 swabs collected immediately after each session using a cotton-tipped applicator first dipped
443 in saline. Samples were smeared onto a clean glass slide and visually classified under a light
444 microscope. For the current study, data from female rats was averaged across all stages of the
445 estrus cycle.

446 Behavioral models

447 We developed separate behavioral models to describe rat's wait time and trial initiation time
448 data. Both wait time and trial initiation time should depend on the value of the environment.
449 For the wait time data, we adapted a model from¹⁶ which described the wait time, WT, in terms
450 of the value of the environment (i.e., the opportunity cost), the delay distribution, and the catch
451 probability (i.e., the probability of the trial being unrewarded). Given an exponential delay
452 distribution, we defined the predicted wait time as

$$\text{WT} = D\tau \log \left(\frac{C}{1-C} \cdot \frac{R - \kappa\tau}{\kappa\tau} \right).$$

453 where τ is the time constant of the exponential delay distribution, C is the probability of reward
454 (1-catch probability), R is the reward on that trial, κ is the opportunity cost, and D is a scaling
455 parameter. In the context of optimal foraging theory and the marginal value theorem, which
456 provided the theoretical foundation for this model, each trial is a depleting “patch” whose value
457 decreases as the rat waits¹⁹. Within a patch, the decision to leave depends on the overall value
458 of the environment, κ , which is stable within trials but can vary across trials and hidden reward
459 states, i.e., blocks.

460 The above equation was shown to be normative for a Markov Decision Process in which
461 the value of the environment was constant for the foreseeable future¹⁶. However, given that
462 the value of the environment changed over blocks in our task, it is possible that this equation
463 is not normative for our case. However, this formulation qualitatively captured features of the
464 data, including the graded dependence of wait times on the catch probability, and sensitivity to
465 reward volumes and blocks. Therefore, we found it to be a useful, if not necessarily normative,
466 process model of behavior.

467 For the trial initiation time, we adapted a model² which describes the optimal trial initiation

468 time, TI , given the value of the environment, κ , as

$$\text{TI} = \frac{D}{\kappa},$$

469 where D is a scale parameter.

470 We initially evaluated two different ways of calculating the value of the environment for
471 these models, which are shared between the wait time and trial initiation time models: a retro-
472 spective and inferential model (see below). We assumed independent log-normal noise for each
473 trial, with a constant variance of 8 seconds for the wait time model and 4 seconds for the trial
474 initiation time model. The log-normal noise model outperformed alternative noise models, such
475 as gamma and ex-Gaussian noise. The noise variance terms were selected from a grid search
476 using data from a subset of animals.

477 **Inferential model**

478 The inferential model has three discrete value parameters (κ_{low} , κ_{mixed} , κ_{high}), each associ-
479 ated with a block. For each trial, the model chooses the κ associated with the most probable
480 block given the rat's reward history. Specifically, for each trial, Bayes' Theorem specifies the
481 following:

$$P(B_t | R_t) \propto P(R_t | B_t)P(B_t).$$

482 where B_t is the block on trial t and R_t is the reward on trial t . The likelihood, $P(R_t | B_t)$, is the
483 probability of the reward for each block, for example,

$$P(R_t | B_t = \text{Low}) = \begin{cases} \frac{1}{3}, & \text{if } R_t = 5, 10, 20 \mu\text{L} \\ 0, & \text{if } R_t = 40, 80 \mu\text{L}. \end{cases}$$

484 To calculate the prior over blocks, $P(B_t)$, we marginalize over the previous block and use the
485 previous estimate of the posterior:

$$P(B_t) = \sum_{B_{t-1}} P(B_t | B_{t-1})P(B_{t-1} | R_{t-1}). \quad (\text{Eq. 1})$$

486 $P(B_t | B_{t-1})$, referred to as the “hazard rate,” incorporates knowledge of the task structure,
487 including the block length and block transition probabilities. For example,

$$P(B_t = \text{Low} | B_{t-1}) = \begin{cases} 1 - H_0, & \text{for } B_{t-1} = \text{Low} \\ H_0, & \text{for } B_{t-1} = \text{Mixed} \\ 0, & \text{for } B_{t-1} = \text{High} \end{cases}$$

488 where $H_0 = 1/40$, to reflect the block length. The model assumed a flat block hazard rate for
489 the following reasons. (1) Since animals broke center fixation on a subset of trials, the actual
490 block duration was highly variable. Based on the distributions of experienced block durations,
491 it is unlikely that rats would have learned a perfect step function hazard rate. (2) The blocks
492 spanned several to tens of minutes, making it unlikely that rats would keep a running tally of
493 trials on such long timescales. (3) Gradual changes in wait times at block transitions are not
494 consistent with the use of a veridical step-function hazard rate. (4) We considered an alternative
495 parameterization in which the veridical step function hazard rate was blurred with a Gaussian,
496 but this would have required a number of nontrivial design choices, such as whether the trial
497 counter should be reset after “misinferred” block transitions, regardless of when they occurred
498 in the actual block. (5) Wait times reflected misinferred blocks based on a constant block hazard
499 rate (Extended Data Fig. 9), suggesting that this simplification was a reasonable approximation
500 of the inference process. Including H_0 as an additional free parameter did not improve the
501 performance of the wait time model evaluated on held-out test data in a subset of rats (data not
502 shown), so H_0 was treated as a constant term.

503 The model selected a fixed value of the environment associated with the most likely block.
504 This formulation is related to an established approximation for solving partially-observable
505 Markov decision processes (POMDPs) known as the Most Likely State algorithm⁶⁵. This al-
506 gorithm is well-studied, has precedence in the literature as a heuristic approximation for the
507 full posterior distribution over states, and may be biologically plausible as it is computationally
508 tractable compared to more complex solutions to POMDPs.

509 **Belief state model**

510 Like the inferential model (above), the belief state model has three distinct value parameters
511 and calculates the probability of being in each block using Bayes Rule. However, rather than
512 selecting a single value associated with the most probable block, the model uses the sum of
513 each value, weighted by that probability, that is,

$$\kappa_t = \sum_{B_t} P(B_t | R_t) \kappa_{B_t}.$$

514 While this model uses the full posterior distribution over states, model comparison found that
515 it was comparable to the simpler Most Likely State algorithm (above; Extended Data Fig. 8).
516 In fact, when the belief distributions were stable (e.g., in adaptation blocks), these models were
517 identical. For that reason, we exclusively used the Most Likely State model (above) for this
518 paper.

519 There may be alternative normative strategies for this task given different sets of assump-
520 tions. For instance, assuming an infinite time horizon, one might compute the average kappa
521 under the Markov process determining block transitions, starting from the current state. With a
522 sufficiently long time horizon, this average will be dominated by the steady-state distribution of
523 the Markov process, which would predict no contextual modulation of wait times. Given that
524 the rats exhibited strong contextual effects, this strategy is not consistent with their behavior.
525 We therefore did not explore such a model in the current manuscript.

526 **Inferential model with lambda parameter**

527 To account for potentially sub-optimal inference across rats, we developed a second in-
528 ferential model. This model also uses Bayes rule to calculate the block probabilities, except
529 with a sub-optimal prior, $\text{Prior}_{\text{subopt}}$. Specifically, we introduce a parameter, λ , that generates
530 the sub-optimal prior by weighting between the true, optimal prior ($P(B_t)$, Eq. 1), and a flat,

531 uninformative prior ($\text{Prior}_{\text{flat}}$, uniformly $1/3$), that is,

$$\text{Prior}_{\text{subopt}} = \lambda P(B_t) + (1 - \lambda) \text{Prior}_{\text{flat}}.$$

532 When $\lambda = 1$, this model reduces to the optimal inferential model, and when $\lambda = 0$, this model
533 uses a flat prior and the block probabilities are driven by the likelihood.

534 **Retrospective model**

535 The retrospective model has a single, trial-varying κ variable which represents the recency-
536 weighted average of all previous rewards. This average depends on the learning rate parameter
537 α with the recursive equation

$$\kappa_{t+1} = \kappa_t + \alpha_t \delta_t,$$

538 where κ_t is the value of the environment on trial t , r_t is the reward on trial t , $\delta_t = r_t - \kappa_t$ is the
539 reward prediction error (RPE), and α_t is a dynamic learning rate given by $\alpha_t = G \cdot \alpha_0$. In order
540 to capture the dynamics of the trial initiation times around block transitions, we included a gain
541 term, G_t on the learning rate, which is inversely related to the trial-by-trial change in the mixed
542 block probability from by the inferential model, given by

$$G_t = \frac{1}{1 - |P(B_t = \text{Mixed}|R_t) - P(B_{t-1} = \text{Mixed}|R_{t-1})|}.$$

543 We used trial-by-trial changes in the mixed block probability as a summary statistic of changes
544 in the full posterior distribution. Given the distribution of rewards and the transition structure
545 between blocks, there is always some ambiguity about whether the hidden state is a mixed
546 block, and the posterior block probabilities sum to one. Therefore, changes in the mixed block
547 probability reflect changes in the full posterior on every trial.

548 The dynamic learning rate we implemented is consistent with previous work showing that
549 humans and animals can adjust their learning rates depending on the volatility and uncertainty
550 in the environment^{36–38}. Other models using either (1) a single, static learning rate ($G = 1$), or

551 (2) a dynamic learning rate where the gain term was the unsigned reward prediction error on
552 that trial ($G = |\delta_t|$) were unable to capture the observed trial initiation time dynamics at block
553 transitions (Extended Data Fig. 7).

554 **Fitting and evaluating models**

555 We used MATLAB's constrained minimization function, `fmincon`, to minimize the sum of
556 the negative log likelihoods with respect to the model parameters. 100 random seeds were used
557 in the maximum likelihood search for each rat; parameter values with the maximum likelihood
558 of these seeds were deemed the best fit parameters. Before fitting to rat's data, we confirmed
559 that our fitting procedure was able to recover generative parameters (Extended Data Fig. 13).
560 When evaluating model performance fit to rat data, we performed 5-fold cross-validation and
561 evaluated the predictive power of the model on the held-out test sets. To compare the different
562 models, we used Bayesian Information Criterion (BIC), $BIC = \log(n) \cdot k + 2 \cdot nLL$, where n
563 is the number of trials, k is the number of parameters, and nLL is the negative log-likelihood of
564 the best-fit model evaluated on all data. We confirmed the model comparison by also comparing
565 Akaike Information Criterion ($AIC = 2 \cdot k + 2 \cdot nLL$ where k is the number of parameters and
566 nLL is the negative log-likelihood of the best-fit model evaluated on all data) and cross-validated
567 negative log-likelihood, which gave similar results to BIC.

568 We only fit models to the rats' wait time data. This is because the distribution of trial
569 initiation times was generally heavy-tailed, and seemed to reflect multiple processes on different
570 interacting timescales (e.g., reward sensitivity on short timescales, attention, motivation, and
571 satiety on longer timescales). These processes made it challenging to fit the data with a single
572 process model. Therefore, we used the inferential and retrospective trial initiation time models
573 to generate qualitative predictions that we could compare to the rats' data.

574 **Statistical analyses**

575 **Wait time and trial initiation times: sensitivity to reward blocks**

576 For all analyses, we removed wait times that were one standard deviation above the pooled-
577 session mean. Without thresholding, the contextual effects are qualitatively similar. Outlier wait
578 times tend to occur in low blocks, likely due to attentional or motivational lapses. Therefore, the
579 main difference is that the wait time curves in low blocks are both flatter and longer compared
580 to the thresholded data (Extended Data Fig. 14). When assessing whether a rat's wait time
581 differed by blocks, we compared each rat's wait time on catch trials offering 20 μ L in high
582 and low blocks using a non-parametric Wilcoxon rank-sum test, given that the wait times are
583 roughly log-normally distributed. We defined each rat's wait time ratio as the average wait
584 time on 20 μ L catch trials in high blocks/low blocks. For trial initiation times, we compared all
585 trial initiation times for each block, again using a non-parametric Wilcoxon rank-sum test. We
586 defined each rat's trial initiation time ratio as the average trial initiation time in high blocks/low
587 blocks.

588 Trial initiation times were bimodally distributed, with the different modes reflecting whether
589 previous trials were rewarded or not. Unrewarded trials included opt-out trials and trials where
590 rats prematurely terminated center fixation ("violation trials"). Analyzing these trial types sep-
591 arately showed that trial initiation times following unrewarded trials were modulated by blocks
592 in a similar pattern as the wait times, with rats initiating trials more quickly in high compared
593 to low blocks (Extended Data Fig. 2). While we used all behavioral trials for analyses of
594 trial initiation times throughout the manuscript, we note that trial initiation times following
595 rewarded trials exhibited a different pattern (Extended Data Fig. 2), consistent with previous
596 studies showing that response outcomes gate behavioral strategies^{27,28}. Specifically, following
597 rewarded trials, there was a weak positive correlation between reward magnitude and trial ini-
598 tiation time, in contrast to the strong negative correlation we observed following unrewarded

599 trials. We interpret the positive correlation as potentially reflecting micro-satiety effects. How-
600 ever, as these effects were weak, most of the variance in the trial initiation times were driven by
601 those following unrewarded trials.

602 To assess block effects across the population, we first z-scored each rat's wait time on all
603 catch trials and trial initiation time on all trials. For wait times, we computed the average z-
604 scored wait time on catch trials offering 20 μ L in high and low blocks for each rat, and compared
605 across the population using a paired Wilcoxon sign-rank test. Similarly for trial initiation times,
606 we averaged all z-scored trial initiation times for high and low blocks for each rat, and compared
607 across the population using a paired Wilcoxon sign-rank test.

608 To assess the effects of catch probability on wait times, we trained cohorts of rats with
609 different catch probabilities. The cohorts varied in size: $N = [3, 183, 61, 151, 39]$ for catch
610 probability = [0.1, 0.15, 0.2, 0.25, 0.35], respectively.

611 **Block transition dynamics**

612 To examine behavioral dynamics around block transitions, for each rat, we first z-scored
613 wait-times for catch trials of each volume separately in order to control for reward volume
614 effects. We then computed the difference in z-scored wait times for each volume, relative to the
615 average z-scored wait time for that volume, in each time bin (trial relative to block transition),
616 before averaging the differences over all volumes (Δ z-scored wait time). For trial initiation
617 times, we z-scored all trial initiation times. In order to remove satiety effects, for each session
618 individually, we regressed trial initiation time against z-scored trial number and subtracted the
619 fit.

620 For each transition type, we averaged the Δ z-scored wait times and trial initiation times
621 based on their distance from a block transition, including violation trials (e.g., averaged all wait
622 times four trials before a block transition). Finally, for each block transition type, we smoothed

623 the average curve for each rat using a 10-point moving average, before averaging over rats.

624 When comparing block transition dynamics in rats with different quality priors, specifically
625 from mixed blocks to high or low, we chose rats in the top or bottom 20th percentile of fit λ 's
626 and averaged each group's block transition dynamics for both wait time and trial initiation time.
627 To compare the normalized dynamics of each group, we fit 4-parameter logistic functions of the
628 following form:

$$y = A + (D - A)/(1 + \exp(-C(x - x_0)))$$

629 to the behavioral curves and compared the four parameters: A (the lower asymptote), D (the
630 upper asymptote), C (the inverse temperature), and x_0 (x -value of the the sigmoid's midpoint).
631 To determine significance for our observed differences, we performed a non-parametric shuffle
632 test. We generated null distributions on differences in the fit parameters by shuffling the labels
633 of the upper and lower percentile λ rats, refitting the logistic to the new shuffled groups' av-
634 erage dynamic curves, and comparing the fit parameters 500 times. We then used these null
635 distributions to calculate p-values for the observed differences in parameters: the area under
636 this distribution evaluated at the actual difference of parameter values (between high and low λ
637 rats) was treated as the p-value.

638 **Trial history effects**

639 To assess wait time sensitivity to previous offers, we focused on 20 μ L catch trials in mixed
640 blocks only. We z-scored the wait times of these trials separately. Next, we averaged wait times
641 depending on whether the previous offer was greater than or less than 20 μ L. For trial initiation
642 times, we used all 20 μ L trials in mixed blocks. We averaged z-scored trial initiation times
643 depending on whether the previous offer was greater or less than 20 μ L. For both wait time
644 and trial initiation time, we defined the sensitivity to previous offers as the difference between
645 average wait time (trial initiation time) for trials with a previous offer less than 20 μ L and

646 trials with a previous offer greater than 20 μ L. We compared wait time and trial initiation time
647 sensitivity to previous offers across rats using a paired Wilcoxon signed-rank test.

648 To capture longer timescale sensitivity across rewards, we regressed previous rewards against
649 wait time and trial initiation time. We focused only on mixed blocks. Additionally, we lin-
650 earized the rewards by taking the binary logarithm of each reward ($\log_2(\text{reward})$). For wait
651 time, we z-scored wait times for catch trials in mixed blocks. Then, we regressed wait times
652 on these trials against the current offer and previous 9 $\log_2(\text{reward})$ offers, including violation
653 trials, along with a constant offset term. Reward offers from a different block (e.g., a previous
654 high block) were given NaN values. For trial initiation times, we again z-scored for mixed
655 block trials only. Then, we regressed against the previous 9 $\log_2(\text{reward})$ offers, not including
656 the current trial, along with a constant offset. Additionally, we set the reward for violation and
657 catch trials to 0, since rats do not receive a reward on these trials.

658 For both wait time and trial initiation time, we used Matlab's builtin regress function to
659 perform the regression. With the coefficients, we found the first non-significant coefficient
660 (coefficient that whose 95% confidence interval contained 0), and set that coefficient and all
661 following coefficients to 0. Finally, we fit a negative exponential decay curve, $y = D \exp -x/\tau$,
662 to each rat's previous trial coefficients (that is, only the previous 9 trial coefficients) for both wait
663 time and trial initiation time and reported the time constant of the exponential decay (tau) for
664 each. If all previous trial coefficients were equal to 0 (as was the case for a vast majority of the
665 wait time coefficients), the time constant was reported as 0. We correlated wait time regression
666 time-constants and trial initiation time regression time-constants using Matlab's builtin corr
667 function.

668 **Learning Dynamics**

669 To assess learning dynamics, we included all sessions after stage 8, not just the sessions
670 that passed criteria for inclusion (above). Because of data limitations examining each session
671 individually (e.g., not every session included both a high and low block), we grouped subse-
672 quent sessions into pairs (i.e., we grouped sessions 1 and 2, sessions 3 and 4, etc.). For each
673 session-pair, we calculated the wait time and trial initiation time ratios as above. To assess the
674 emergence of block effects on wait time data, we regressed wait time for each session against
675 both the current reward and a categorical variable representing the current block identity (1 =
676 low block, 2 = mixed block, 3 = high block). To assess the emergence of previous trial effects
677 on trial initiation time, we regressed trial initiation time for each sessions against the previous
678 reward. We smoothed each regression coefficient over sessions using a 5-session moving av-
679 erage. Finally, we set outlier coefficients (3 scaled median absolute deviations away from a
680 5-point moving median, using Matlab's builtin *isoutlier* function) to NaN. Finally, we averaged
681 regression coefficients over sessions across rats.

682 **Pre-initiation cue task**

683 To modulate the subjective uncertainty in the rat's estimate of state (block) before trial
684 initiation time, we ran a subset of rats on a variation of the task where we cued reward offer
685 before rats initiated a trial ($N = 16$). All other aspects of the task remained identical: reward offer
686 cued played again after the rat initiated the trial, rats waited uncued exponentially-distributed
687 delays for rewards, etc. We included both rats that initially trained on the original task before
688 switching to the pre-initiation cue task ($N = 12$), as well as rats who were trained only on the
689 pre-initiation cue task ($N = 4$). To allow the rats who had started on the original task time to
690 adjust to the new task, we only included data after 30 pre-initiation cue sessions. For the rats
691 who were exclusively trained on the pre-initiation cue task, we included all stage 8 sessions.

692 For all rats, we did not exclude sessions using the wait time criteria (see above).

693 To compare effects for rats who had started on the original task, we performed all analyses
694 for data collected on the original task and on the pre-initiation cue task. First, to confirm that the
695 rats learned that the tone before trial initiation indicated the upcoming reward, we averaged z-
696 scored trial initiation times by the offered reward in mixed blocks. We excluded post-violation
697 trials in the original task session, because those trials repeat the same volume as the previ-
698 ous trial so the rat could conceivably use that to modulate their trial initiation time. All other
699 analyses (sensitivity to the previous reward and previous reward regression) were performed as
700 described above.

701 **References**

- 702 1. Dickinson, A. & Balleine, B. The role of learning in the operation of motivational systems.
703 (2002).
- 704 2. Niv, Y., Daw, N. D., Joel, D. & Dayan, P. Tonic dopamine: opportunity costs and the
705 control of response vigor. *Psychopharmacology* **191**, 507–520 (2007).
- 706 3. Sutton, R. S. & Barto, A. G. *Reinforcement learning: An introduction* (MIT press, 2018).
- 707 4. Pezzulo, G., Rigoli, F. & Chersi, F. The mixed instrumental controller: using value of
708 information to combine habitual choice and mental simulation. *Frontiers in Psychology* **4**,
709 92 (2013).
- 710 5. Gershman, S. J., Horvitz, E. J. & Tenenbaum, J. B. Computational rationality: A converg-
711 ing paradigm for intelligence in brains, minds, and machines. *Science (New York, N.Y.)*
712 **349**, 273–278 (2015).
- 713 6. Dayan, P. How to set the switches on this thing. *Current Opinion in Neurobiology* **22**,
714 1068–1074 (2012).
- 715 7. Keramati, M., Smittenaar, P., Dolan, R. J. & Dayan, P. Adaptive integration of habits
716 into depth-limited planning defines a habitual-goal-directed spectrum. *Proceedings of the*
717 *National Academy of Sciences of the United States of America* **113** (2016).
- 718 8. Daw, N. D., Niv, Y. & Dayan, P. Uncertainty-based competition between prefrontal and
719 dorsolateral striatal systems for behavioral control. *Nature Neuroscience* **8**, 1704–1711
720 (2005).
- 721 9. Keramati, M., Dezfouli, A. & Piray, P. Speed/Accuracy Trade-Off between the Habitual
722 and the Goal-Directed Processes. *PLOS Computational Biology* **7** (2011).
- 723 10. Balleine, B. W. The meaning of behavior: discriminating reflex and volition in the brain.
724 *Neuron* **104**, 47–62 (2019).
- 725 11. Van Der Meer, M., Kurth-Nelson, Z. & Redish, A. D. Information processing in decision-
726 making systems. *The Neuroscientist* **18**, 342–359 (2012).
- 727 12. Redish, A. D., Schultheiss, N. W. & Carter, E. C. The Computational Complexity of Valua-
728 tion and Motivational Forces in Decision-Making Processes. *Current Topics in Behavioral*
729 *Neurosciences* **27**, 313–333 (2016).
- 730 13. Zador, A. *et al.* Catalyzing next-generation Artificial Intelligence through NeuroAI. *Na-*
731 *ture Communications* **14**, 1597 (2023).
- 732 14. Daw, N. D., Gershman, S. J., Seymour, B., Dayan, P. & Dolan, R. J. Model-based influ-
733 ences on humans’ choices and striatal prediction errors. *Neuron* **69** (2011).
- 734 15. Kool, W., Gershman, S. J. & Cushman, F. A. Cost-Benefit Arbitration Between Multiple
735 Reinforcement-Learning Systems. *Psychological Science* **28**, 1321–1333 (2017).

- 736 16. Lak, A. *et al.* Orbitofrontal cortex is required for optimal waiting based on decision con-
737 fidence. *Neuron* **84** (2014).
- 738 17. Khaw, M. W., Glimcher, P. W. & Louie, K. Normalized value coding explains dynamic
739 adaptation in the human valuation process. *Proceedings of the National Academy of Sci-*
740 *ences* **114**, 12696–12701 (2017).
- 741 18. Steiner, A. P. & Redish, A. D. Behavioral and neurophysiological correlates of regret in
742 rat decision-making on a neuroeconomic task. *Nature neuroscience* **17**, 995–1002 (2014).
- 743 19. Charnov, E. L. Optimal foraging, the marginal value theorem. *Theoretical Population Bi-*
744 *ology* **9**, 129–136 (1976).
- 745 20. Stephens, D. W. & Krebs, J. R. in *Foraging theory* (Princeton university press, 2019).
- 746 21. Rigoli, F. Reference effects on decision-making elicited by previous rewards. *Cognition*
747 **192** (2019).
- 748 22. Kawagoe, R., Takikawa, Y. & Hikosaka, O. Expectation of reward modulates cognitive
749 signals in the basal ganglia. *Nature neuroscience* **1**, 411–416 (1998).
- 750 23. Xu-Wilson, M., Zee, D. S. & Shadmehr, R. The intrinsic value of visual information af-
751 fects saccade velocities. *Experimental Brain Research* **196**, 475–481 (2009).
- 752 24. Wang, A. Y., Miura, K. & Uchida, N. The dorsomedial striatum encodes net expected re-
753 turn, critical for energizing performance vigor. *Nature neuroscience* **16**, 639–647 (2013).
- 754 25. Shadmehr, R., Huang, H. J. & Ahmed, A. A. A representation of effort in decision-making
755 and motor control. *Current biology* **26**, 1929–1934 (2016).
- 756 26. Shadmehr, R. & Ahmed, A. A. *Vigor: Neuroeconomics of movement control* (MIT Press,
757 2020).
- 758 27. Hermoso-Mendizabal, A. *et al.* Response outcomes gate the impact of expectations on
759 perceptual decisions. *Nature communications* **11**, 1–13 (2020).
- 760 28. Iigaya, K., Fonseca, M. S., Murakami, M., Mainen, Z. F. & Dayan, P. An effect of sero-
761 tonergic stimulation on learning rates for rewards apparent after long intertrial intervals.
762 *Nature communications* **9**, 1–10 (2018).
- 763 29. Flaherty, C. F. Incentive contrast: A review of behavioral changes following shifts in re-
764 ward. *Animal Learning & Behavior* **10**. ISSN: 1532-5830 (1982).
- 765 30. Constantino, S. M. & Daw, N. D. Learning the opportunity cost of time in a patch-foraging
766 task. *Cognitive, Affective, & Behavioral Neuroscience* **15**, 837–853 (2015).
- 767 31. Verтеchi, P. *et al.* Inference-based decisions in a hidden state foraging task: differential
768 contributions of prefrontal cortical areas. *Neuron* **106**, 166–176 (2020).
- 769 32. Wilson, R. C., Takahashi, Y. K., Schoenbaum, G. & Niv, Y. Orbitofrontal cortex as a
770 cognitive map of task space. *Neuron* **81**, 267–279 (2014).

- 771 33. Jones, J. L. *et al.* Orbitofrontal cortex supports behavior and learning using inferred but
772 not cached values. *Science* **338**, 953–956 (2012).
- 773 34. Davis, H. Transitive inference in rats (*Rattus norvegicus*). *Journal of Comparative Psy-*
774 *chology* **106**, 342 (1992).
- 775 35. Gallistel, C., Mark, T. A., King, A. P. & Latham, P. The rat approximates an ideal detector
776 of changes in rates of reward: implications for the law of effect. *Journal of experimental*
777 *psychology: Animal behavior processes* **27**, 354 (2001).
- 778 36. Behrens, T. E., Woolrich, M. W., Walton, M. E. & Rushworth, M. F. Learning the value
779 of information in an uncertain world. *Nature neuroscience* **10**, 1214–1221 (2007).
- 780 37. Nassar, M. R. *et al.* Rational regulation of learning dynamics by pupil-linked arousal sys-
781 tems. *Nature neuroscience* **15**, 1040–1046 (2012).
- 782 38. Grossman, C. D., Bari, B. A. & Cohen, J. Y. Serotonin neurons modulate learning rate
783 through uncertainty. *Current Biology* **32**, 586–599.e7 (2022).
- 784 39. Gershman, S. J. & Niv, Y. Learning latent structure: carving nature at its joints. *Current*
785 *opinion in neurobiology* **20**, 251–256 (2010).
- 786 40. Miranda, B., Malalasekera, W. M. N., Behrens, T. E., Dayan, P. & Kennerley, S. W. Com-
787 bined model-free and model-sensitive reinforcement learning in non-human primates.
788 *PLoS computational biology* **16** (2020).
- 789 41. Bromberg-Martin, E. S., Matsumoto, M., Nakahara, H. & Hikosaka, O. Multiple timescales
790 of memory in lateral habenula and dopamine neurons. *Neuron* **67**, 499–510 (2010).
- 791 42. Drummond, N & Niv, Y. Model-based decision making and model-free learning. *Current*
792 *biology : CB* **30** (2020).
- 793 43. Balleine, B. W. & Dickinson, A. Effects of Outcome Devaluation on the Performance of a
794 Heterogeneous Instrumental Chain. *International Journal of Comparative Psychology* **18**
795 (2005).
- 796 44. Freidin, E. & Kacelnik, A. Rational choice, context dependence, and the value of infor-
797 mation in European starlings (*Sturnus vulgaris*). *Science* **334**, 1000–1002 (2011).
- 798 45. Hayden, B. Y., Pearson, J. M. & Platt, M. L. Neuronal basis of sequential foraging deci-
799 sions in a patchy environment. *Nature neuroscience* **14**, 933–939 (2011).
- 800 46. Kolling, N., Behrens, T. E., Mars, R. B. & Rushworth, M. F. Neural mechanisms of for-
801 aging. *Science* **336**, 95–98 (2012).
- 802 47. Kahneman, D. & Tversky, A. in *Handbook of the fundamentals of financial decision mak-*
803 *ing: Part I* 99–127 (World Scientific, 2013).
- 804 48. Kőszegi, B. & Rabin, M. A model of reference-dependent preferences. *The Quarterly*
805 *Journal of Economics* **121**, 1133–1165 (2006).

- 806 49. Dayan, P., Niv, Y., Seymour, B. & Daw, N. D. The misbehavior of value and the discipline
807 of the will. *Neural networks* **19**, 1153–1160 (2006).
- 808 50. Sweis, B. M. *et al.* Sensitivity to "sunk costs" in mice, rats, and humans. *Science (New*
809 *York, N.Y.)* **361** (2018).
- 810 51. Starkweather, C. K., Babayan, B. M., Uchida, N. & Gershman, S. J. Dopamine reward
811 prediction errors reflect hidden-state inference across time. *Nature neuroscience* **20**, 581–
812 589 (2017).
- 813 52. Khalvati, K., Kiani, R. & Rao, R. P. Bayesian inference with incomplete knowledge ex-
814 plains perceptual confidence and its deviations from accuracy. *Nature communications* **12**,
815 5704 (2021).
- 816 53. Lak, A., Nomoto, K., Keramati, M., Sakagami, M. & Kepecs, A. Midbrain dopamine
817 neurons signal belief in choice accuracy during a perceptual decision. *Current Biology* **27**,
818 821–832 (2017).
- 819 54. Bromberg-Martin, E. S., Matsumoto, M., Hong, S. & Hikosaka, O. A pallidus-habenula-
820 dopamine pathway signals inferred stimulus values. *Journal of neurophysiology* **104**, 1068–
821 1076 (2010).
- 822 55. Feher da Silva, C. & Hare, T. A. Humans primarily use model-based inference in the
823 two-stage task. *Nature Human Behaviour* **4**, 1053–1066 (2020).
- 824 56. Miller, K. J., Botvinick, M. M. & Brody, C. D. Dorsal hippocampus contributes to model-
825 based planning. *Nature neuroscience* **20**, 1269–1276 (2017).
- 826 57. Polanía, R., Woodford, M. & Ruff, C. C. Efficient coding of subjective value. *Nature*
827 *neuroscience* **22**, 134–142 (2019).
- 828 58. Louie, K. & Glimcher, P. W. Efficient coding and the neural representation of value. *An-*
829 *nals of the New York Academy of Sciences* **1251**, 13–32 (2012).
- 830 59. Tymula, A. & Glimcher, P. Expected subjective value theory (ESVT): A representation of
831 decision under risk and certainty. *Available at SSRN 2783638* (2021).
- 832 60. Barlow, H. B. *et al.* Possible principles underlying the transformation of sensory messages.
833 *Sensory communication* **1**, 217–233 (1961).
- 834 61. Padoa-Schioppa, C. Range-adapting representation of economic value in the orbitofrontal
835 cortex. *Journal of Neuroscience* **29**, 14004–14014 (2009).
- 836 62. Weber, A. I., Krishnamurthy, K. & Fairhall, A. L. Coding principles in adaptation. *Annual*
837 *review of vision science* **5**, 427–449 (2019).
- 838 63. Kobayashi, S., de Carvalho, O. P. & Schultz, W. Adaptation of reward sensitivity in or-
839 bitofrontal neurons. *Journal of Neuroscience* **30**, 534–544 (2010).
- 840 64. Heffner, H. E., Heffner, R. S., Contos, C. & Ott, T. Audiogram of the hooded Norway rat.
841 *Hearing research* **73**, 244–247 (1994).

- 842 65. Cassandra, A. R. *Exact and approximate algorithms for partially observable Markov de-*
843 *cision processes* (Brown University, 1998).

844 **Acknowledgments**

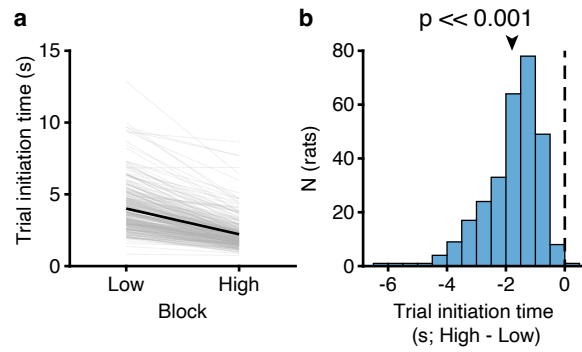
845 We thank Paul Glimcher, Catherine Hartley, Roozbeh Kiani, Kenway Louie, Kevin Miller,
846 Cristina Savin and members of the Constantinople lab for feedback. We thank Madori Spiker,
847 Daljit Kaur, Mitzi Adler-Wachter, Royall McMahan Ward, and Luke Chen for animal training.

848 **Funding:** This work was supported by a K99/R00 Pathway to Independence Award (R00MH111926),
849 an Alfred P. Sloan Fellowship, a Klingenstein-Simons Fellowship in Neuroscience, an NIH Di-
850 rector's New Innovator Award (DP2MH126376), an NSF CAREER Award, R01MH125571,
851 and a McKnight Scholars Award to C.M.C. A.M. was supported by 5T90DA043219 and F31MH130121.
852 A.M. and S.S.S. were supported by 5T32MH019524.

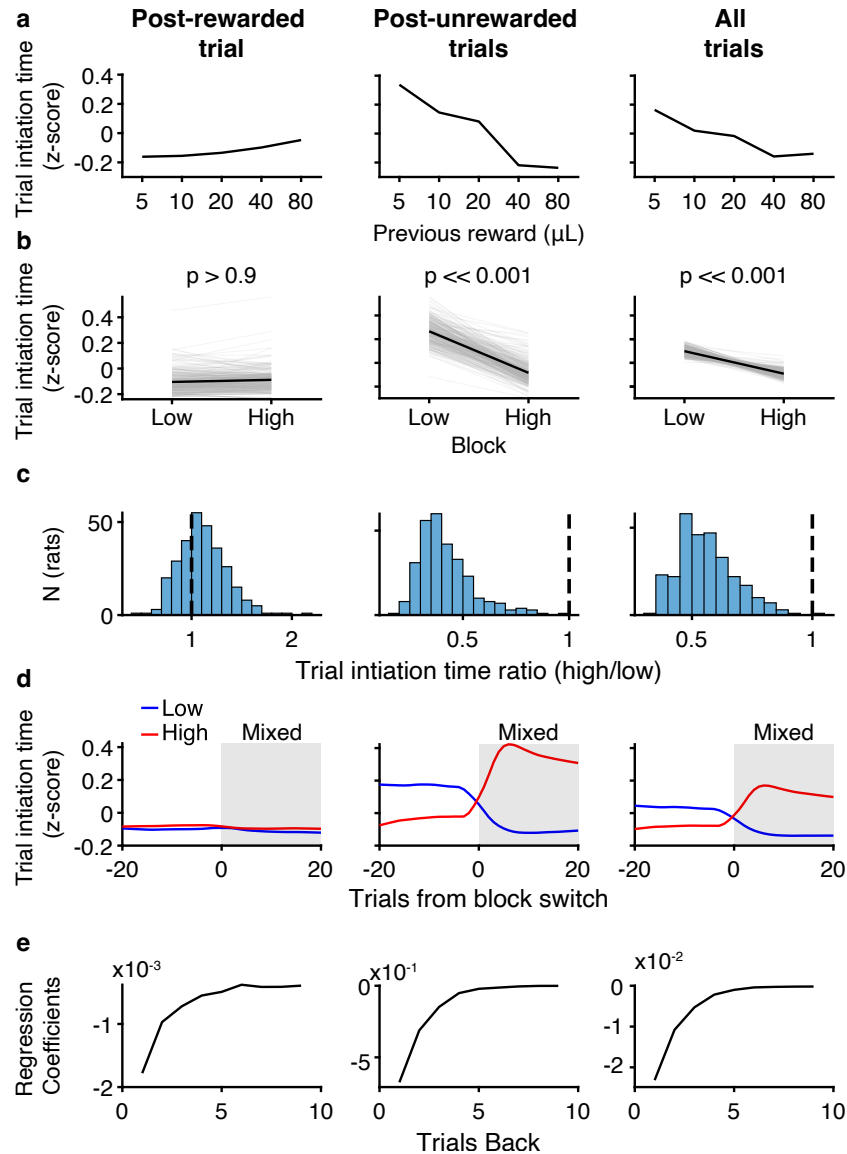
853 **Author Contributions**

854 C.M.C. designed the task. V.B. wrote training software. V.B. and C.M.C. developed soft-
855 ware infrastructure for the high-throughput facility and data management. A.M. and C.M.C.
856 analyzed the data. S.S.S. contributed to behavioral experiments. A.M. prepared the figures.
857 C.M.C. and A.M. wrote the manuscript. C.M.C supervised the project.

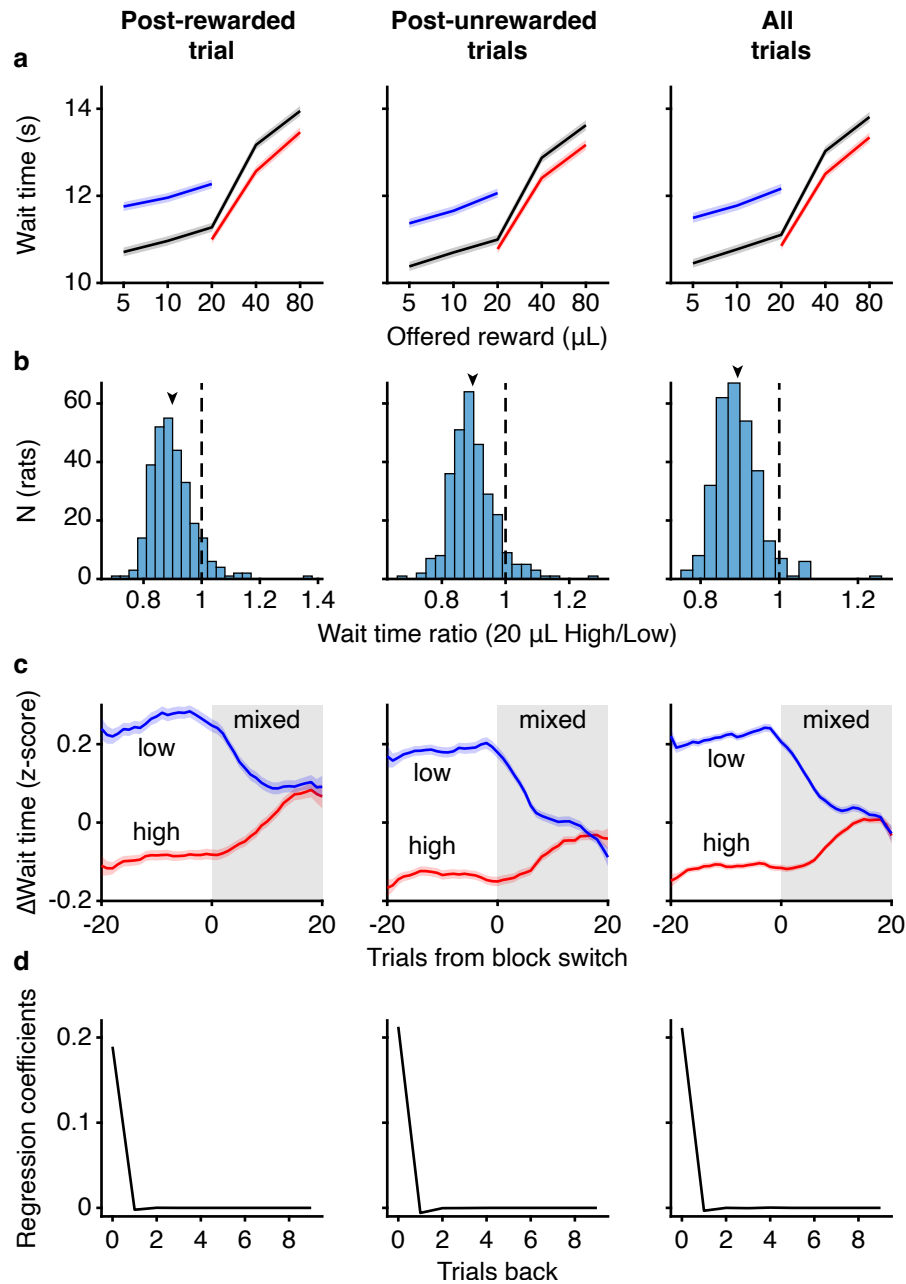
858 **Supplementary materials**



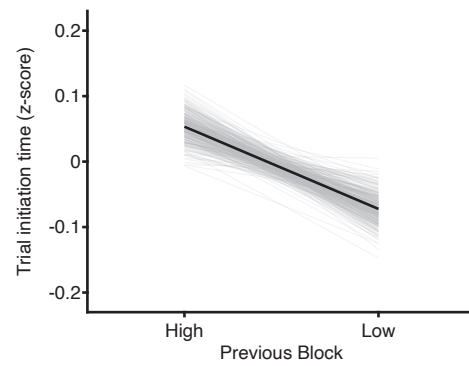
Extended Data Fig. 1: **Trial initiation time in units of seconds** **a.** Trial initiation time by block ($N = 291$). Data are replotted from Fig. 1h but in units of seconds. **b.** Trial initiation time difference (high - low) across all rats. (Wilcoxon signed-rank test, $N = 291$).



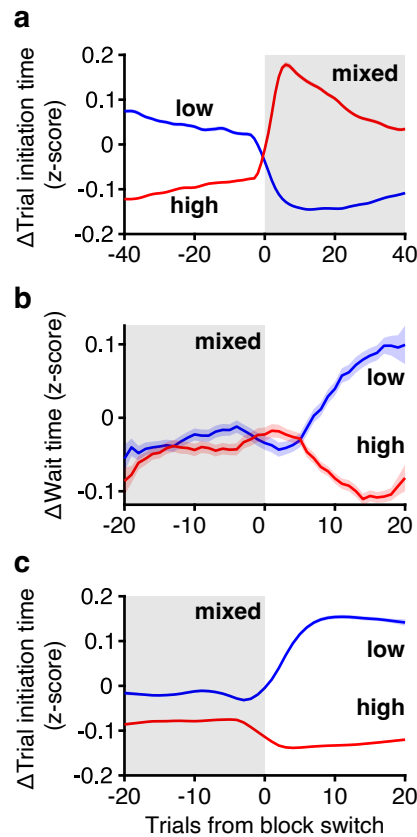
Extended Data Fig. 2: **Trial initiation times depend on previous trial outcome.** **a.** Trial initiation time by previous reward in mixed blocks for (left) post-rewarded trials, (center) post-unrewarded trials, and (right) all trial. **b.** Trial initiation time averaged over block (Wilcoxon Signed-rank test, $N = 291$). **c.** Trial initiation time ratio (mean trial initiation time in high blocks/low blocks, $N = 291$). **d.** Mean change trial initiation times from low or high blocks to mixed blocks, $N = 291$. **e.** Previous trial regression coefficients in mixed blocks, $N = 291$.



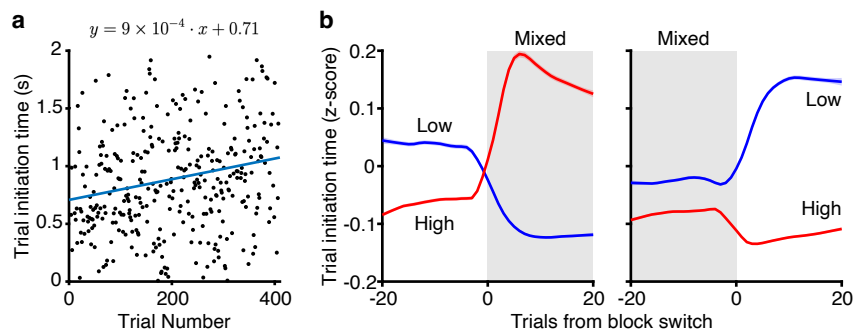
Extended Data Fig. 3: **Wait times are not affected by previous trial outcome.** **a.** Average wait time by volume for each block conditioned on whether the previous trial was (left) rewarded or (center) unrewarded, and (right) all trials ($N = 291$). **b.** Wait time ratios (wait time for 20 μL High/Low) across rats ($N = 291$). **c.** Wait time dynamics transitioning from low (blue) or high (red) blocks into mixed blocks ($N = 291$). **d.** Reward



Extended Data Fig. 4: **Average trial initiation time in mixed blocks conditioned on the previous block** ($N = 291$).

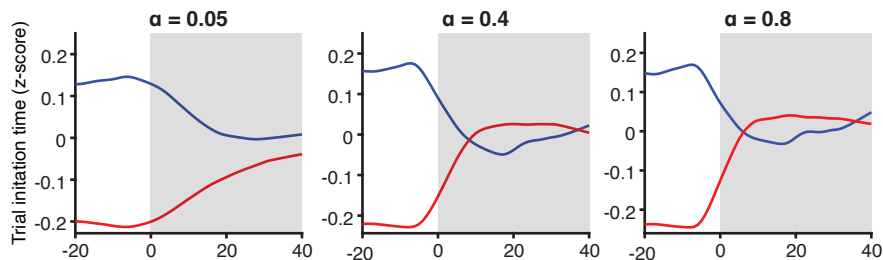


Extended Data Fig. 5: **Dynamics of wait times (top) and trial initiation times (bottom) at transitions from mixed to high (red) or low (blue) blocks.** **a.** Data are replotted from Fig. 2b, but with expanded x-axis limits. Trial initiation times still maintain contrast effects 40 trials into mixed blocks. **b.** Wait time transitions from mixed to high (red) and low (blue) blocks. **c.** Trial initiation time transitions from mixed to high (red) and low (blue). Block labels refer to the block at trial 0 after the mixed block. Colors are flipped relative to Fig. 2b because a current low block (blue here) is always preceded by a high block (red in Fig. 2b).

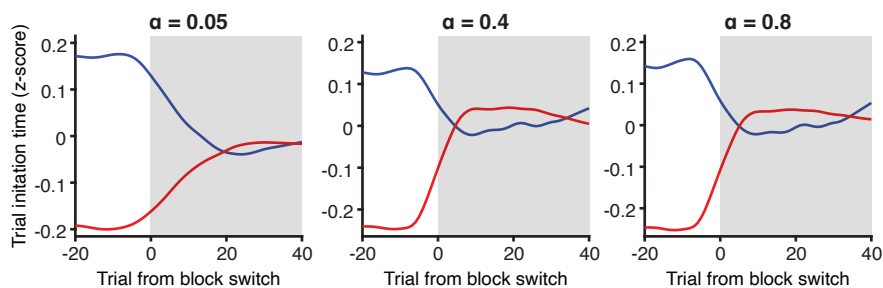


Extended Data Fig. 6: **Satiety effects for trial initiation time are modest and do not qualitatively affect results a.** Trial initiation time for an example session as a function of trial number. Line is least-squares regression. **b.** Trial initiation times block transition plots without detrending. Results are qualitatively similar to Fig. 2.

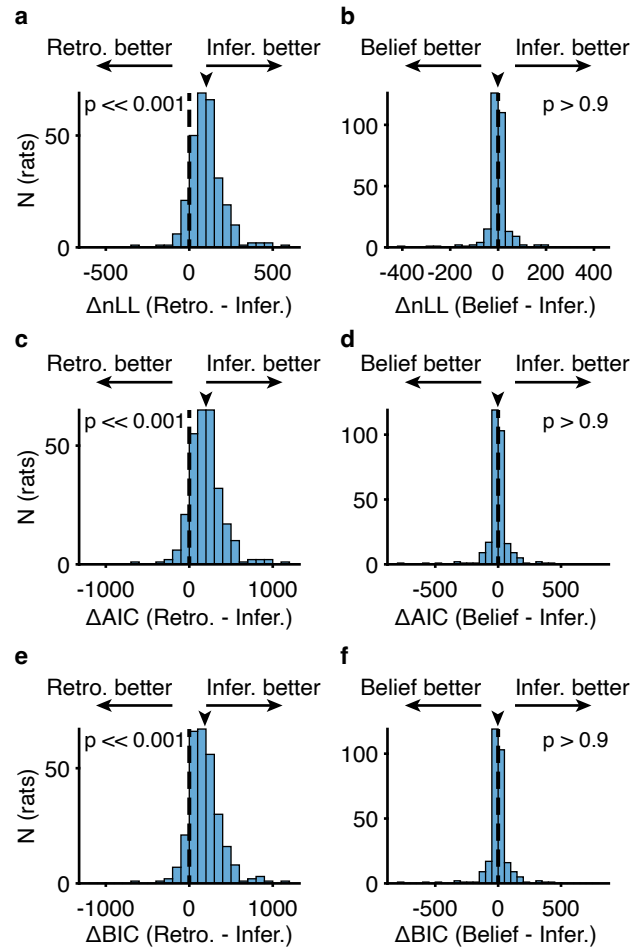
Vanilla learning rate model: a single, static learning rate



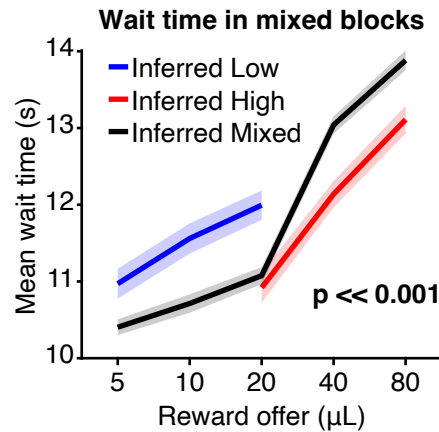
RPE-gain learning rate: learning rate gain = |RPE|



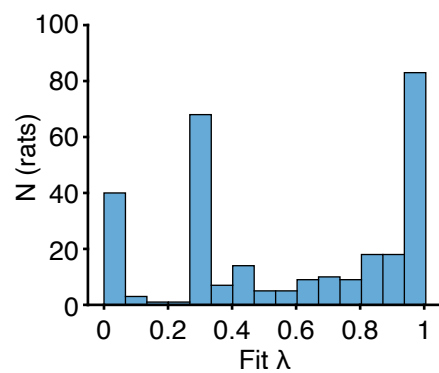
Extended Data Fig. 7: **Alternative retrospective models fail to capture both fast and slow trial initiation time dynamics at block transitions.** Trial initiation time model transitions from low (blue) or high (red) blocks to mixed blocks. Top: A “vanilla” learning rate model with a single, static learning rate. Bottom: a dynamic learning rate model where learning rate gain is equal to the unsigned RPE of that trial.



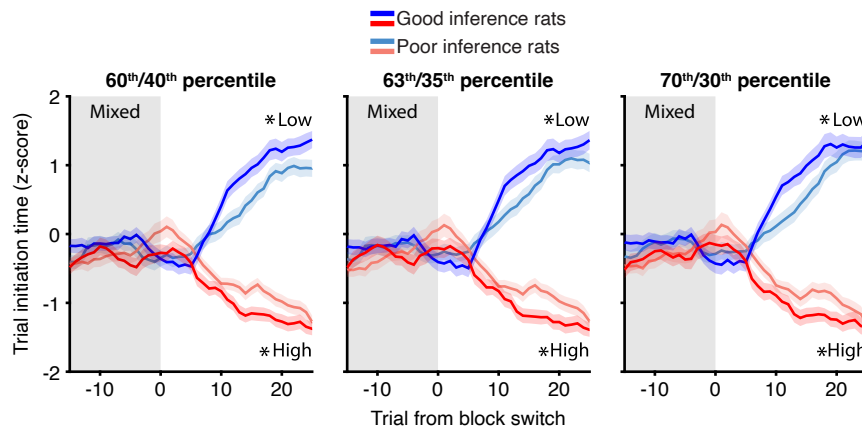
Extended Data Fig. 8: **Model comparison for wait times favors inferential over retrospective model, but does not distinguish between inferential and belief state models.** **a-b.** Cross-validated negative log-likelihood comparing inferential model and (a.) retrospective or (b.) belief state model. **c-d.** Akaike information criterion (AIC) comparing inferential model and (c.) retrospective or (d.) belief state model. **e-f.** Bayesian information criterion (BIC) comparing inferential model and (e.) retrospective or (f.) belief state model. For each, Wilcoxon signed-rank test, $N = 291$



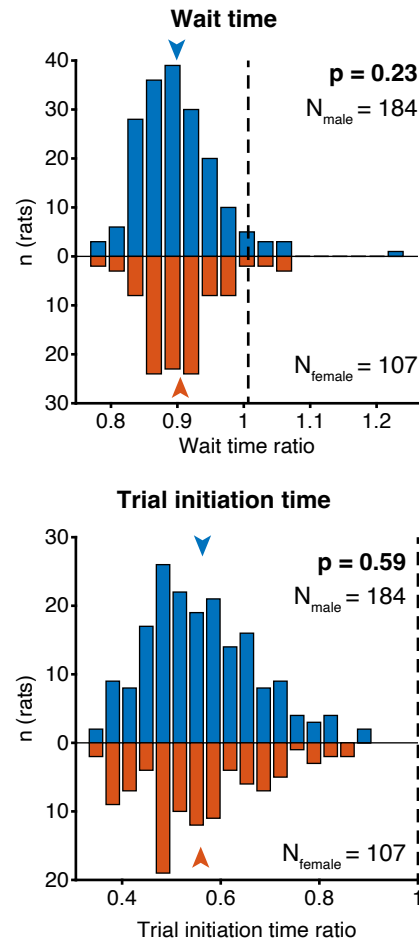
Extended Data Fig. 9: **Inferential model identifies mistaken inferences during mixed blocks across rats.** **a.** Average wait time curves conditioned by model-inferred block in mixed blocks only in held-out test set across rats. **b.** Wait time ratio (wait time on 20 μL inferred high/low trials) is modulated by inferred block ($p \ll 0.001$, Wilcoxon Signed-rank test, $N = 291$)



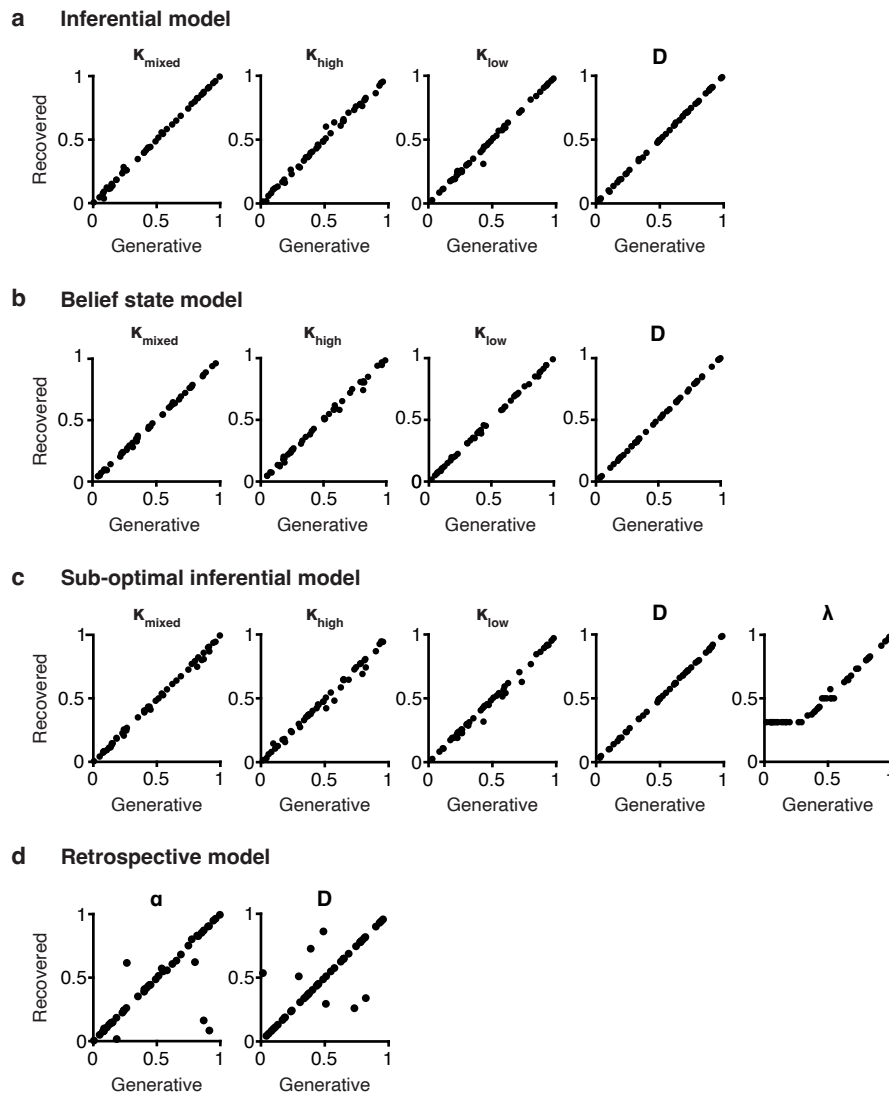
Extended Data Fig. 10: **Sub-optimal inferential model with lambda.** Distribution of λ fit over rats ($N = 291$).



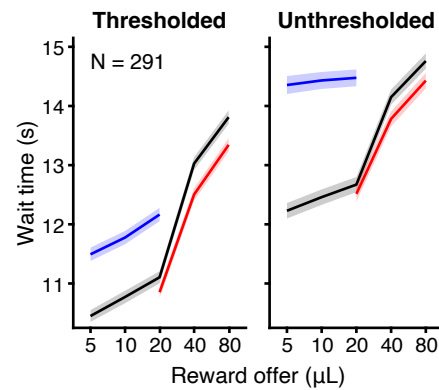
Extended Data Fig. 11: **Differential wait time dynamics based on λ from sub-optimal Bayes model are robust across a range of percentiles.**



Extended Data Fig. 12: **Males and females have comparable wait time ratios (top) and trial initiation time ratios (bottom).** Wait time $p = 0.23$, Wilcoxon Rank-sum test, $N = 184$ males, 107 females. Trial initiation time $p = 0.59$, Wilcoxon Rank-sum test, $N = 184$ males, 107 females.



Extended Data Fig. 13: **Models are able to recover generative parameters.** $N = 48$ random parameter sets.



Extended Data Fig. 14: **Wait time curves without threshold (right) have qualitatively similar context effects, but longer average wait times.** Wait times one standard deviation above the pooled session mean were excluded for most analyses in this study (left). Including all wait times preserved the contextual effects, but resulted in longer average wait times, as the mean is particularly sensitive to outliers. Outlier wait times tended to occur in low blocks, likely due to attentional or motivational lapses. Therefore, the main difference between the thresholded and unthresholded data is that the wait time curves in low blocks are both flatter and longer in the unthresholded data.