

Rats rapidly switch between retrospective and inferential value computations

Andrew Mah¹, 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.

There are many ways to compute value. For instance, animals can compute value by learning from the past or by imagining future outcomes, but it is unclear if or how these computations interact. We used high-throughput training to collect statistically powerful datasets from 240 rats performing a temporal wagering task with hidden reward states. Rats adjusted how quickly they initiated trials and how long they waited for rewards across states, balancing effort and time costs against expected rewards. Statistical modeling revealed that animals computed the value of the environment differently when initiating trials versus when deciding how long to wait for rewards, even though these decisions were only seconds apart. This work reveals that sequential decisions use parallel value computations on single trials.

Main Text

The value of the environment, or how much reward it is expected to yield, determines animals' motivational states and sets their expectations for error-based learning (1–3). But how

18 are values computed? Reinforcement learning systems can store or “cache” values of states,
19 actions, or outcomes that are learned directly from experience, or they can compute values us-
20 ing a learned model of the environment to simulate possible futures (3). These different value
21 computations have distinct tradeoffs, and a central question is how neural systems decide which
22 computations to use or whether/how to combine them (4–8). However, it is difficult to deter-
23 mine the value computations that subjects use, especially over behaviorally relevant timescales
24 of seconds. In standard two-alternative forced choice tasks, the behavioral read-out is a binary
25 choice, and the underlying values driving choice are obscure. State-of-the-art methods for re-
26 vealing how values are computed use regression models that pool data over entire behavioral
27 sessions (9), or pre-determined subsets of trials (10), thereby obscuring moment-by-moment
28 changes in value computations. Therefore, whether or how multiple value computations inter-
29 act on rapid timescales in the same subject is unclear.

30 **Rats’ deliberative and motivational decisions are sensitive to the value of** 31 **the environment.**

32 We developed a temporal wagering task for rats, in which they were offered one of several
33 water rewards on each trial, the volume of which (5, 10, 20, 40, 80 μ L) was indicated by a tone
34 (Fig. 1A). The reward was assigned randomly to one of two ports, indicated by an LED. The
35 rat could wait for an unpredictable delay to obtain the reward, or at any time could terminate
36 the trial by poking in the other port (“opt-out”). Wait times were defined as how long rats
37 waited before opting out. Trial initiation times were defined as the time from opting-out or
38 consuming reward to initiating a new trial. Reward delays were drawn from an exponential
39 distribution, and on 15-25 percent of trials, rewards were withheld to force rats to opt-out,
40 providing a continuous behavioral readout of subjective value (Fig. 1B) (11–13). We used
41 a high-throughput facility to train 240 rats using computerized, semi-automated procedures.

42 The facility generated statistically powerful datasets (median = 30,842 behavioral trials, 65
43 sessions).

44 The task contained latent structure: rats experienced blocks of 40 completed trials (hidden
45 states) in which they were presented with low (5, 10, or 20 μ L) or high (20, 40, or 80 μ L) re-
46 wards (12). These were interleaved with “mixed” blocks which offered all rewards (Fig. 1C).
47 20 μ L was present in all blocks, so comparing behavior on trials offering this reward revealed
48 contextual effects (i.e., effects of hidden states). The hidden states differed in their average re-
49 ward and therefore in their opportunity costs, or what the rat might miss out on by continuing to
50 wait. According to foraging theories, the opportunity cost is the long-run average reward, or the
51 value of the environment (14). In accordance with these theories (14, 15), rats adjusted how long
52 they were willing to wait for rewards in each block, and on average waited \sim 10 percent less
53 time for 20 μ L in high blocks, when the opportunity cost was high, compared to in low blocks (p
54 \ll 0.001, Wilcoxon signed-rank test, $N = 240$; Fig. 1D-F). These are strong contextual effects
55 compared to previous studies (12, 16).

56 Trial initiation times were modulated by blocks in a similar pattern as the wait times, with
57 rats initiating trials more quickly in high compared to low blocks ($p \ll$ 0.001, Wilcoxon
58 signed-rank test, $N = 240$; Fig. 1G-I). Previous work suggests that this pattern optimally bal-
59 ances the costs of vigor against the benefits of harvesting reward in environments with different
60 reward rates (2, 17). Therefore, both the trial initiation times, which reflect motivation, and the
61 wait times, which reflect deliberating between waiting and opting-out, were modulated by the
62 value of the environment.

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

64 Surprisingly, wait and trial initiation times exhibited dramatically different dynamics at
65 block transitions. In mixed blocks, the wait times following high and low blocks converged to a

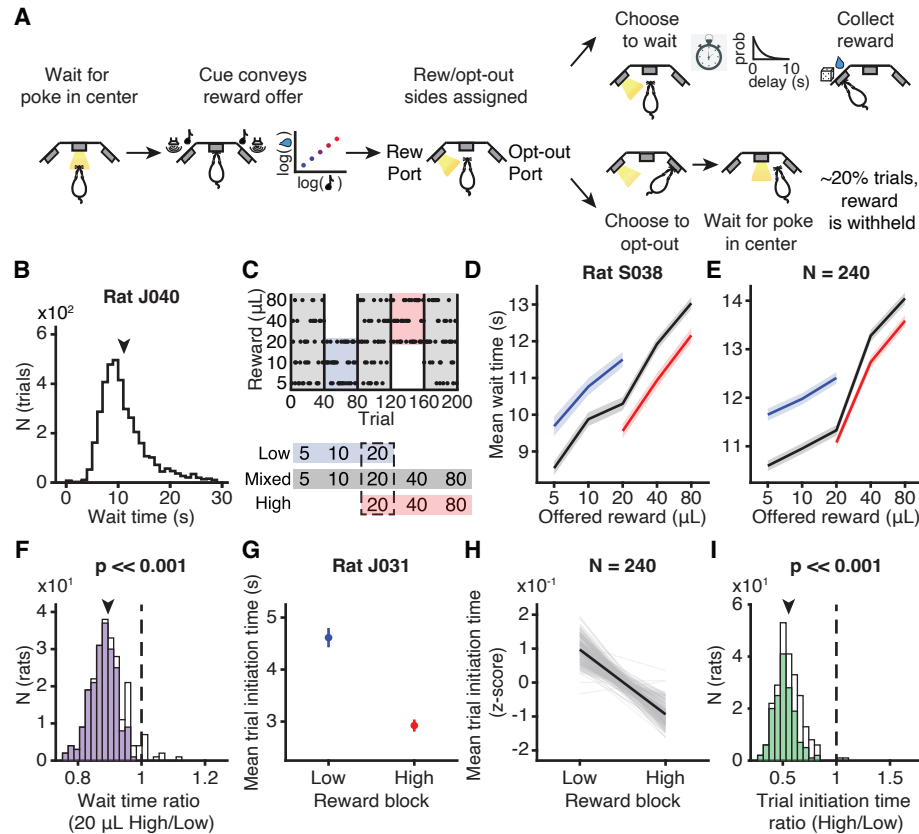


Figure 1: Wait time and trial initiation time were modulated by the value of the environment. **A.** Schematic of behavioral paradigm. **B.** Distribution of wait times for one rat. **C.** Block structure of task. **D-E.** Average wait time on catch trials by reward in each block for (D) one rat and (E) averaged across rats. **F.** Wait time ratio (average wait time for 20 μ L 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, $N = 240$. **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, $N = 240$.

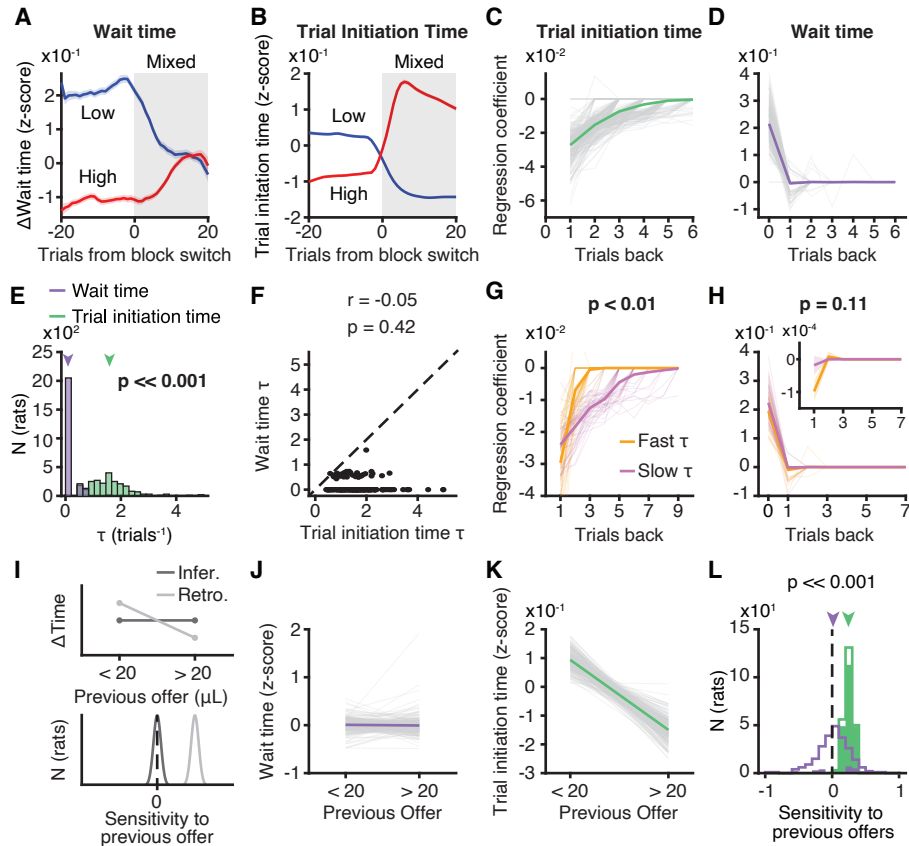


Figure 2: Wait and trial initiation times use distinct estimates of the value of the environment. **A-B.** Mean change in wait times (A) and trial initiation times (B) from low or high blocks to mixed blocks, $N = 240$. Data are mean \pm S.E.M. **C-D.** Regression coefficients for (C) trial initiation 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, $p \ll 0.001$, Wilcoxon sign-rank test, $N = 240$, and (F) uncorrelated, $r = 0.08$, $p = 0.18$, Pearson linear correlation, $N = 240$. **G-H.** Fast or slow initiation time τ ($<20^{\text{th}}$ or $>80^{\text{th}}$) meaningfully divided rats based on their initiation time regression coefficients (G; $p \ll 0.01$, one-tailed permutation test, $N = 47$), but not wait time coefficients (H; $p = 0.1$, one-tailed permutation test, $N = 47$). **I.** Predictions for sensitivity to previous offers (behavior conditioned on previous offer $<20\mu\text{L}$ - $>20\mu\text{L}$) for fixed (light) versus sequentially-updated (dark) estimates of environmental value, consistent with inferential and retrospective strategies, respectively. **J.** Wait time on 20 μL catch trials in mixed blocks conditioned on previous reward offer. Difference is significant ($p < 0.05$) in only 28/240 rats, Wilcoxon rank-sum test. **K.** Trial initiation time in mixed blocks conditioned on previous reward offer. Difference is significant ($p < 0.05$) in 212/240 rats, Wilcoxon rank-sum test. **L.** Sensitivity to previous offers for wait time (purple) and trial initiation time (green). $p \ll 0.001$, Wilcoxon sign-rank test, $N = 240$. Colored bars are individual rats with $p < 0.05$, Wilcoxon rank-sum test.

66 common value, regardless of the previous block type, suggesting the use of a fixed estimate of
67 environmental value in mixed blocks (Fig. 2A). Trial initiation times, however, showed longer
68 timescale effects such that initiation times in mixed blocks strongly depended on the previous
69 block identity (Fig. 2B). These longer timescale dynamics, which are reminiscent of incen-
70 tive contrast effects (18), were also evident in the transitions from mixed blocks into high/low
71 blocks for trial initiation times, but not wait times (fig. S1), indicating that trial initiation and
72 wait times utilize distinct estimates of the value of the environment.

73 To better characterize their temporal dynamics, we regressed the trial initiation and wait
74 times against rewards offered on previous trials. We included current rewards as regressors
75 in the wait time model, and restricted this analysis to mixed blocks only. Examination of the
76 regression coefficients revealed qualitatively different dynamics, in which the wait times were
77 explained by the reward offered on the current trial, but the trial initiation times reflected an
78 exponentially weighted effect of previous rewards, consistent with a model-free temporal dif-
79 ference learning rule (Fig. 2C,D). We fit exponential curves to the previous trial coefficients
80 for each rat, and found that the distributions of exponential decay time constant parameters (τ)
81 were significantly different for the trial initiation and wait times ($p \ll 0.01$, Wilcoxon sign-
82 rank test, $N = 240$; Fig. 2E). Moreover, τ parameters were not correlated across models ($r =$
83 0.08 , $p = 0.18$, Pearson linear correlation, $N = 240$, Fig. 2F).

84 To leverage individual variability across rats, we compared rats with fast and slow temporal
85 integration for trial initiation times (τ from exponential fit to regression coefficients < 20 th or
86 > 80 th percentiles). There were differences in temporal integration for trial initiation times,
87 but not wait times, for these groups (Fig. 2G-H, trial initiation time $p \ll 0.001$, wait time $p =$
88 0.5 , permutation test, $N = 111$). Collectively, these data suggest that within a block, wait times
89 use a fixed estimate of the value of the environment, whereas trial initiation times are sensitive
90 to previous rewards (Fig. 2C,D). Indeed, for almost all rats (89%), wait times for 20 μ L offers

91 in mixed blocks were not significantly different if they were preceded by rewards that were
92 smaller or larger than $20\mu\text{L}$ ($p > 0.05$, Wilcoxon rank-sum test, $N = 212/240$). However, for
93 89% of rats, trial initiation times were significantly modulated by previous rewards, suggesting
94 fixed and incrementally updated estimates of the value of the environment, respectively ($p <$
95 0.05 , Wilcoxon rank-sum test, $N = 212/240$, Fig. 2I-L).

96 **Computational modeling reveals distinct value computations for sequential** 97 **decisions.**

98 Our data suggest that rats' sequential decisions (when to initiate trials and how long to wait
99 for rewards) reflect different value computations. We developed behavioral models for wait and
100 trial initiation times, inspired by foraging theories (14). The wait time model implemented a
101 trial value function that scaled with the offered reward and decayed to reflect reward probability
102 over time (11). The model's predicted wait time was when the value function fell below the
103 value of the environment (opportunity cost) on each trial (Fig. 3A). Different versions of the
104 model estimated the value of the environment using different computations.

105 Analysis of rats' trial initiation times suggests that they estimate the value of the environ-
106 ment as a running average of rewards (Fig. 2C) (2, 12, 19). We refer to this computation as
107 retrospective, as it reflects past experience (20). Alternatively, rats' wait times reflected the
108 use of discrete estimates of block value (Fig. 2A,D,J). Therefore, rats might infer the current
109 block (20–24), and use fixed estimates of block value based on that inference. We refer to this
110 computation as inferential, since it requires hidden state inference.

111 The inferential model selected the most likely block using Bayes' Rule with a prior that
112 incorporated reward history and knowledge of the block transition structure. This model reca-
113 pitulated the rats' wait times converging to a common value in mixed blocks (Fig. 3B-C). This
114 reflects the model's use of a fixed estimate of the value of the environment in each block.

115 In the retrospective case, the value of the environment was estimated as a recency-weighted
116 average of offered rewards according to a temporal-difference learning rule (Fig. 3C). A static
117 learning rate was unable to capture the rats' behavior (fig. S2). Previous work has shown that
118 animals adjust their learning rates depending on the volatility in the environment, since it is
119 advantageous to learn faster in dynamic environments (25–27). Therefore, our model scaled the
120 learning rate by the trial-by-trial change in the inferential model's beliefs about the hidden state
121 (derivative of the posterior, see Methods).

122 We fit these models to rats' wait times. By several model comparison metrics, wait times
123 were better fit by the inferential model that used hidden state inference to select block-specific
124 estimates of the value of the environment ($p \ll 0.001$, Wilcoxon signed-rank test, $N = 240$;
125 Fig. 3F, fig. S3), consistent with that model reproducing the wait time dynamics (Fig. 2A,3B).
126 We also used the model to identify trials in mixed blocks where the rats were likely to make
127 mistaken inferences. The rats' wait times reflected these mistaken inferences, further indicating
128 that their wait times were well-described by the inferential model (fig. S4).

129 We also developed a “belief state” model that estimated the value of the environment as
130 the sum of block-specific values weighted by their posterior probabilities. These models make
131 qualitatively similar predictions about the average wait times. In fact, when the posterior beliefs
132 are stable, which is often the case, the belief state and inferential models are identical, and model
133 comparison did not favor one model over the other (data not shown).

134 While the inferential model captured rats' wait times, the retrospective model captured two
135 key features of their trial initiation times, which we modeled as inversely proportional to the
136 value of the environment (Fig. 2D-E) (2). First, with a sufficiently small learning rate (<0.1 ,
137 fig. S2), the model integrated reward history on long timescales such that trial initiation times in
138 mixed blocks depended on the previous block identity. Second, the dynamic learning rate cap-
139 tured the rapid behavioral dynamics at block transitions. We explored versions of the dynamic

140 learning rate that did not reflect inference, including using the unsigned reward prediction error
141 or a running average of reward prediction errors (27). However, these models could not cap-
142 ture both short and long timescale dynamics at block transitions (fig. S2). This suggests that
143 trial initiation times reflect a retrospective computation that is influenced by subjective belief
144 distributions (25, 26).

145 To leverage individual differences, we turned to the inferential model of wait times. We
146 added a parameter, λ , that controlled the extent to which the model used an optimal prior, $\lambda = 1$,
147 versus an uninformative prior, $\lambda = 0$ (Fig. 3F; fig. S5). We divided the rats into groups with low
148 or high values of λ ($\lambda < 20$ th or > 80 th percentiles), and compared the parameters of logistic
149 functions fit to the average wait time dynamics for these groups. Rats with optimal and poor
150 inference exhibited significantly different dynamics at transitions from mixed into low blocks,
151 indicated by different inverse temperature parameters, but not into high blocks, (mix to low,
152 $p < 0.05$, mix to high, $p = 0.08$, one-tailed permutation test, $N = 180$ Fig. 3G). This suggests
153 that λ may have captured variability in rats' priors over low blocks in particular. There was no
154 difference in the dynamics of trial initiation times for those same groups of rats (mixed to low:
155 $p = 0.3$, mixed to high: $p = 0.2$, one-tailed permutation test, $N = 180$; Fig. 3G).

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

157 Structure learning is the process of learning the hidden structure of environments, including
158 latent states and transition probabilities between them (28). If wait and trial initiation times dif-
159 ferentially required knowledge of latent task structure, they should exhibit different dynamics
160 over training. In the final stage of training, when rats were introduced to the hidden states, their
161 wait times for 20 μ L gradually became sensitive to the reward block (Fig. 4A). We observed
162 a gradual increase in the magnitude of reward and block regression coefficients that mirrored
163 the behavioral sensitivity to hidden states (Fig. 4B). In contrast, trial initiation times exhibited

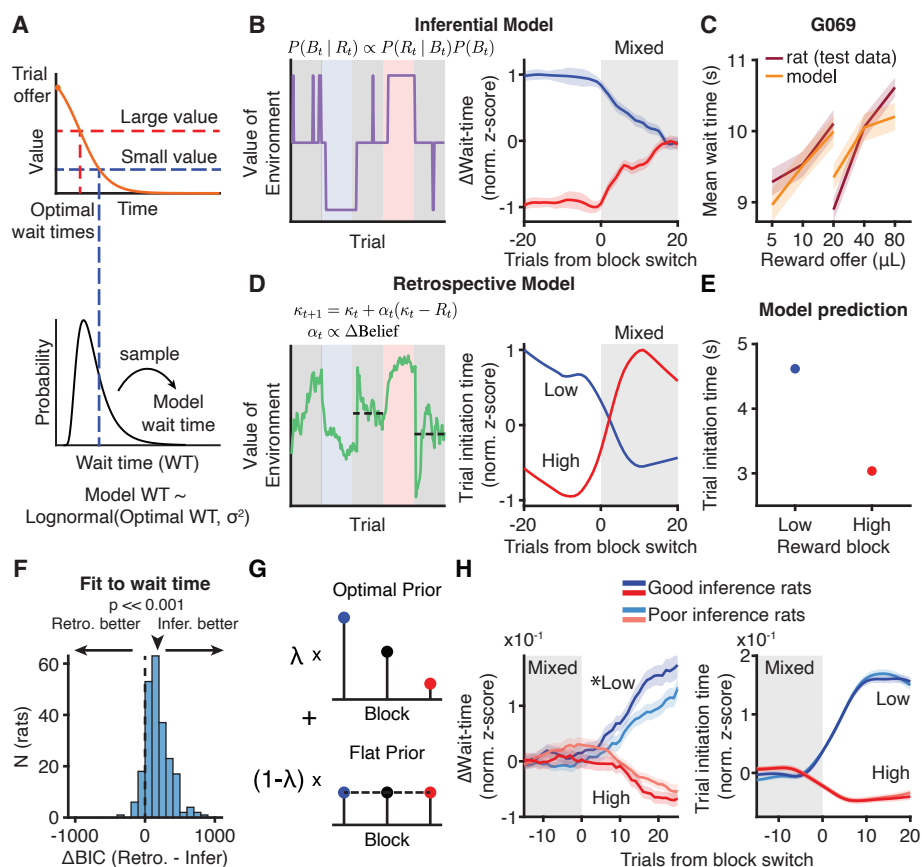


Figure 3: Computational modeling reveals distinct value computations for wait time and trial initiation **A.** Model schematic. **B.** Example opportunity cost and wait time dynamics from inferential model. **C.** Inferential model fit to rats can capture wait time behavior in held-out test data. **D.** Example opportunity cost and wait time dynamics from retrospective model. **E.** Retrospective model can qualitatively capture trial initiation time behavior. **F.** Model comparison using Δ BIC prefers inferential model compared to retrospective model when fit to wait time data ($p < 0.001$, Wilcoxon Signed-rank test, $N = 240$) **G.** Schematic for sub-optimal inference model **H.** Transitions from mixed to low (blue) or high (red) blocks for (G) wait time or (H) trial initiation time separated by quality of inference ($\lambda < 20\text{th}$ or $> 80\text{th}$ percentile). * $p < 0.05$, one-tailed non-parametric shuffle test comparing logistic fit parameters, $N = 47$. Data are mean \pm S.E.M.

164 block sensitivity on the first session in the final training stage (Fig. 4A). This sensitivity was
165 comparable early and late in training, consistent with animals using previous rewards to a simi-
166 lar extent at these timepoints (Fig. 4C). These data suggest that block sensitivity for wait times,
167 but not trial initiation times, required learned knowledge of hidden task states, and that these
168 decisions reflected computations with distinct learning dynamics.

169 The modest increase in trial initiation time block sensitivity over training is consistent with
170 the gradual use of a dynamic learning rate that reflected learned knowledge of the blocks. A
171 hallmark of the dynamic learning rate was the “overshoot” after transitions from high to mixed
172 blocks (difference between maximum trial initiation time after transitioning and the trial ini-
173 tiation time 20 trials post-transition; Fig. 2B). The overshoot became more prominent with
174 training (Fig. 4D), on a similar timescale as block sensitivity for wait times (Fig. 4E), suggest-
175 ing a shared mechanism.

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

177 Why would animals use a retrospective computation at trial initiation, but rely on an inferen-
178 tial computation as rats deliberated just 1-2 seconds later? In non-human primates, the decision
179 to initiate trials can also reflect retrospectively computed values that differ from the values gov-
180 erning the subsequent choice (29, 30). One possibility is that motivation and approach behavior
181 rely on neural circuits that do not support inference (31). Another possibility is that actions
182 more distal to rewards are more likely to be retrospective, because there are more steps required
183 to mentally simulate outcomes for forward-looking strategies like planning (32, 33). According
184 to either hypothesis, the decision of when to initiate a trial is inherently retrospective.

185 Theoretical work in reinforcement learning has suggested that the brain should select the
186 strategy that is the fastest and most accurate when taking into account uncertainty (8, 34). There-
187 fore, perhaps trial initiation times are retrospective because the rats’ subjective beliefs about the

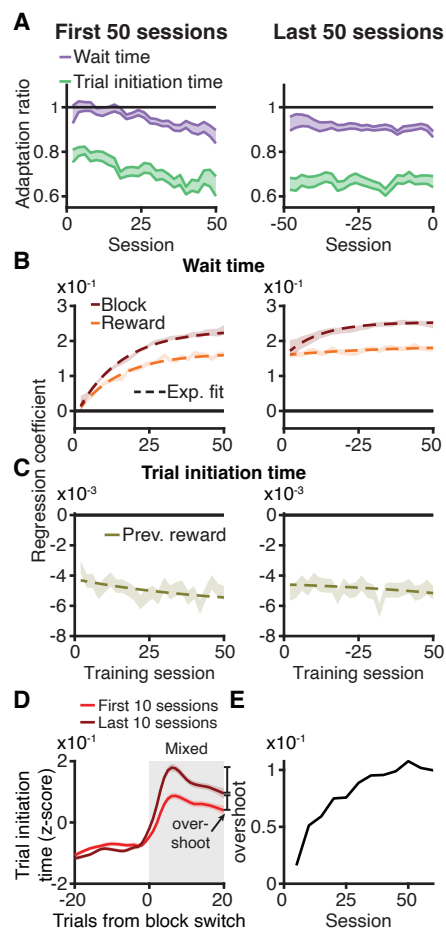


Figure 4: Block sensitivity for wait times requires structure learning. **A.** Wait time adaptation ratio (average wait time for 20 μL in high/low blocks) evolved over training, while trial initiation time ratio (average in high/low blocks) was below 1 on first session. **B.** Linear regression coefficients 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.

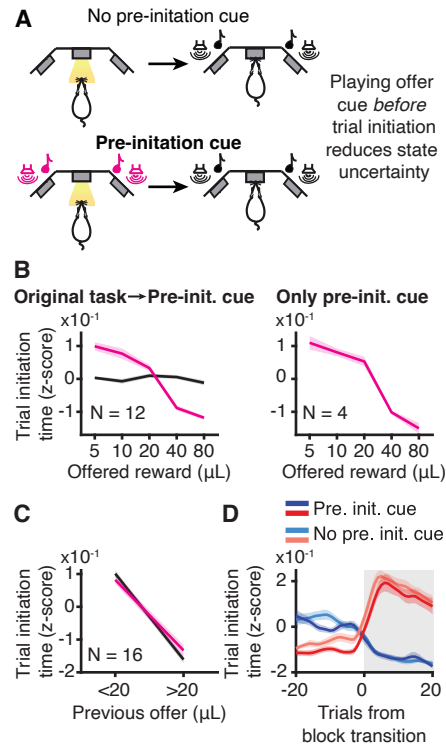


Figure 5: Value computations for motivation do not depend on state uncertainty. A. Schematic of pre-initiation cue experiment. **B.** Trial initiation time varied as a function of offered volume for rats that trained on the original task before transitioning to pre-initiation cue task (left) and for 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. 13/16 rats had $p < 0.05$, Wilcoxon Rank-sum test, $N = 16$. **D.** Trial initiation times in mixed blocks depended on previous block type in pre-initiation cue task.

188 inferred state have more uncertainty before they hear the reward offer. Model simulations of a
 189 Bayes' optimal observer did show that the reward offer reduced the uncertainty of subjective
 190 beliefs about the hidden state (comparing variance of prior to variance of posterior, $p \ll 0.001$,
 191 Wilcoxon sign-rank test).

192 To test this hypothesis, we modified the task so that some rats heard the reward cue before
 193 they initiated the trial, when the center light turned on; they heard the tone again at trial initia-
 194 tion, as in the standard task (Fig. 5A). Their trial initiation times became sensitive to the offered

195 reward (Fig. 5B). However, trial initiation times for 20 μ L in mixed blocks were still modu-
196 lated by the previous reward, consistent with the use of incrementally updated estimates of the
197 value of the environment within a block ($p < 0.05$ for 13/16 rats; Fig. 5C). Moreover, how
198 quickly they initiated trials in mixed blocks continued to depend on the previous block identity
199 (Fig. 5D). These data indicate that there may be something inherently retrospective about the
200 motivational decision to initiate a trial.

201 **Discussion**

202 We used high-throughput training to collect statistically powerful datasets and leverage in-
203 dividual variability across hundreds of animals. Consistent with previous work, rats adjusted
204 their behavior as we varied the richness of the environment in a way consistent with foraging
205 theories (14, 19, 35–37), and behavioral economic theories of reference dependence (38, 39).
206 Notably, we found that animals used multiple, parallel computations to estimate the richness of
207 the environment, and rapidly switched between these computations on single trials, indicating
208 that value computations vary on fine timescales (seconds). Our data are consistent with evidence
209 for multiple decision-making systems that rely on distinct neural circuits (40–43). While ani-
210 mals' decisions of how long to wait for rewards relied on hidden state inference, the decision of
211 when to initiate the trial was governed by a retrospective computation that calculated the value
212 of the environment as the running average of rewards. Reducing state uncertainty before the
213 trial did not change the value computations governing trial initiation times, suggesting that this
214 decision may be inherently retrospective, although influenced by subjective belief distributions
215 via a dynamic learning rate.

216 Recent work in psychology and machine learning has characterized how parallel value com-
217 putations might be combined (4–8, 29). For instance, in multi-step decision tasks, interaction
218 effects in regression models are thought to reflect the use of combined retrospective and inferen-

219 tial value estimates (9, 10), and hybrid strategies for computing values have been approximated
220 as a weighted average of retrospective and inference-based values (29). Our findings add to this
221 body of work. Instead of simply combining or averaging values that were computed in different
222 ways, rats seemed to coordinate their dynamics: changes in subjective beliefs about inferred
223 states acted as a gain on retrospective value learning rates.

224 It may be counterintuitive that the retrospective computation produced faster dynamics at
225 block transitions than hidden state inference (Fig. 1E,I). Two features of the models explain this
226 observation. First, the inferential model selects the block with the maximum posterior proba-
227 bility. This argmax operation nonlinearly thresholds whether changes in the posterior produce
228 changes in the inferred state. In contrast, the retrospective model’s estimate of the value of the
229 environment is directly influenced by graded, “subthreshold” changes in the posterior via the
230 dynamic learning rate. Subthreshold changes in the posterior necessarily precede changes that
231 cross threshold for inferring a state change. Second, the inferential model’s prior is recursive:
232 the posterior on one trial becomes the prior on the next trial. This means that the prior accu-
233 mulates information over trials to infer state changes, instead of making them instantaneously.
234 Indeed, individual differences in the informativeness of rats’ priors predicted the dynamics of
235 their inferred state changes (Fig. 3G).

236 The contextual effects we observed likely reflect efficient coding of value (12, 44–46). Ac-
237 cording to the efficient coding hypothesis, to represent stimuli efficiently, neurons should be
238 tuned to stimulus distributions that animals are most likely to encounter in the world (47).
239 Recent studies have shown that biases in value-based decision-making, including the contex-
240 tual effects observed here, reflect efficient value coding (12, 44, 45). Previous studies exam-
241 ined how neurons “adapted” to reward or stimulus distributions over blocks of trials or ses-
242 sions, implying gradual, experience-dependent adjustments in behavioral sensitivity and neural
243 tuning (12, 48, 49). Our findings suggest that if animals have learned the reward or stimulus

244 distributions associated with a particular state, they can condition their subjective value rep-
245 resentations on that inferred state, perhaps via discrete, state-dependent adjustments in neural
246 sensitivity (50). A major future question is how multi-regional neural circuits represent be-
247 lief distributions for hidden state inference, and condition rapid adjustments in efficient neural
248 representations of value on inferred states.

249 **Methods**

250 **Subjects**

251 A total of 240 Long-evans rats (148 male, 92 female) between the ages of 6 and 24 months
252 were used for this study (*Rattus norvegicus*). The Long-evans cohort also included ADORA2A-
253 Cre ($N=10$), ChAT-Cre ($N=2$), DRD1-Cre ($N=3$), and TH-Cre ($N=12$). Animal use procedures
254 were approved by the New York University Animal Welfare Committee (UAWC #2021-1120)
255 and carried out in accordance with National Institutes of Health standards.

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

262 **Behavioral training**

263 Rats were trained in a high-throughput behavioral facility in the Constantinople lab using
264 a computerized training protocol. They were trained in custom operant training boxes with
265 three nose ports. Each nose port was 3-D printed, and the face was protected with an epoxied
266 stainless steel washer (McMaster-Carr #92141A056). All ports contained a visible light emit-

267 ting diode (LED; Digikey #160-1850-ND), and an infrared LED (Digikey #365-1042-ND) and
268 infrared photodetector (Digikey #365-1615-ND) that enabled detection of when a rat broke the
269 infrared beam with its nose. Additionally, the side ports contained stainless steel lick tubes
270 (McMaster-Carr #8988K35, cut to 1.5mm) that delivered water via solenoid valves (Lee Com-
271 pany #LHDA1231115H). There was a speaker mounted above each side port that enabled de-
272 livery of stereo sounds (Bohlender Graebener). The behavioral task was instantiated as a finite
273 state machine on an Arduino-based behavioral system with a Matlab interface (Bpod State Ma-
274 chine r2, Sanworks), and sounds were delivered using a low-latency analog output module
275 (Analog Output Module 4ch, Sanworks) and stereo amplifier (Lepai LP-2020TI).

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

284 **Sound Calibration**

285 We calibrated sounds using a hand-held Precision Sound Level Meter with a 1/2" micro-
286 phone (Bruel & Kjaer, Type 2250). The microphone was calibrated with a sound level calibrator
287 (Bruel & Kjaer, Type 4230). Tones of different frequencies (1, 2, 4, 8, 16kHz) were presented
288 for 10 seconds each; these tones were selected because they are in the trough of the behav-
289 ioral audiogram for rats (51). They are also on a logarithmic scale and thus should be equally
290 discriminable to the animals. We adjusted the auditory gain in software for each frequency

291 stimulus to match the sound pressure level to 70dB in the rig, measured when the microphone
292 was proximal to the center poke.

293 **Task Logic**

294 LED illumination from the center port indicated that the animal could initiate a trial by
295 poking its nose in that port - upon trial initiation the center LED turned off. While in the center
296 port, rats needed to maintain center fixation for a duration drawn uniformly from [0.8, 1.2]
297 seconds. During the fixation period, a tone played from both speakers, the frequency of which
298 indicated the volume of the offered water reward for that trial [1, 2, 4, 8, 16kHz, indicating
299 5, 10, 20, 40, 80 μ L rewards]. Following the fixation period, one of the two side LEDs was
300 illuminated, indicating that the reward might be delivered at that port; the side was randomly
301 chosen on each trial. This event (side LED ON) also initiated a variable and unpredictable delay
302 period, which was randomly drawn from an exponential distribution with mean = 2.5 seconds.
303 The reward port LED remained illuminated for the duration of the delay period, and rats were
304 not required to maintain fixation during this period, although they tended to fixate in the reward
305 port. When reward was available, the reward port LED turned off, and rats could collect the
306 offered reward by nose poking in that port. The rat could also choose to terminate the trial
307 (opt-out) at any time by nose poking in the opposite, un-illuminated side port, after which a
308 new trial would immediately begin. On a proportion of trials (15-25%), the delay period would
309 only end if the rat opted out (catch trials). If rats did not opt-out within 100s on catch trials, the
310 trial would terminate.

311 The trials were self-paced: after receiving their reward or opting out, rats were free to
312 initiate another trial immediately. However, if rats terminated center fixation prematurely, they
313 were penalized with a white noise sound and a time out penalty (typically 2 seconds, although
314 adjusted to individual animals). Following premature fixation breaks, the rats received the same

315 offered reward, in order to disincentivize premature terminations for small volume offers.

316 We introduced semi-observable, hidden-states in the task by including uncued blocks of
317 trials with varying reward statistics (12): high and low blocks, which offered the highest three
318 or lowest three rewards, respectively, and were interspersed with mixed blocks, which offered
319 all volumes. There was a hierarchical structure to the blocks, such that high and low blocks
320 alternated after mixed blocks (e.g., mixed-high-mixed-low, or mixed-low-mixed-high). The first
321 block of each session was a mixed block. Blocks transitioned after 40 successfully completed
322 trials. Because rats prematurely broke fixation on a subset of trials, in practice, block durations
323 were variable.

324 **Criteria for including behavioral data**

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

339 **Shaping**

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

345 Stages 2-5 progressively introduced the full set of reward volumes and corresponding au-
346 ditory cues. Rats continued to receive deterministic rewards with no delay after maintaining a
347 1 second center poke. Each stage added one additional reward that could be selected on each
348 trial- stage 2 added 40 μL , stage 3 added 5 μL , stage 4 added 80 μL , and stage 5 added 10 μL .
349 Each stage progressed after 400 successfully completed trials. All subsequent stages used all 5
350 reward volumes.

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

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

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

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

364 **Training for male and female rats**

365 We collected data from both male and female rats (160 male, 114 female). Male and female
366 rats were trained in identical behavioral rigs with the same shaping procedure described above.
367 Early cohorts of female rats experienced the same reward set as the males. However, female
368 rats are smaller, and they consumed less water and performed substantially fewer trials than
369 the males. Therefore, to obtain sufficient behavioral trials from them, reward offers for female
370 rats were slightly reduced while maintaining the logarithmic spacing: [4, 8, 16, 32, 64 μ L]. For
371 behavioral analysis, reward volumes were treated as equivalent to the corresponding volume for
372 the male rats (e.g., 16 μ L trials for female rats were treated the same as 20 μ L trials for male
373 rats). The auditory tones were identical to those used for male rats. We did not observe any
374 significant differences between the male and female rats, in terms of the degree of wait time
375 adaptation, and the qualitative nature of behavioral dynamics at block transitions (fig. S6).

376 We tracked most female rats' stages in the estrus cycle using vaginal cytology, with vaginal
377 swabs collected immediately after each session using a cotton-tipped applicator first dipped

378 in saline. Samples were smeared onto a clean glass slide and visually classified under a light
379 microscope. For the current study, data from female rats was averaged across all stages of the
380 estrus cycle.

381 **Behavioral models**

382 We developed separate behavioral models to describe rat's wait time and trial initiation time
383 data. Both wait time and trial initiation time should depend on the value of the environment. For
384 the wait time data, we adapted a model from (1) which described the optimal wait time, WT_{opt} ,
385 in terms of the value of the environment (i.e., the opportunity cost), the delay distribution, and
386 the catch probability (i.e., the probability of the trial being unrewarded). Given an exponential
387 delay distribution, we defined the optimal wait time as

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

388 where τ is the time constant of the exponential delay distribution, C is the probability of reward
389 (1-catch probability), R is the reward on that trial, κ is the opportunity cost, and D is a scaling
390 parameter. For the trial initiation time, we adapted a model from (2) which describes the optimal
391 trial initiation time, TI_{opt} , given the value of the environment, κ , as

$$TI_{\text{opt}} = \frac{D}{\kappa},$$

392 where D is a scale parameter.

393 We initially evaluated two different ways of calculating the value of the environment for
394 these models, which are shared between the wait time and trial initiation time models: a retro-
395 spective and inferential model (see below). We assumed independent log-normal noise for each
396 trial, with a constant variance of 8 seconds for the wait time model and 4 seconds for the trial
397 initiation time model. The log-normal noise model outperformed alternative noise models, such

398 as gamma and ex-Gaussian noise. The noise variance terms were selected from a grid search
399 using data from a subset of animals.

400 **Inferential model**

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

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

405 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
406 probability of the reward for each block, for example,

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

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

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

409 $P(B_t | B_{t-1})$, referred to as the "hazard rate," incorporates knowledge of the task structure,
410 including the block length and block transition probabilities. For example,

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

411 where $H_0 = 1/40$, to reflect the block length. The model assumed a flat block hazard rate
412 for the following reasons. (1) Since animals broke center fixation on a subset of trials, the
413 actual block duration was highly variable. Based on the distributions of experienced block
414 durations, it is unlikely that rats would have learned a perfect step function hazard rate. (2) The

415 blocks spanned several to tens of minutes, making it unlikely that rats would keep a running
416 tally of trials on such long timescales. (3) Gradual changes in wait times at block transitions
417 are not consistent with the use of a veridical step-function hazard rate. (4) We considered an
418 alternative parameterization in which the veridical step function hazard rate was blurred with a
419 Gaussian, but this would have required a number of nontrivial design choices, such as whether
420 the trial counter should be reset after “misinferred” block transitions, regardless of when they
421 occurred in the actual block. (5) Wait times reflected misinferred blocks based on a constant
422 block hazard rate (fig. S4), suggesting that this simplification was a reasonable approximation
423 of the inference process. Including H_0 as an additional free parameter did not improve the
424 performance of the wait time model evaluated on held-out test data in a subset of rats (data not
425 shown), so H_0 was treated as a constant term.

426 **Belief state model**

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

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

431 **Inferential model with lambda parameter**

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

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

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

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

439 **Retrospective model**

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

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

443 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
444 reward prediction error (RPE), and α_t is a dynamic learning rate given by $\alpha_t = G \cdot \alpha_0$. In order
445 to capture the dynamics of the trial initiation times around block transitions, we included a gain
446 term, G_t on the learning rate, which is inversely related to the trial-by-trial change in the mixed
447 block probability from by the inferential model, given by

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

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

453 The dynamic learning rate we implemented is consistent with previous work showing that
454 humans and animals can adjust their learning rates depending on the volatility and uncertainty
455 in the environment (25–27). Other models using either (1) a single, static learning rate ($G = 1$),

456 or (2) a dynamic learning rate where the gain term was the unsigned reward prediction error on
457 that trial ($G = |\delta_t|$) were unable to capture the observed trial initiation time dynamics at block
458 transitions (fig. S2).

459 **Fitting and evaluating models**

460 We used MATLAB's constrained minimization function, `fmincon`, to minimize the sum of
461 the negative log likelihoods with respect to the model parameters. 5-10 random seeds were used
462 in the maximum likelihood search for each rat; parameter values with the maximum likelihood
463 of these seeds were deemed the best fit parameters. Before fitting to rat's data, we confirmed
464 that our fitting procedure was able to recover generative parameters (fig. S7). When evaluating
465 model performance fit to rat data, we performed 5-fold cross-validation and evaluated the pre-
466 dictive power of the model on the held-out test sets. To compare the different models, we used
467 Bayesian Information Criterion (BIC), $BIC = \log(n) \cdot k + 2 \cdot nLL$, where n is the number of
468 trials, k is the number of parameters, and nLL is the negative log-likelihood of the best-fit model
469 evaluated on all data. We confirmed the model comparison by also comparing Akaike Informa-
470 tion Criterion (AIC) and cross-validated negative log-likelihood, which gave similar results to
471 BIC.

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

478 **Statistical analyses**

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

480 For all analyses, we removed wait times that were one standard deviation above the pooled-
481 session mean. When assessing whether a rat's wait time differed by blocks, we compared each
482 rat's wait time on catch trials offering 20 μ L in high and low blocks using a non-parametric
483 Wilcoxon rank-sum test, given that the wait times are roughly log-normally distributed. We
484 defined each rat's wait time ratio as the average wait time on 20 μ L catch trials in high blocks/low
485 blocks. For trial initiation times, we compared all trial initiation times for each block, again
486 using a non-parametric Wilcoxon rank-sum test. We defined each rat's trial initiation time ratio
487 as the average trial initiation time in high blocks/low blocks.

488 Trial initiation times were bimodally distributed, with the different modes reflecting whether
489 previous trials were rewarded or not. Unrewarded trials included opt-out trials and trials where
490 rats prematurely terminated center fixation ("violation trials"). Analyzing these trial types sep-
491 arately showed that trial initiation times following unrewarded trials were modulated by blocks
492 in a similar pattern as the wait times, with rats initiating trials more quickly in high compared
493 to low blocks (fig. S8). While we used all behavioral trials for analyses of trial initiation times
494 throughout the manuscript, we note that trial initiation times following rewarded trials exhibited
495 a different pattern (fig. S8), consistent with previous studies showing that response outcomes
496 gate behavioral strategies (52, 53). Specifically, following rewarded trials, there was a weak
497 positive correlation between reward magnitude and trial initiation time, in contrast to the strong
498 negative correlation we observed following unrewarded trials. We interpret the positive corre-
499 lation as potentially reflecting micro-satiety effects. However, as these effects were weak, most
500 of the variance in the trial initiation times were driven by those following unrewarded trials.

501 To assess block effects across the population, we first z-scored each rat's wait time on all
502 catch trials and trial initiation time on all trials. For wait times, we computed the average z-

503 scored wait time on catch trials offering 20 μ L in high and low blocks for each rat, and compared
504 across the population using a paired Wilcoxon sign-rank test. Similarly for trial initiation times,
505 we averaged all z-scored trial initiation times for high and low blocks for each rat, and compared
506 across the population using a paired Wilcoxon sign-rank test.

507 **Block transition dynamics**

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

516 For each transition type, we averaged the Δ z-scored wait times and trial initiation times
517 based on their distance from a block transition, including violation trials (e.g., averaged all wait
518 times four trials before a block transition). Finally, for each block transition type, we smoothed
519 the average curve for each rat using a 10-point moving average, before averaging over rats.

520 When comparing block transition dynamics in rats with different quality priors, specifically
521 from mixed blocks to high or low, we chose rats in the top or bottom 40th percentile of fit λ 's
522 and averaged each group's block transition dynamics for both wait time and trial initiation time.
523 We then normalized each curve by subtracting the average wait or initiation time value before
524 the block transition. To compare the normalized dynamics of each group, we fit 3-parameter
525 logistic functions of the following form:

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

526 to the behavioral curves and compared the three parameters: D (the upper asymptote), C (the
527 inverse temperature), and x_0 (x -value of the the sigmoid's midpoint). To determine significance
528 for our observed differences, we performed a non-parametric shuffle test. We generated null
529 distributions on differences in the fit parameters by shuffling the labels of the upper and lower
530 percentile λ rats, refitting the logistic to the new shuffled groups' average dynamic curves, and
531 comparing the fit parameters 500 times. We then used these null distributions to calculate p-
532 values for the observed differences in parameters: the area under this distribution evaluated
533 at the actual difference of parameter values (between high and low λ rats) was treated as the
534 p-value.

535 **Trial history effects**

536 To assess wait time sensitivity to previous offers, we focused on 20 μL catch trials in mixed
537 blocks only. We z-scored the wait times of these trials separately. Next, we averaged wait times
538 depending on whether the previous offer was greater than or less than 20 μL . For trial initiation
539 times, we used all 20 μL trials in mixed blocks. We averaged z-scored trial initiation times
540 depending on whether the previous offer was greater or less than 20 μL . For both wait time
541 and trial initiation time, we defined the sensitivity to previous offers as the difference between
542 average wait time (trial initiation time) for trials with a previous offer less than 20 μL and
543 trials with a previous offer greater than 20 μL . We compared wait time and trial initiation time
544 sensitivity to previous offers across rats using a paired Wilcoxon signed-rank test.

545 To capture longer timescale sensitivity across rewards, we regressed previous rewards against
546 wait time and trial initiation time. We focused only on mixed blocks. Additionally, we lin-
547 earized the rewards by taking the binary logarithm of each reward ($\log_2(\text{reward})$). For wait
548 time, we z-scored wait times for catch trials in mixed blocks. Then, we regressed wait times
549 on these trials against the current offer and previous 9 $\log_2(\text{reward})$ offers, including violation

550 trials, along with a constant offset term. Reward offers from a different block (e.g., a previous
551 high block) were given NaN values. For trial initiation times, we again z-scored for mixed
552 block trials only. Then, we regressed against the previous 9 $\log_2(\text{reward})$ offers, not including
553 the current trial, along with a constant offset. Additionally, we set the reward for violation and
554 catch trials to 0, since rats do not receive a reward on these trials.

555 For both wait time and trial initiation time, we used Matlab's builtin regress function to
556 perform the regression. With the coefficients, we found the first non-significant coefficient (co-
557 efficient that whose 95% confidence interval contained 0), and set that coefficient and all fol-
558 lowing coefficients to 0. Finally, we fit a negative exponential decay curve, $y = D \exp -x/\tau$,
559 to each rat's previous trial coefficients (that is, only the previous 9 trial coefficients) for both
560 wait time and trial initiation time and reported the time constant of the exponential decay (tau)
561 for each. If all previous trial coefficients were equal to 0 (as was the case for a vast majority of
562 the wait time coefficients), the time constant was reported as NaN. We correlated wait time re-
563 gression time-constants and trial initiation time regression time-constants using Matlab's builtin
564 corr function.

565 **Learning Dynamics**

566 To assess learning dynamics, we included all sessions after stage 8, not just the sessions
567 that passed criteria for inclusion (above). Because of data limitations examining each session
568 individually (e.g., not every session included both a high and low block), we grouped subse-
569 quent sessions into pairs (i.e., we grouped sessions 1 and 2, sessions 3 and 4, etc.). For each
570 session-pair, we calculated the wait time and trial initiation time ratios as above. To assess the
571 emergence of block effects on wait time data, we regressed wait time for each session against
572 both the current reward and a categorical variable representing the current block identity (1 =
573 low block, 2 = mixed block, 3 = high block). To assess the emergence of previous trial effects

574 on trial initiation time, we regressed trial initiation time for each sessions against the previous
575 reward. We smoothed each regression coefficient over sessions using a 5-session moving av-
576 erage. Finally, we set outlier coefficients (3 scaled median absolute deviations away from a
577 5-point moving median, using Matlab's builtin *isoutlier* function) to NaN. Finally, we averaged
578 regression coefficients over sessions across rats.

579 **Pre-initiation cue task**

580 To modulate the subjective uncertainty in the rat's estimate of state (block) before trial
581 initiation time, we ran a subset of rats on a variation of the task where we cued reward offer
582 before rats initiated a trial ($N=16$). All other aspects of the task remained identical: reward offer
583 cued played again after the rat initiated the trial, rats waited uncued exponentially-distributed
584 delays for rewards, etc. We included both rats that initially trained on the original task before
585 switching to the pre-initiation cue task ($N = 12$), as well as rats who were trained only on the
586 pre-initiation cue task ($N = 4$). To allow the rats who had started on the original task time to
587 adjust to the new task, we only included data after 30 pre-initiation cue sessions. For the rats
588 who were exclusively trained on the pre-initiation cue task, we included all stage 8 sessions.
589 For all rats, we did not exclude sessions using the wait time criteria (see above).

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

References

- 598
- 599 1. A. Dickinson, B. Balleine (2002).
 - 600 2. Y. Niv, D. Joel, P. Dayan, *Trends in Cognitive Sciences* **10**, 375 (2006).
 - 601 3. R. S. Sutton, A. G. Barto, *Reinforcement learning: An introduction* (MIT press, 2018).
 - 602 4. G. Pezzulo, F. Rigoli, F. Chersi, *Frontiers in Psychology* **4**, 92 (2013).
 - 603 5. S. J. Gershman, E. J. Horvitz, J. B. Tenenbaum, *Science (New York, N.Y.)* **349**, 273 (2015).
 - 604 6. P. Dayan, *Current Opinion in Neurobiology* **22**, 1068 (2012).
 - 605 7. M. Keramati, P. Smittenaar, R. J. Dolan, P. Dayan, *Proceedings of the National Academy*
606 *of Sciences of the United States of America* **113**, 12868 (2016).
 - 607 8. N. D. Daw, Y. Niv, P. Dayan, *Nature Neuroscience* **8**, 1704 (2005). Number: 12 Publisher:
608 Nature Publishing Group.
 - 609 9. N. D. Daw, S. J. Gershman, B. Seymour, P. Dayan, R. J. Dolan, *Neuron* **69**, 1204 (2011).
 - 610 10. W. Kool, S. J. Gershman, F. A. Cushman, *Psychological Science* **28**, 1321 (2017).
 - 611 11. A. Lak, *et al.*, *Neuron* **84**, 190 (2014).
 - 612 12. M. W. Khaw, P. W. Glimcher, K. Louie, *Proceedings of the National Academy of Sciences*
613 **114**, 12696 (2017). Publisher: Proceedings of the National Academy of Sciences.
 - 614 13. A. P. Steiner, A. D. Redish, *Nature neuroscience* **17**, 995 (2014).
 - 615 14. E. L. Charnov, *Theoretical Population Biology* **9**, 129 (1976).
 - 616 15. D. W. Stephens, J. R. Krebs, *Foraging theory* (Princeton university press, 2019).

- 617 16. F. Rigoli, *Cognition* **192**, 104034 (2019).
- 618 17. R. Shadmehr, A. A. Ahmed, *Vigor: Neuroeconomics of movement control* (MIT Press,
619 2020).
- 620 18. C. F. Flaherty, *Animal Learning & Behavior* **10**, 409 (1982).
- 621 19. S. M. Constantino, N. D. Daw, *Cognitive, Affective, & Behavioral Neuroscience* **15**, 837
622 (2015).
- 623 20. P. Verтеchi, *et al.*, *Neuron* **106**, 166 (2020).
- 624 21. R. C. Wilson, Y. K. Takahashi, G. Schoenbaum, Y. Niv, *Neuron* **81**, 267 (2014).
- 625 22. J. L. Jones, *et al.*, *Science* **338**, 953 (2012).
- 626 23. H. Davis, *Journal of Comparative Psychology* **106**, 342 (1992).
- 627 24. C. Gallistel, T. A. Mark, A. P. King, P. Latham, *Journal of experimental psychology: Ani-
628 mal behavior processes* **27**, 354 (2001).
- 629 25. T. E. Behrens, M. W. Woolrich, M. E. Walton, M. F. Rushworth, *Nature neuroscience* **10**,
630 1214 (2007).
- 631 26. M. R. Nassar, *et al.*, *Nature neuroscience* **15**, 1040 (2012).
- 632 27. C. D. Grossman, B. A. Bari, J. Y. Cohen, *Current Biology* **32**, 586 (2022).
- 633 28. S. J. Gershman, Y. Niv, *Current opinion in neurobiology* **20**, 251 (2010).
- 634 29. B. Miranda, W. M. N. Malalasekera, T. E. Behrens, P. Dayan, S. W. Kennerley, *PLoS
635 computational biology* **16**, e1007944 (2020).
- 636 30. E. S. Bromberg-Martin, M. Matsumoto, H. Nakahara, O. Hikosaka, *Neuron* **67**, 499 (2010).

- 637 31. M. Van Der Meer, Z. Kurth-Nelson, A. D. Redish, *The Neuroscientist* **18**, 342 (2012).
- 638 32. N. Drummond, Y. Niv, *Current biology : CB* **30** (2020). Publisher: Curr Biol.
- 639 33. B. W. Balleine, A. Dickinson, *International Journal of Comparative Psychology* **18** (2005).
- 640 34. M. Keramati, A. Dezfouli, P. Piray, *PLOS Computational Biology* **7**, e1002055 (2011).
641 Publisher: Public Library of Science.
- 642 35. E. Freidin, A. Kacelnik, *Science* **334**, 1000 (2011).
- 643 36. B. Y. Hayden, J. M. Pearson, M. L. Platt, *Nature neuroscience* **14**, 933 (2011).
- 644 37. N. Kolling, T. E. Behrens, R. B. Mars, M. F. Rushworth, *Science* **336**, 95 (2012).
- 645 38. D. Kahneman, A. Tversky, *Handbook of the fundamentals of financial decision making:*
646 *Part I* (World Scientific, 2013), pp. 99–127.
- 647 39. B. Kőszegi, M. Rabin, *The Quarterly Journal of Economics* **121**, 1133 (2006).
- 648 40. A. D. Redish, N. W. Schultheiss, E. C. Carter, *Current Topics in Behavioral Neurosciences*
649 **27**, 313 (2016).
- 650 41. P. Dayan, Y. Niv, B. Seymour, N. D. Daw, *Neural networks* **19**, 1153 (2006).
- 651 42. B. W. Balleine, *Neuron* **104**, 47 (2019).
- 652 43. B. M. Sweis, *et al.*, *Science (New York, N.Y.)* **361**, 178 (2018).
- 653 44. R. Polanía, M. Woodford, C. C. Ruff, *Nature neuroscience* **22**, 134 (2019).
- 654 45. K. Louie, P. W. Glimcher, *Annals of the New York Academy of Sciences* **1251**, 13 (2012).
- 655 46. A. Tymula, P. Glimcher, *Available at SSRN 2783638* (2021).

- 656 47. H. B. Barlow, *et al.*, *Sensory communication* **1**, 217 (1961).
- 657 48. C. Padoa-Schioppa, *Journal of Neuroscience* **29**, 14004 (2009).
- 658 49. A. I. Weber, K. Krishnamurthy, A. L. Fairhall, *Annual review of vision science* **5**, 427
659 (2019).
- 660 50. S. Kobayashi, O. P. de Carvalho, W. Schultz, *Journal of Neuroscience* **30**, 534 (2010).
- 661 51. H. E. Heffner, R. S. Heffner, C. Contos, T. Ott, *Hearing research* **73**, 244 (1994).
- 662 52. A. Hermoso-Mendizabal, *et al.*, *Nature communications* **11**, 1 (2020).
- 663 53. K. Iigaya, M. S. Fonseca, M. Murakami, Z. F. Mainen, P. Dayan, *Nature communications*
664 **9**, 1 (2018).

665 **Acknowledgments**

666 We thank Paul Glimcher, Catherine Hartley, Roozbeh Kiani, Kenway Louie, Kevin Miller,
667 Cristina Savin and members of the Constantinople lab for helpful discussions. We thank Madori
668 Spiker, Daljit Kaur, Mitzi Adler-Wachter, Royall McMahon Ward, and Luke Chen for animal
669 training.

670 Funding: This work was supported by a K99/R00 Pathway to Independence Award (R00MH111926),
671 an Alfred P. Sloan Fellowship, a Klingenstein-Simons Fellowship in Neuroscience, an NIH Di-
672 rector's New Innovator Award (DP2MH126376), an NSF CAREER Award, R01MH125571,
673 and a McKnight Scholars Award to CMC. AM was supported by 5T90DA043219, 5T32MH019524,
674 and F31MH130121.

675 **Supplementary materials**

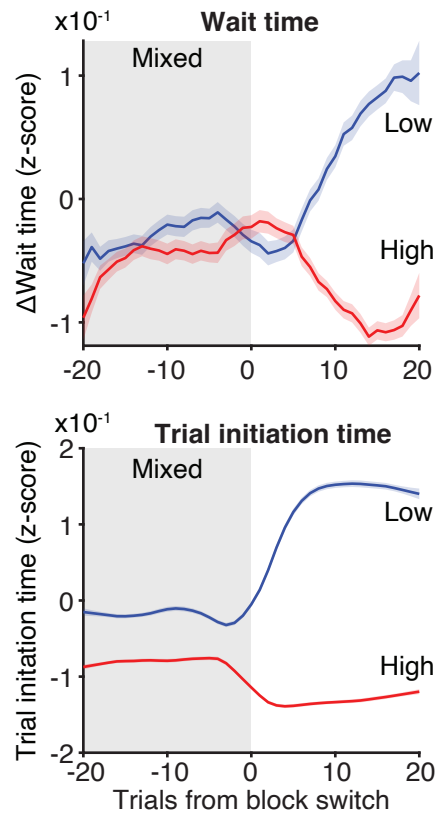
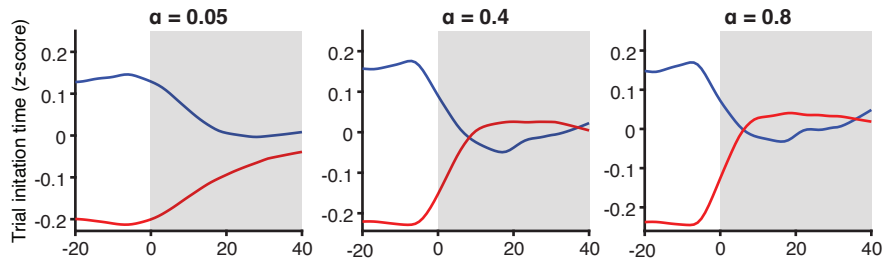


Fig. S1: Dynamics of wait times (top) and trial initiation times (bottom) at transitions from mixed to high (red) or low (blue) blocks. Long timescale effects were observable for trial initiation times but not wait times: even by the end of the mixed block, how quickly rats initiated trials depended on the previous block identity.

Vanilla learning rate model: a single, static learning rate



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

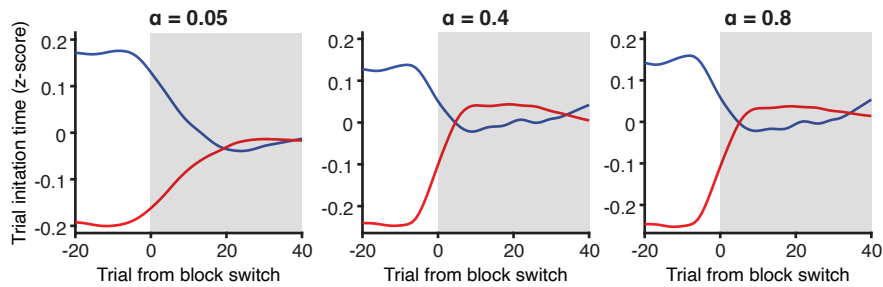


Fig. S2: Alternative retrospective models fail to capture both fast and slow trial initiation time dynamics at block transitions. Trial initiation time model transitions from low (blue) or high (red) blocks to mixed blocks. Top: A “vanilla” learning rate model with a single, static learning rate. Bottom: a dynamic learning rate model where learning rate gain is equal to the unsigned RPE of that trial.

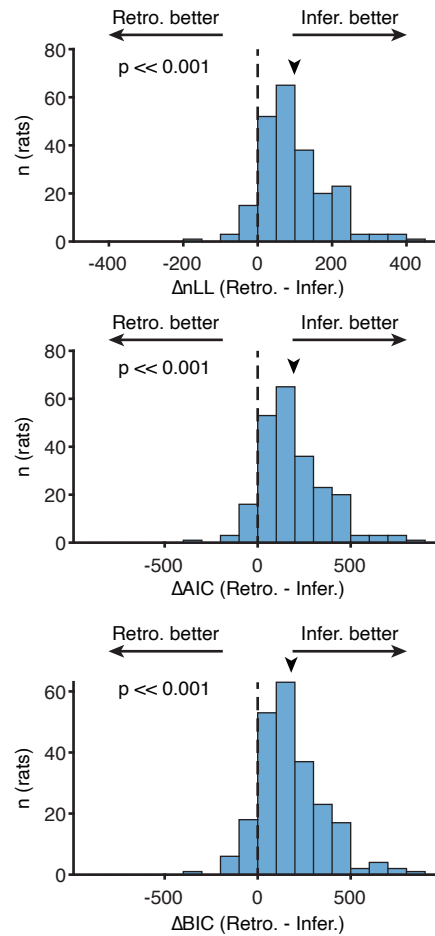


Fig. S3: All forms of model comparison favor the inferential over retrospective model as a description of rat's wait times. Wilcoxon signed-rank test, $N = 240$

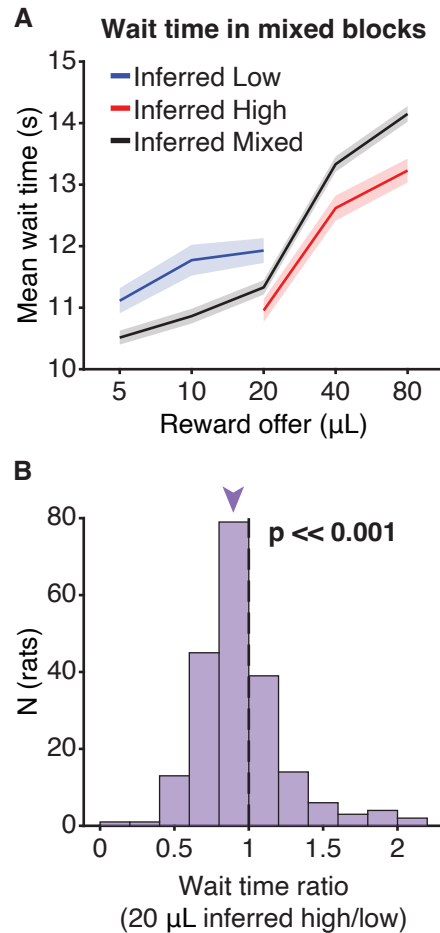


Fig. S4: Inferential model identifies mistaken inferences during mixed blocks across rats. **A.** Average wait time curves conditioned by model-inferred block in mixed blocks only in held-out test set across rats. **B.** Wait time ratio (wait time on 20 μL inferred high/low trials) is modulated by inferred block ($p \ll 0.001$, Wilcoxon Signed-rank test, $N = 240$)

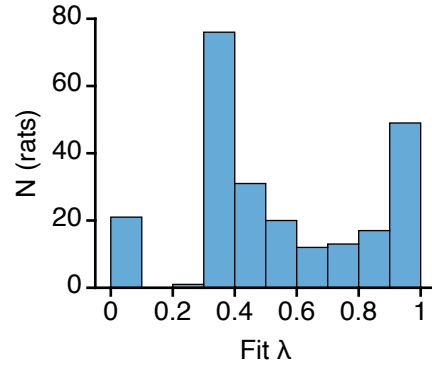


Fig. S5: **Sub-optimal inferential model with lambda.** Distribution of λ fit over rats.

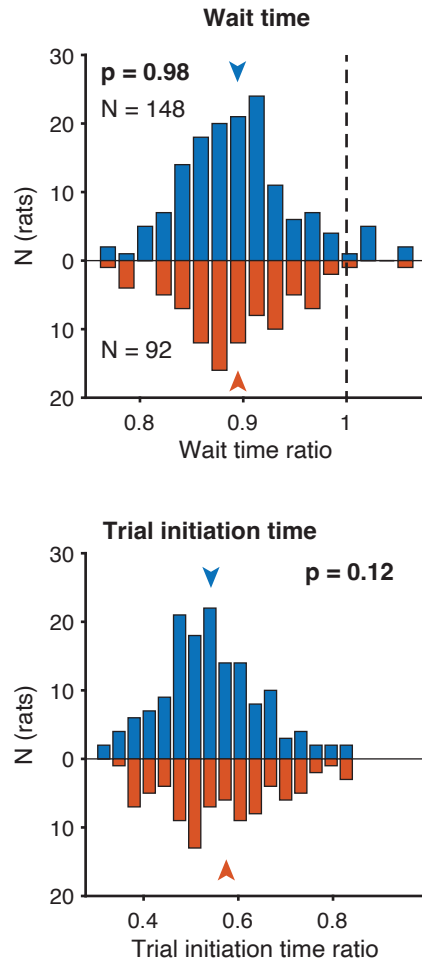


Fig. S6: **Males and females have comparable wait time ratios (top) and trial initiation time ratios (bottom).** Wait time $p = 0.98$, Wilcoxon Rank-sum test, N = 148 males, 92 females. Trial initiation time $p = 0.12$, Wilcoxon Rank-sum test, N = 148 males, 92 females.

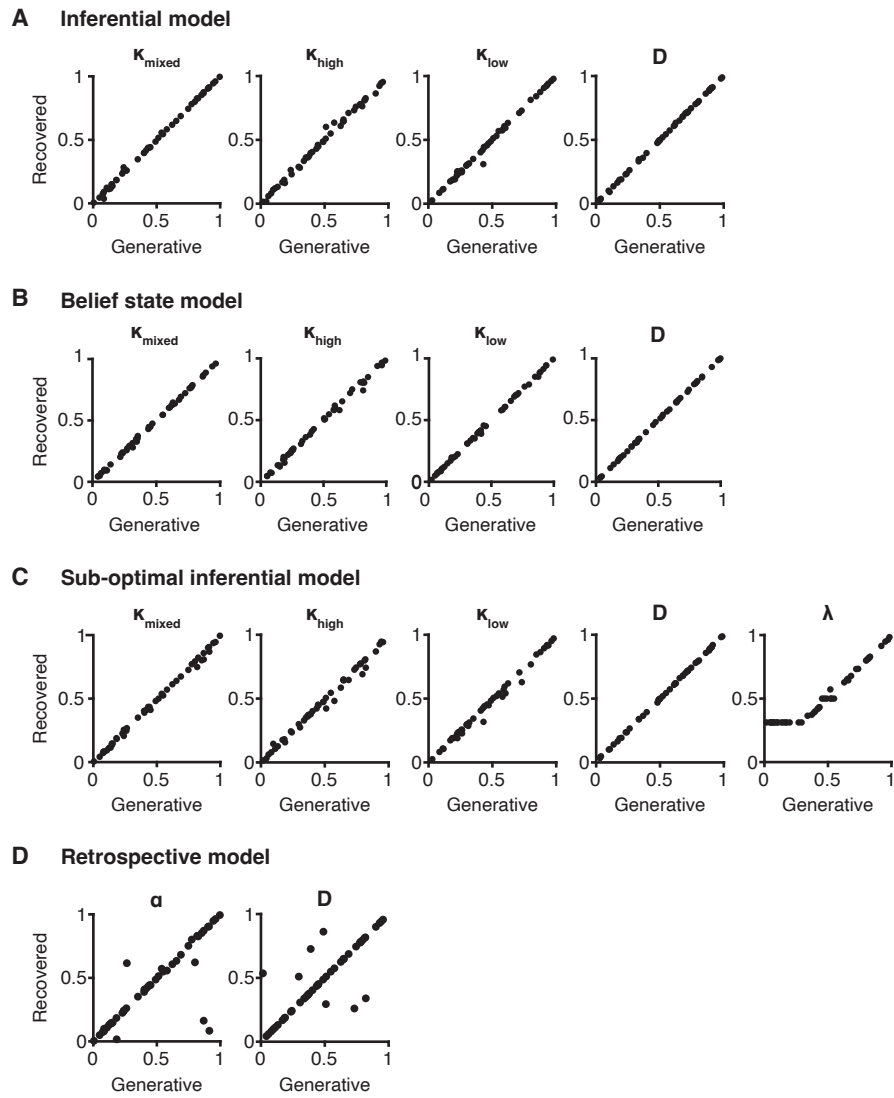


Fig. S7: Models are able to recover generative parameters. $N = 48$ random parameter sets.

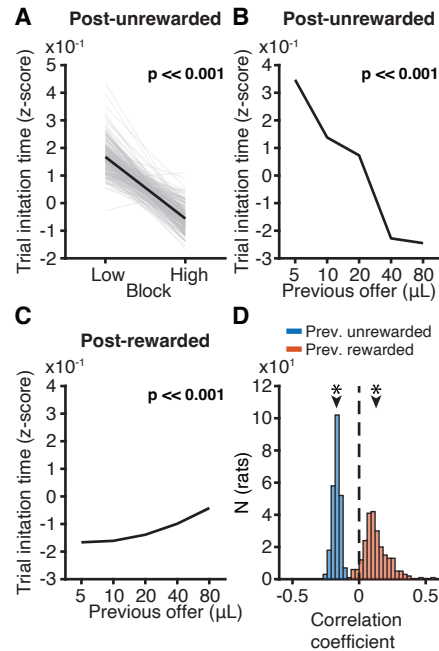


Fig. S8: **Trial initiation times depend on previous trial outcome.** **A.** Trial initiation times after unrewarded trials were faster in high blocks compared to low blocks ($p << 0.001$, Wilcoxon Signed-rank test, $N = 240$). **B.** Trial initiation times after unrewarded trials were negatively modulated by the previous offer in mixed blocks (linear regression slope < 0 , $p << 0.001$, Student's t-test, $N = 240$). **C.** Trial initiation times were slightly slower following larger volume rewarded trials (linear regression slope > 0 , $p < 0.05$, Student's t-test, $N = 240$). **D.** Correlation coefficient between previous reward offer and trial initiation time across rats differed both in sign and magnitude following rewarded and unrewarded trials ($p << 0.001$, Wilcoxon signed-rank test, $N = 240$).