

A cognitive process model captures near-optimal confidence-guided waiting in rats

J Tyler Boyd-Meredith^{1,2}, Alex T Piet³, Chuck D Kopec¹, and Carlos D Brody^{1,4,*}

¹Princeton Neuroscience Institute, Princeton University, Princeton, United States.

²Sainsbury Wellcome Centre, University College London, London, UK.

³Allen Institute, Seattle, Washington, United States.

⁴Howard Hughes Medical Institute, Princeton University, Princeton, United States.

*correspondence should be addressed to: Carlos D Brody (brody@princeton.edu)

1 Abstract

2 Rational decision-makers invest more time pursuing rewards they are more confident
3 they will eventually receive. A series of studies have therefore used willingness to wait
4 for delayed rewards as a proxy for decision confidence. However, interpretation of
5 waiting behavior is limited because it is unclear how environmental statistics influence
6 optimal waiting, and how sources of internal variability influence subjects' behavior. We
7 trained rats to perform a confidence-guided waiting task, and derived expressions for
8 optimal waiting that make relevant environmental statistics explicit, including travel
9 time incurred traveling from one reward opportunity to another. We found that rats
10 waited longer than fully optimal agents, but that their behavior was closely matched
11 by optimal agents with travel times constrained to match their own. We developed
12 a process model describing the decision to stop waiting as an accumulation to bound
13 process, which allowed us to compare the effects of multiple sources of internal variability
14 on waiting. Surprisingly, although mean wait times grew with confidence, variability
15 did not, inconsistent with scalar invariant timing, and best explained by variability in
16 the stopping bound. Our results describe a tractable process model that can capture
17 the influence of environmental statistics and internal sources of variability on subjects'
18 decision process during confidence-guided waiting.

19 Introduction

20 A decision maker's estimate of the probability that a decision is correct given the evidence
21 is referred to as decision confidence^{1,2}. Confidence is critical for learning improvements
22 in decision policy in response to feedback³, for deciding whether to act or gather more
23 information⁴, and for determining how long to wait for an expected outcome before
24 seeking reward elsewhere⁵. However, study of the neural underpinnings of confidence has
25 been limited by the difficulty of measuring confidence in animal subjects.

26 Recent work^{6,5} has developed a promising assay of decision confidence by asking how
27 long subjects are willing to wait for a rewarding outcome after a decision before moving on
28 to another reward opportunity. This temporal post-decision wagering paradigm has been
29 used to study confidence in olfactory^{5,7}, auditory⁷, visual⁸, and mnemonic⁹ decisions,

30 and has been used to study hallucinations¹⁰. Temporal post-decision wagers have also
31 been used to probe learning about environmental reward statistics¹¹.

32 The premise of these confidence-guided waiting studies is that reward-rate-maximizing
33 (“optimal”) agents are willing to wait longer for delayed rewards when they are more
34 confident in their decisions. Modulations in willingness to wait are therefore understood to
35 reflect variations in decision confidence. However, optimal waiting behavior also depends
36 on the opportunity cost of opting to continue waiting for reward rather than starting a
37 new trial, which is set by the maximum achievable environmental reward rate^{5,12,13,14,15}.
38 Previous work on this task did not explicitly define the environmental reward rate and
39 therefore cannot specify how to choose the optimal average willingness to wait in a given
40 environment, which is the first-order statistic that needs to be optimized in order to
41 maximize reward rate. A complete account of optimal behavior in this task requires a
42 definition of the reward rate that specifies all relevant environmental statistics affecting
43 the environmental reward rate. Without such an account, it is not possible to determine
44 whether animal subjects performing this task are behaving optimally.

45 Here, we trained rats to perform an auditory evidence accumulation task¹⁶ requiring
46 binary decisions followed by a temporal wager⁵. We developed an expression for the
47 reward rate in the task, which allowed us to find the reward-rate-maximizing waiting
48 policy. Doing so made explicit a key environmental statistic: the travel time that is
49 incurred when moving on from one reward opportunity to the next. We found that rats
50 performing the task spent longer waiting for rewards than optimal agents who maximized
51 reward rate on matched datasets. This finding was consistent with the observation of
52 “overharvesting” in foraging studies^{17,14}. However, when we measured each individual
53 subject’s travel times and treated these as constraints on agents’ behavior, the rats’
54 overall average willingness to wait was not different from the optimal agents’, suggesting
55 that their waiting behavior was approximately optimal.

56 In addition to finding near-optimal overall average waiting behavior, the rats’ also
57 showed modulation of wait times by decision confidence, as has been seen previously^{5,7,8,10},
58 consistent with optimal agents. However, it is not clear how the near-optimal behavior
59 we observed can be executed algorithmically in the brain. Nor is it clear how that
60 decision might evolve in time and be influenced by sources of internal variability other
61 than confidence. To develop a candidate model of the waiting decision process, we
62 used the sequential probability ratio test (SPRT) to derive a decision variable that
63 could achieve optimal waiting via an accumulation to bound process, as is often used in
64 decision-making tasks^{18,19,20}. The decision variable was initialized at a point encoding
65 the decision confidence and evolved with a linear drift toward a single fixed bound that
66 encoded the estimate of the environmental reward rate. The drift in this model came
67 from the observation that as time elapses without reward after a decision, the odds that
68 the trial will be rewarded eventually decrease. Under the model, waiting continues until
69 the moment that the decision variable crosses the bound at which point the current trial
70 is abandoned.

71 The process model allowed us to compare various mechanistic sources of noise that
72 might effect the decision process¹⁶. We considered variability in the drift rate that
73 would produce the property of scale invariance often observed in the literature on
74 timing judgments^{21,22,23}, whereby the standard deviation of timing judgments grows
75 proportionally to the interval being timed. We also considered diffusion noise that would
76 corrupt the decision variable in each time step and cause the standard deviation of
77 timing judgments to scale with the square root of the interval being timed, as in the
78 drift diffusion model²⁴. Finally, we considered variability in the setting of the bound,

79 which would cause variability in timing judgments to be constant across all intervals
80 being timed. Surprisingly, we found that, whereas scale invariant timing noise had been
81 assumed to dominate during this task⁵, the data was most consistent with variability
82 in the bound. We speculate that the dominance of this source of variability may arise
83 from continual learning of the bound setting based on a recency weighted average of the
84 reward history, as has been observed in previous studies^{14,11} and is used in models of
85 foraging as evidence accumulation²⁵.

86 Our results lay out a more complete theory of optimal behavior in temporal post-
87 decision wagering tasks and present a process model for estimating the moment-to-moment
88 cognitive state of subjects during the waiting decision process. Taken together, these
89 contributions increase the interpretive value of post-decision temporal wagers for studies
90 of confidence and learning.

91 Results

92 Evidence accumulation task with confidence-guided waiting

93 We trained rats (n=16) to perform an auditory evidence accumulation task¹⁶ with
94 randomly delayed reward delivery (Fig. 1a), as in Lak et al.⁵. The task requires two
95 decisions of interest. First, the animal should decide which of two reward ports is more
96 likely to provide a water reward given the auditory stimulus. Then, the animal must
97 decide how long to wait for reward to be delivered before moving on to the next trial. In
98 principle, this decision could be made at the time of the port choice, but may also be
99 characterized as a series of decisions to wait at the chosen port for another timestep or
100 abandon the port and move on.

101 **Port choice** Rats performed the task in a chamber containing an array of three nose
102 ports. Rats initiated trials by poking their nose into a central nose port, which triggered
103 stimulus playback. The stimulus consisted of two trains of auditory clicks, generated
104 from two different Poisson rates, played from speakers on either side of the rat's head.
105 The rat's task was to listen to the click trains and then, after a "go" cue, report which
106 click train had the larger number of clicks by poking it's nose into the reward port on the
107 side associated with the higher click rate. Trials where the rat withdrew from the center
108 port before the "go" cue were labeled "center poke violations," invalidated, and the rat
109 was moved on to the next trial after a brief white noise stimulus.

110 There were two versions of this task, referred to as the location task and the frequency
111 task. Each rat was trained to perform one of the two tasks (location task, n=9; frequency
112 task, n=7). In the location task, one click train was played from a speaker to the left of
113 the center port and the other click train was played from a speaker to the right of the
114 center port. Rats were rewarded for choosing the port on the same side as the speaker
115 that emitted the greater number of clicks. In the frequency task, the clicks were played
116 in stereo, but the clicks in the two click trains were played at different frequencies. The
117 high frequency click train instructed rightward choices, whereas the low frequency click
118 train instructed leftward choices. The click trains depicted in Fig 1a correspond to the
119 location task. Trial difficulty was controlled by varying the rates of the two click trains.

120 **Wait time decision** After reporting a decision at the reward port, rats did not receive
121 immediate feedback. Instead, after correct choices, reward delivery was delayed until
122 an experimenter-determined reward time, t_r . On a subset of correct trials (mean: 7.1%

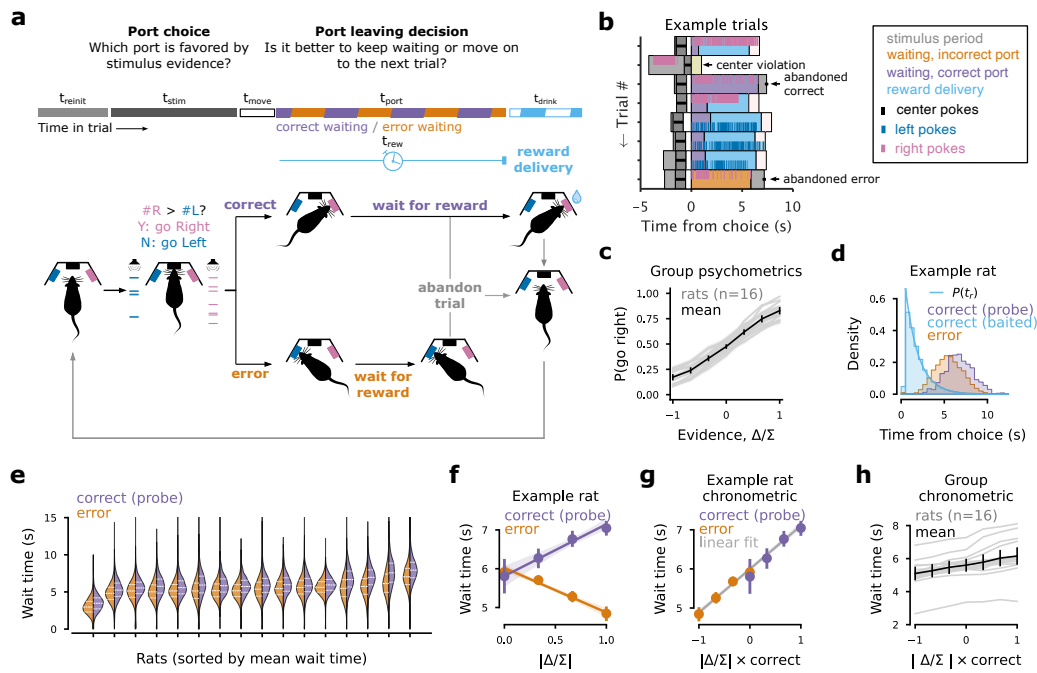


Figure 1: Evidence accumulation task with delayed rewards (A) Schematic of the task structure. Rats first make a port choice based on an auditory stimulus. On non-probe trials, if the port choice is correct, a reward is scheduled for delivery. The rat waits at the chosen port until either the reward is delivered or the rat decides to abandon the current trial and move on to the next. (B) Example trials from an example rat aligned to the time of the port choice. Correct port choices (purple bars) often lead to rewards (cyan bars), but sometimes the animal abandons the trial and center pokes to start a new trial before reward is delivered (an example is annotated). Error trials (orange bar) are never signalled and the rat eventually has to abandon the trial (an example is annotated). If the animal fails to hold it's nose in the center port during the stimulus period, the trial is considered a “center violation” (an example is annotated), increasing the time to the next possibly rewarded trial. (C) Probability of choosing the rightward reward port is plotted as a function of the evidence, operationalized as the click difference ($\#R - \#L$), normalized by the total number of clicks ($\#R + \#L$), denoted Δ/Σ . Each rat is shown as a gray trace and the average (with 95% confidence intervals) of the gray curves is shown in black. (D) Wait time distributions for an example rat conditioned on whether the decision was incorrect (shaded orange), correct with a baited reward (shaded cyan), or a correct probe trial (shaded purple). The reward delay distribution used to draw reward delivery time on non-probe trials is underlaid (cyan trace). (E) Violin plots showing each rat's wait times for error trials (orange) and correct probe trials (purple). Medians are plotted as dashed white lines with 25th and 75th as dotted white lines. (F) Mean wait time for an example rat in correct probe trial (purple) and error trials (orange) as a function of the absolute evidence strength, $|\Delta/\Sigma|$. Data are overlaid on linear fits to the correct and error trials separately. Errorbars are bootstrapped 95% confidence intervals. (G) Wait time chronometric curve showing the mean wait times from F plotted as a function of the strength of evidence favoring the chosen option, $|\Delta/\Sigma| \times \text{correct}$. Data are overlaid on a linear fit to all the data. (H) Wait time chronometric curve for each rat computed as in g (gray traces) along with the mean (with 95% confidence intervals) of all the gray traces (black trace).

123 $\pm 1.6\%$), rewards were omitted to provide an uncensored report of the rat's maximum
 124 willingness to wait on that trial. We refer to these trials as “probe” trials (but, note that

125 “catch” trials is another common terminology). On error trials, feedback was omitted
126 entirely. Rats were allowed to give up waiting for reward and move on to a new trial at
127 any time after making the port choice by withdrawing from the reward port and poking
128 into the center port. A series of example trials is shown in Figure 1b.

129 In these experiments, reward delays were drawn from an exponential distribution, so
130 that the density of reward delays was given by

$$P(t_r = t) = \frac{1}{\tau} e^{-(t_r - t_{r,\min})/\tau} \quad (1)$$

131 with time constant $\tau = 1.5$, and a minimum reward delay $t_{r,\min} \in (.05, .5)$. The
132 exponential distribution has a flat hazard rate on the interval $t \in (t_{r,\min}, \infty)$. This means
133 that, given that reward is set to be delivered on a trial, but hasn’t been delivered so far
134 by time t_w , the probability of receiving reward in the next time step is constant. We will
135 write the hazard rate of the reward distribution as

$$P(r_w | \neg R_w, R_\infty) = 1/\tau \quad (2)$$

136 where $r_w \equiv t_r \in (t_w, t_w + \delta t)$ is used to indicate the event that reward is delivered in
137 the infinitesimal timestep δt beginning with time t_w , the sum $R_w \equiv \sum_{i=0}^{w-1} r_i$ is used
138 to indicate whether reward is baited at some time before t_w , and the negation, $\neg R_w$,
139 indicates that no reward is baited before t_w . In this notation, R_∞ is used to indicate
140 whether reward is set to be delivered eventually in the trial. The resulting mean reward
141 delay is $\langle t_r \rangle \approx t_{r,\min} + \tau$ for trials where reward was baited.

142 **Rat behavior** All rats included in the study learned to perform the task with at least
143 60% accuracy (group mean: 74.4% correct trials). We computed each rat’s psychometric
144 function for the port choice (Figure 1c). The psychometric function was defined as the
145 probability of making a rightward choice given the stimulus evidence favoring rightward
146 choice. Stimulus evidence favoring rightward choice is defined as the click difference
147 normalized by the total number of clicks, $\Delta/\Sigma \equiv \frac{\#R - \#L}{\#R + \#L}$, where $\#R$ represents the
148 number of clicks favoring rightward choice and $\#L$ represents the number of clicks favoring
149 a leftward choice on a given trial.

150 We measured time spent waiting for reward at the side port in three trial types of
151 interest: error trials, correct probe trials (where no reward is baited), and correct trials
152 where reward is baited. We excluded trials where the animal took more than 2 seconds
153 to initiate a new trial by center poking after leaving the chosen port. This is a standard
154 criterion used to focus analysis on trials where the animal is engaged in the task⁵. The
155 example rat was willing to wait long enough to receive reward on most correct trials,
156 so the distribution of waiting times on correct trials where reward was baited closely
157 resembles the reward delay distribution (Figure 1d).

158 On trials where the rat waited long enough to receive reward, the full duration that
159 the rat would have been willing to wait is unknown, because reward delivery censors our
160 observation of the full willingness to wait. We used the probe trials to measure how long
161 rats were willing to wait on correct trials. On both error trials and correct probe trials,
162 rats were willing to wait much longer than the typical reward delays (Fig 1d,e). This
163 held across all rats who learned the task (Fig 1e). Additionally, all rats waited longer at
164 the choice port after correct choices on probe trials than after errors (Fig. 1d,e; $p < .01$,
165 rank-sum test, 16/16 rats), indicating that waiting was guided by an internal estimate of
166 decision accuracy.

167 To measure the modulation of wait time by the stimulus evidence, we plotted the
168 example rat's average wait time as a function of the absolute stimulus strength, $|\Delta/\Sigma|$,
169 separately for correct probe trials and error trials (Fig 1f). We expected correct trials with
170 strong signal to be the trials in which the animal has the highest confidence on average.
171 Indeed, wait times were longest for correct trials where the evidence most strongly favored
172 the choice made by the animal, as has been seen previously^{5,7,10}. Correspondingly, wait
173 times were shortest in the trials where the evidence most strongly favored the alternative
174 not chosen by the animal. We expect these to be the trials where the animal has the
175 lowest confidence on average.

176 To create an axis along which both confidence and wait time should increase mono-
177 tonically, we used the strength of the evidence supporting the option chosen by the rat,
178 $|\Delta/\Sigma| \times \text{correct}$. This quantity takes positive values when the animal makes a correct
179 choice and negative values when the choice is incorrect. When we plot the example
180 rat's wait times against this axis, we see a graded increase in wait time as a function of
181 evidence supporting the choice (Fig 1g). The data is overlaid on a linear fit to the data,
182 which has a significantly positive slope (Pearson's $r = .22$, $p < .01$). We computed wait
183 time as a function of evidence supporting choice for all rats (Fig 1h) and computed linear
184 fits to each rat. All rats had a significant, positive relationship between waiting time and
185 the strength of evidence for the chosen option ($p < .01$ for 16/16 rats). This indicates
186 that all of our rats modulated waiting times by their decision confidence.

187 Overall reward rate maximization depends on travel time

188 Previous work⁵ has shown that in order to maximize reward rate in decision tasks with
189 delayed reward, subjects should be willing to wait longer when they are more confident in
190 their decisions. However, the trial-by-trial modulation of waiting time by confidence alone
191 is not enough to maximize the long term average reward rate. To maximize the long term
192 average reward rate, subjects must also find an appropriate overall average willingness to
193 wait. This value depends on a variety of other environmental statistics that influence
194 the environmental reward rate. However, previous work has not developed an explicit
195 expression for the environmental reward rate in the task, so it has not been possible
196 to test whether rats learn this first-order optimization of overall wait time. In another
197 study of rats performing confidence-guided waiting for delayed rewards, Stolyarova et al.⁸
198 noted that their rats' overall wait times were long relative to the average time of reward
199 delivery, as is true in our rats. The authors interpreted this observation as likely being a
200 deviation from optimality in rat behavior. This would be consistent with previous studies
201 in human¹⁴ and animal subjects^{17,26} which report a bias, referred to as "overharvesting,"
202 toward spending more time than would be optimal on a given reward opportunity before
203 moving on to the next. Here, we develop a definition of the reward rate that makes all
204 relevant environmental statistics explicit. We can then determine the optimal average
205 willingness to wait for a given environment, making it possible to test whether subjects
206 achieve optimal behavior in the task.

207 To develop an expression for reward rate in our task, we make use of optimal foraging
208 theory^{12,13}, which describes the optimal time an agent should spend in a series of "patches"
209 containing depleting, continuous rewards before traveling to the next patch. In each trial,
210 we think of the rat's nose poke into the chosen reward port as an entry into a "patch."
211 We refer to the time spent at the port as t_{port} . We refer to the elapsed time between
212 leaving the reward port on a given trial and entering a reward port on the next trial
213 as "travel time" and note it in equations as t_0 (Fig. 2a). In this task, the travel time

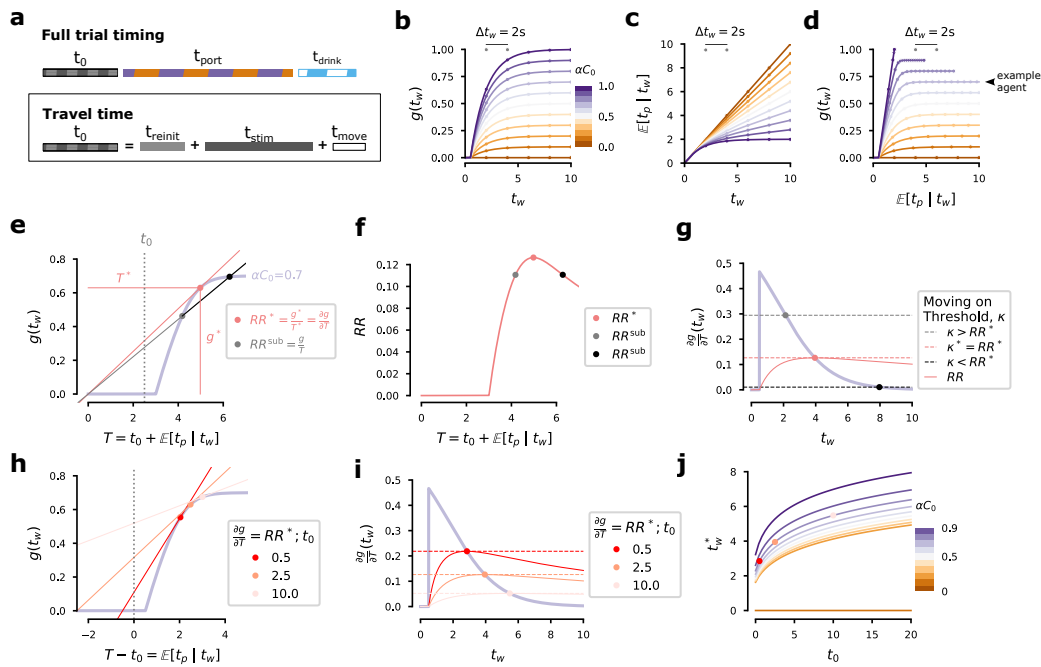


Figure 2: Across trial reward maximization depends on travel time (A) Task timing can be broken into travel time, t_0 , time at the port, t_{port} , and time spent drinking reward, t_{drink} (top). Travel time includes all periods between the end of the waiting/drinking period and the start of the next period of waiting, including trial reinitiation time, t_{reinit} , stimulus playback time, t_{stim} , and movement time, t_{move} (bottom). (B) Expected reward per trial, $g(t_w)$, if willing to wait for time t_w . Colormap shows probability of eventual reward, αC_0 . Points indicate 2s increments of t_w . (C) Expected time spent at the port per trial, $\mathbb{E}[t_{\text{port}} | t_w]$, plotted as in B. (D) Expected reward in a trial, $g(t_w)$, as a function of expected time at the port, $\mathbb{E}[t_{\text{port}} | t_w]$, plotted as in B and C. The effect of an additional 2s increment in t_w now depends on t_w and αC_0 . (E-J) Consider reward maximization for an example agent with $\alpha C_0 = .67$ on every trial. (E) Expected reward per trial, $g(t_w)$, plotted as a function of trial time, $T(t_w) = t_0 + \mathbb{E}[t_{\text{port}} | t_w]$ (solid purple trace). The reward rate is $RR = \frac{g(t_w)}{T(t_w)}$ and is maximized when $RR^* = \frac{\partial g}{\partial T}(t_w^*)$ (red point) at $T(t_w^*)$. The red trace from the origin through this point has the highest slope of any line from the origin through the purple trace. All other values of $T(t_w)$ are suboptimal (e.g., gray and black points achieve the reward rate RR^{sub} , which is the slope of the gray and black traces). Here, t_0 is set to 2.5s. (F) The reward rate for the example agent is plotted as a function of $T(t_w)$ (dashed red trace) with the maximum reward rate marked (red point) along with the reward rate achieved if not willing to wait long enough (gray point) or willing to wait too long (black point). (G) The instantaneous reward expectation within a trial after time t_w passes without receiving reward, $\frac{\partial g}{\partial T} = P(r_w | \neg R_w)$ is plotted as a function of t_w (solid purple trace). The session reward rate from F is shown for comparison (solid red trace). Reward rate is maximized if the agent sets a moving on threshold, $\kappa = RR^*$ (dashed red trace). Suboptimal reward is achieved when κ is not set to RR^* (e.g., gray and black traces). (H) Probability of reward, $g(t_w) = P(R_w)$, plotted as a function of expected port time, $T(t_w) - t_0 = \mathbb{E}[t_{\text{port}} | t_w]$ (rather than as a function of T as in E). Reward maximizing solutions are marked for three values of t_0 , including the same from E, one (darkest red trace) smaller, and one larger (lightest red trace). (I) Instantaneous reward expectation plotted as in G with the reward rates and optimal settings of κ for the three levels of t_0 used in H. (J) Optimal wait time, t_w^* , as a function of travel time, t_0 , for all levels of αC_0 used in panels b-d, except $\alpha C_0 = 1$, which corresponds to $t_w^* = \infty$. Example levels of t_0 are marked for the example level of αC_0 (red points).

214 includes the period of trial reinitiation from leaving the reward port on the previous trial
215 to entering the center port on the current trial (labeled t_{reinit} in Fig 2a), plus the stimulus
216 period during which the rat hears the stimulus to inform the port choice (labeled t_{stim}
217 in Fig 2a), and the time it takes to move from the center port to the chosen reward
218 port (labeled t_{move} in Fig 2a). Unlike in the classical foraging theory where reward is
219 continuous and the agent has perfect information about the patch identity, rewards in
220 our task are stochastic, limited to at most one per trial, and the subject has only partial
221 information about whether it is in a rewarded or unrewarded patch. The optimal strategy
222 for such a task has been described for environments where rewarded and unrewarded
223 patches occur with equal probability¹⁵. We generalize that theory to arbitrary initial
224 probability of being in the rewarded patch.

225 The partial information about patch type comes from the agent’s decision confidence,
226 an estimate of the probability that the port choice was correct at the time of the decision
227 given the available perceptual evidence, which we write as

$$C_0 \equiv P(\text{correct} \mid \text{percept}). \quad (3)$$

228 If the agent believes that all correct choices are rewarded, then their initial estimate of
229 the probability that the choice will eventually be rewarded is C_0 . If the agent believes
230 that correct choices are only rewarded in some fraction, α , of non-probe trials, then the
231 probability that the choice will be rewarded eventually is

$$P(R_\infty) = P(R_\infty \mid \text{correct})P(\text{correct} \mid \text{percept}) = \alpha C_0. \quad (4)$$

232 Later, we will see that the agent’s posterior belief about whether the port choice will be
233 rewarded (rewarded or unrewarded) falls over time.

234 We will simplify the expression for reward rate as a function of the subject’s willingness
235 to wait for reward across trials by beginning with the case of an agent who has no trial
236 variation in decision confidence (i.e., $P(\text{correct} \mid \text{percept}) = P(\text{correct})$). This agent
237 should therefore be willing to wait the same amount of time, t_w , on every trial. We can
238 write the expected overall reward rate for an agent willing to wait until time t_w as the
239 ratio of expected reward per trial, $g(t_w)$, and expected time per trial, $T_{\text{total}}(t_w)$:

$$RR_{\text{total}} = \frac{g(t_w)}{T_{\text{total}}(t_w)} = \frac{g(t_w)}{t_0 + \mathbb{E}[t_{\text{port}} \mid t_w] + \mathbb{E}[t_{\text{drink}} \mid t_w]} \quad (5)$$

240 where t_0 is the “travel” time between leaving the chosen side port on one trial and nose
241 poking at a side port on the next trial, $\mathbb{E}[t_{\text{port}} \mid t_w]$ is the expected time spent at the side
242 port, and $\mathbb{E}[t_{\text{drink}} \mid t_w]$ is the expected time spent consuming reward. While $\mathbb{E}[t_{\text{drink}} \mid t_w]$
243 affects the overall reward rate, it can be ignored for the reward maximization process
244 (see Supplemental Information for derivation).

245 Both expected reward on each trial, $g(t_w)$, and expected time at the port, $\mathbb{E}[t_{\text{port}} \mid t_w]$,
246 depend implicitly on the probability that reward will be delivered eventually if the agent
247 waits long enough ($P(R_\infty)$; equation 4). Expected reward per trial rises exponentially as a
248 function of willingness to wait toward an asymptote at αC_0 (Figure 2b; see Supplemental
249 Information for mathematical details). Expected time at the port per trial increases as
250 a function of willingness to wait, but does not asymptote except in the case that all
251 trials are eventually rewarded ($\alpha C_0 = 1$; Figure 2c; see Supplemental Information for
252 mathematical details). Otherwise, greater willingness to wait increases expected time
253 spent at the port on each trial (in the extreme case where no trials are rewarded, the
254 expected time at the trial is equal to willingness to wait). Figure 2d shows how expected

255 reward per trial increases as the expected time at the port increases when t_w is varied,
 256 combining the information in Figures 2b and c.

257 We can maximize the total reward rate in equation 5 by computing its derivative and
 258 setting it to zero, which yields

$$\frac{\partial g}{\partial T_{\text{total}}}(t_w^*) = RR_{\text{total}}^* \quad (6)$$

259 where t_w^* and RR^* are the reward-rate-maximizing willingness to wait and the corre-
 260 sponding reward rate at that optimal t_w^* (see supplement for derivation). This is a
 261 generalization of the marginal value theorem¹² to the case of stochastic rewards¹⁵. That
 262 is, equation 6 states that the prescribed rule for maximizing reward rate is to be willing
 263 to wait for reward until the time, t_w^* , when the derivative of expected reward rate within
 264 the trial, $\frac{\partial g}{\partial T_{\text{total}}}(t_w^*)$, falls to the level of the maximum achievable reward rate across trials,
 265 RR_{total}^* . The latter quantity is the opportunity cost of continuing to wait for reward
 266 rather than beginning a new trial.

267 As noted above, we can simplify the computation by ignoring the reward consumption
 268 time in the denominator of equation 5. Instead, we will maximize the expected reward
 269 per time spent pursuing (not consuming) reward, $T(t_w)$:

$$RR = \frac{g(t_w)}{T(t_w)} = \frac{g(t_w)}{t_0 + \mathbb{E}[t_{\text{port}} | t_w]}, \quad (7)$$

270 which is maximized when

$$\frac{\partial g}{\partial T}(t_w^*) = RR^*. \quad (8)$$

271 (Note that previous work⁵ assumed the optimality condition $\frac{\partial g}{\partial T}(t_w^*) = RR_{\text{total}}^*$, which
 272 produces suboptimal behavior when $t_{\text{drink}} \neq 0$, as is the case in our data.)

273 This reward maximization rule can be understood graphically by plotting expected
 274 reward in a trial, $g(t_w)$, as a function of expected time pursuing reward in the trial, $T(t_w)$,
 275 for an example agent (Fig 2d,e). The reward rate for any choice of t_w will be equal to
 276 the slope of a line that passes from the origin through the point $(T(t_w), g(t_w))$. The
 277 maximum possible slope (i.e., maximum possible reward rate) is achieved when this line
 278 is tangent to the reward rate curve satisfying equation 8 (Fig 2e,f). In standard optimal
 279 foraging theory, the forager receives continuous reward and gives up and moves on at
 280 a time under its full control, $\mathbb{E}[t_{\text{port}} | t_w] = t_w$, which would mean that $\frac{\partial T}{\partial t_w} = 1$ and
 281 equation 8 reduces to the marginal value theorem. However, because our task provides
 282 at most one reward per trial, the agent must estimate the expected rate of reward in
 283 each trial through experience and then set an upper bound, t_w , on the time it will spend
 284 waiting for reward delivery.

285 **Optimizing wait time** Now that we have found the condition under which reward
 286 rate is maximized (equation 8), we are able to find t_w^* for a given set of environmental
 287 statistics. To do so, we first compute the derivative in left hand side of equation 8, which
 288 we write as

$$\frac{\partial g}{\partial T}(t_w) = \lim_{\delta t \rightarrow 0} \frac{g(t_w + \delta t) - g(t_w)}{T(t_w + \delta t) - T(t_w)}. \quad (9)$$

289 For an agent that has already waited for time t_w , this quantity takes one of two values.
 290 If reward has already been delivered, the derivative is zero ($g(t_w) = g(t_w + \delta t) = 1$), and

291 the agent should move on to the next trial as soon as it finishes consuming reward. In
 292 the second case, the agent has not yet received the reward ($g(t_w) = 0$). In this case, the
 293 expected reward after waiting for an additional time step is the probability of receiving
 294 reward in the next time step given that it hasn't been delivered so far, $P(r_w | \neg R_w)$.
 295 This quantity is the hazard rate of the distribution of reward delays including the trials
 296 in which no reward is baited. We refer to this quantity as the instantaneous reward
 297 expectation after waiting for time t_w without reward. We can write it as a product of
 298 the reward hazard rate for trials where reward is baited (equation 2) and the posterior
 299 probability that reward will be delivered in a trial given that it has not been delivered so
 300 far:

$$P(r_w | \neg R_w) = P(r_w | \neg R_w, R_\infty)P(R_\infty | \neg R_w). \quad (10)$$

301 We refer to the second term as the agent's same posterior belief that reward will be delivered
 302 on a given trial after waiting for time t_w without receiving reward. We write this quantity
 303 using Bayes' rule

$$P(R_\infty | \neg R_w) = \frac{P(\neg R_w | R_\infty)P(R_\infty)}{P(\neg R_w)} \quad (11)$$

304 and evaluate it for the distribution used in our experiment

$$P(R_\infty | \neg R_w) = \frac{\alpha C_0 e^{-(t_w - t_{r,\min})/\tau}}{1 - \alpha C_0 + \alpha C_0 e^{-(t_w - t_{r,\min})/\tau}} \quad (12)$$

305 (see Supplemental Information materials for detailed derivation). This quantity has the
 306 value αC_0 at the time of choice ($t_w = 0$) and falls to 0 as time passes. Note that when
 307 the agent is unaware of the probe trials (i.e., the agent estimates $\alpha = 1$), equation 12
 308 is equal to the posterior belief that the port choice was correct, the posterior decision
 309 confidence, after waiting for time t_w without reward.

310 Substituting equations 2 and 12 into equation 10, we get the instantaneous reward
 311 expectation after waiting for time t_w without receiving reward

$$P(r_w | \neg R_w) = \frac{1}{\tau} \cdot \frac{\alpha C_0 e^{-(t_w - t_{r,\min})/\tau}}{1 - \alpha C_0 + \alpha C_0 e^{-(t_w - t_{r,\min})/\tau}} \quad (13)$$

312 for $t_w \geq t_{r,\min}$ (instantaneous reward expectation is 0 for $t_w < t_{r,\min}$). Note that
 313 equation 13 is equivalent to equation 5 in Lak et al.⁵ if we substitute $C = \alpha C_0$ and
 314 $t = t_w - t_{r,\min}$. However, our derivation clarifies that even though the reward hazard
 315 rate is fixed in the task, there is a decrease in instantaneous reward expectation while
 316 waiting for reward that can be attributed to a decrease in the posterior belief that the
 317 trial will be rewarded as time passes without reward delivery. Later, we will make use of
 318 this observation to develop a model for describing the port-leaving decision as a process
 319 that unfolds in time.

320 Instantaneous reward expectation (equation 13) for the example agent is plotted as
 321 a function of elapsed time without reward in Figure 2g. To find t_w^* , the agent needs to
 322 estimate the instantaneous reward expectation as a function of t_w and choose a moving
 323 on threshold, κ , whose optimal value is RR^* (Fig 2g). When $\kappa < RR^*$, the agent is
 324 impatient and receives a below average reward rate, and when $\kappa > RR^*$, the agent wastes
 325 time at the reward port that would be better spent starting a new trial. Choosing the
 326 appropriate threshold will lead to optimal waiting with

$$t_w^* = t_{r,\min} + \tau \left(\log \frac{\alpha C_0}{1 - \alpha C_0} - \log \frac{RR^* \tau}{1 - RR^* \tau} \right) \quad (14)$$

327 (see Supplementary Information for detailed derivation). This is equivalent to equation 6
328 in Lak et al.⁵ if we again substitute $C = \alpha C_0$, set $t_{r,\min} = 0$, and substitute $\kappa = RR^*$.
329 But note that in Lak et al.⁵, κ is defined in words as the “environmental reward rate”
330 whereas we have developed an explicit expression for RR^* , which clarifies that the correct
331 RR^* in this expression is not the total environmental reward rate (equation 5), but rather
332 the reward per time spent pursuing reward (equation 7). Moreover, because we have
333 provided an expression for RR that makes all relevant environmental statistics explicit,
334 we can now compute t_w^* for a given experiment, which was not possible previously.

335 The optimal strategy for maximizing reward is influenced by the confidence on a given
336 trial and all of the factors that influence the maximum achievable reward rate, including
337 the probe trial fraction, the reward delivery time constant τ , and the travel time, t_0 . As
338 t_0 increases, the maximum possible reward rate decreases and the value of t_w^* increases.
339 The reward optimization procedure is depicted for three example levels of travel time in
340 Figure 2h,i.

341 We found the optimal willingness to wait, t_w^* , as a function of travel time, t_0 , for all
342 the levels of αC_0 by using root finding to solve equation 8 (Figure 2j; see Methods for
343 details). The amount of time that a reward rate maximizing agent is willing to wait in
344 this task increases monotonically as travel time increases for all levels of αC_0 (except 0
345 and 1, where the agent should either not be willing to wait at all, or should always be
346 willing to wait until the reward is delivered, respectively).

347 Rats maximize reward rate after accounting for travel time

348 Now that we can compute t_w^* for an agent with fixed confidence across trials, we can
349 test whether our rats achieved the maximum possible reward rate across trials. To do
350 this, we compared rats’ willingness to wait, averaged across trials, to t_w^* , the willingness
351 to wait that would maximize the reward rate for an agent with fixed confidence across
352 trials. To estimate the rats’ average willingness to wait across trials, we computed the
353 average wait time for correct probe trials and for a subset of error trials, subsampled so
354 that the proportion of error trials used in this analysis matched the proportion in the full
355 dataset (which also includes non-probe trials; Fig 3c,d). To compute t_w^* for each subject’s
356 dataset, we estimated the necessary terms from that subject’s data: α was the fraction
357 of non-probe trials in the rat’s dataset, C_0 was the fraction of correct trials, and t_0 was
358 estimated from the mean travel time for the rat (after excluding the longest 1% of travel
359 times, because the rats occasionally fully disengaged from the task for long periods of
360 time; Fig 3a, b). Using these terms allowed us to evaluate both sides of equation 8 and
361 compute t_w^* by root finding (see Methods for details).

362 A fully optimal agent performing this task should spend as little time traveling from
363 one reward opportunity to the next. However, our rats spent more time than necessary
364 traveling between reward opportunities. This was due to the self-paced nature of the
365 task and exacerbated by center poke violations, which caused trials to be invalidated,
366 further delaying time to the next reward opportunity. We reasoned that it may be very
367 difficult for the rats to further minimize travel time. Among other things, decreasing
368 travel times would require reducing the center port violation rate, which the animal is
369 presumably already incentivized to do as much as possible. Therefore, we treated travel
370 time as a constraint experienced by the animal and used the agents with matched travel
371 times to ask whether the rats maximized reward rate given this constraint. To compare
372 each rat to an agent who had also minimized travel time, we also computed t_w^* for an
373 agent whose average t_0 was set to the value of the shortest travel time achieved by the

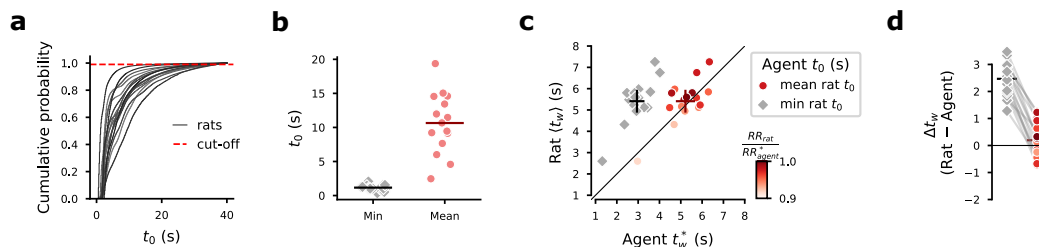


Figure 3: Rats maximize reward rate after accounting for travel times (A) Cumulative distribution of travel times, t_0 , for each rat (black traces). We excluded the longest 1% of travel times for each rat to compute the mean travel time (dashed red line). (B) Minimum t_0 for each rat (gray diamonds) and mean included t_0 for each rat (red points). Means of the presented values are shown as horizontal lines for the minimum and mean t_0 values. (C) Each rat’s willingness to wait is plotted as a function of the optimal willingness to wait for an agent with the same reward delivery statistics, trial accuracy, and the same mean travel time as the rat, but without trial-by-trial variations in confidence (red points). Each rat’s willingness to wait is also plotted as a function of the optimal willingness to wait for an agent as described, but with travel time equal to the minimum travel time achieved by the rat (gray diamonds). Rat willingness to wait is estimated by computing the mean wait time in correct probe trials and in a subsample of error trials (subsampling so that the proportion of error trials in this analysis is equal to the proportion of error trials for the full dataset when all correct trials are included). Optimal willingness to wait was determined by root finding using equation 8. The mean and 95% confidence intervals are shown as crosses for each group. The shade of the red points indicates the fraction of the reward maximizing agent’s reward rate achieved by the rat. For the comparison between the rat and the agent with matched mean t_0 , the shade of red of the points indicates the fraction of the maximized agent reward rate achieved by the rat. (D) Difference between the rat data and the agent data for the rats’ travel times (red points) and for the shorter travel times (gray diamonds). Colormap is the same as in C. The mean difference for each group is marked with a horizontal line.

374 rat across sessions (Fig 3b).

375 We found that the rats’ average willingness to wait was not different from the reward
 376 maximizing agents’ with matched travel times ($p = .25$, paired t-test; Fig 3c,d). However,
 377 rats’ wait times were much longer than those of the agents optimized with short travel
 378 times ($p = 6.15 \times 10^{-11}$, paired t-test; Fig 3c,d). Subjects “overharvested” relative to fully
 379 optimal agents, as has been seen in previous studies of foraging behaviors^{17,14}. However,
 380 when travel time is treated as constrained, and behavior is optimized over t_w alone,
 381 subjects’ overall willingness to wait was near-optimal, approximately maximizing their
 382 overall reward rate.

383 Process model for optimal confidence-modulated waiting

384 To understand how the port-leaving decision might be implemented in the brain, we
 385 developed a process model that described the decision to stop waiting as an accumulation
 386 to bound process. This model provides us with a cognitively tractable algorithm that
 387 can achieve optimal waiting and model the cognitive state of the animal during this task,
 388 which may be useful for studies of neural recordings in the task. It also allows us to
 389 capture variability in waiting that may be explained by sources of internal variability

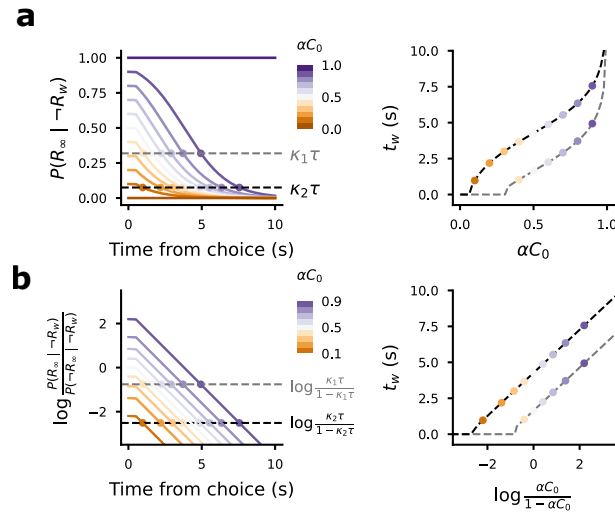


Figure 4: **Optimal waiting can be implemented with linear drift from a confidence-dependent initial point to a fixed bound** (A) Optimal wait time procedure as presented in fig 2g,i, but scaled by τ , so that we track the evolution of $P(R_\infty | \neg R_w)$ to a bound $\kappa\tau$. Two settings of κ are shown (gray and black dashed traces) along with the bound hitting times for different levels of αC_0 (left) and the resulting wait times are plotted against αC_0 (right). (B) Equivalent model with linear drift from an initial point $x_0 = \log \frac{\alpha C_0}{1 - \alpha C_0}$, which terminates at a bound $Z = \log \frac{\kappa\tau}{1 - \kappa\tau}$. This process leads to the same waiting times as in A. Colormaps are the same as in A, but note that $C_0 = 0$ and $C_0 = 1$ do not appear in panel B, because they start at $-\infty$ and $+\infty$, respectively, and never hit the bound.

390 other than variations in decision confidence.

391 To develop such a model, we used the sequential probability ratio test¹⁸ to derive a
 392 tractable update rule, a linear drift with time, for a decision variable that can be used
 393 to produce optimal wait times. From equations 8 and 10, we know that optimal policy
 394 is to stop waiting when the instantaneous reward expectation falls to the level of the
 395 maximum reward rate in the environment, or equivalently, when the posterior belief that
 396 reward will be delivered eventually, $P(R_\infty | \neg R_w)$ falls to the maximum reward rate in
 397 the environment scaled by τ

$$P(R_\infty | \neg R_w) = RR^* \tau. \quad (15)$$

398 We can describe an agent who uses this strategy, but does not necessarily choose the
 399 optimal bound by replacing RR^* with a parameter κ whose optimal value is $\kappa^* = RR^*$.
 400 This process is shown in Figure 4a.

401 To produce a decision variable that is tractable to update, we define x_w as the log
 402 odds of eventual reward delivery, given that reward has not been delivered by time t_w :

$$x_w = \log \frac{P(R_\infty | \neg R_w)}{P(\neg R_\infty | \neg R_w)}. \quad (16)$$

403 To find an update rule that integrates the information from the passage of time without
 404 reward into x_w , we decompose x_w into two terms representing the previous value x_{w-1}

405 and an update Δx when timestep w elapses without reward (note that if reward is
406 delivered in timestep w , the process ends):

$$\begin{aligned} x_w &= \log \left(\frac{P(R_\infty | \neg R_{w-1}) P(\neg r_w | R_\infty, \neg R_{w-1})}{P(\neg R_\infty | \neg R_{w-1}) P(\neg r_w | \neg R_\infty, \neg R_{w-1})} \right) \\ &= x_{w-1} + \log (P(\neg r_w | R_\infty, \neg R_{w-1}) \Delta t) \\ \Delta x &= \log (P(\neg r_w | R_\infty, \neg R_{w-1}) \Delta t) \end{aligned} \quad (17)$$

407 In the time before the earliest possible reward delivery ($t_w < t_{r,\min}$), the update term is
408 0 and x_w is constant, afterward x drifts at the hazard rate of the reward distribution

$$\Delta x = \log \left(1 - \frac{1}{\tau} \Delta t \right) = -\frac{1}{\tau} \Delta t$$

409 where we have used $\log(1 - n) \approx -n$ for $|n| \ll 1$. Taking the timestep to zero, we get
410 the linear drift dynamics

$$dx = -\frac{1}{\tau} dt, \quad (18)$$

411 which we can combine with equation 4 to write x_w as a function of it's initial value x_0 :

$$x_w = x_0 - (t_w - t_{r,\min})/\tau \quad x_0 = \log \frac{\alpha C_0}{1 - \alpha C_0}. \quad (19)$$

412 By equations 15 and 16, we know that the agent should stop waiting when x_w hits a
413 bound specified by

$$Z^* = \log \frac{RR^* \tau}{1 - RR^* \tau}. \quad (20)$$

414 If the waiting process terminates when x hits the bound Z^* , we achieve the reward
415 maximizing wait times equivalent to equation 14:

$$t_w^* = t_{r,\min} + \frac{Z^* - x_0}{A^*} \quad (21)$$

416 where A is a drift rate whose optimal value is $A^* = -1/\tau$. The evolution of x_w and
417 equivalence of this waiting process with that of Figure 4a is shown in Figure 4b. This
418 expression for optimal willingness to wait is equal to equation 6 in Lak et al.⁵ after
419 setting $t_{r,\min} = 0$ and $C = \alpha C_0$. But, now also has an algorithmic interpretation that
420 may be possible to implement in the brain. In words, optimal wait time decisions can be
421 made by initializing a decision variable at a value that is set by the decision confidence
422 and evolving it toward a fixed bound that is set by the overall reward rate and reward
423 delivery timing in the task. The drift rate toward that bound is set by the reward hazard
424 rate.

425 Contributions of different sources of timing noise to waiting process

426 In studies of timing judgments, subjects often exhibit the phenomenon of scale invariance
427 in which the standard deviation of timing estimates increases linearly in proportion to
428 the interval to be timed^{21,22}. A previous model⁵ of confidence-guided waiting behavior
429 assumed that scale invariant timing noise was the dominant source of noise affecting wait
430 times. However, this has not been directly tested and it is not clear that timing in this

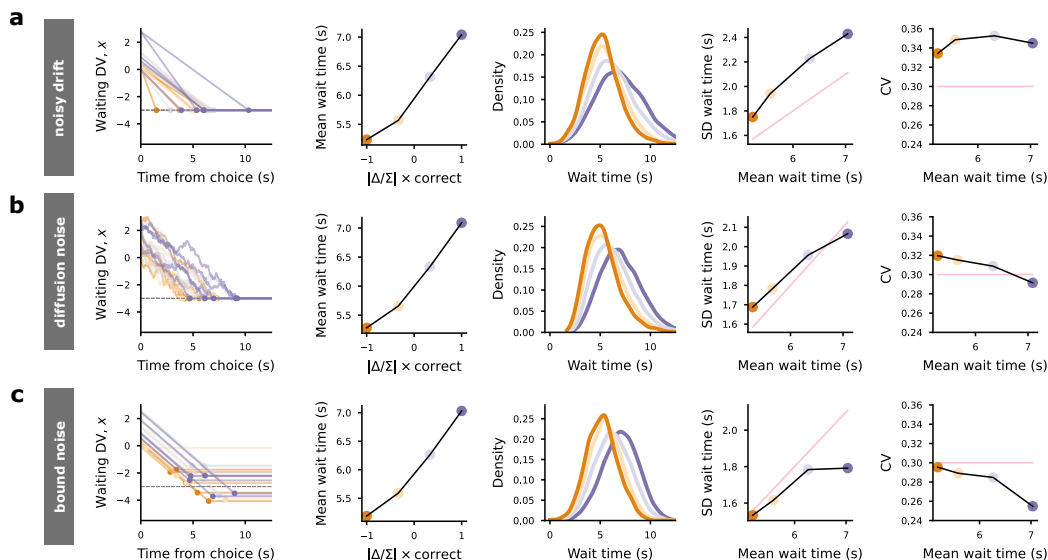


Figure 5: **Candidate timing noise models (A-C)** We consider three candidate models for adding noise to the wait time decision. In all three models, the agent makes left/right port choice based on a stimulus, which is corrupted by perceptual noise and creates an accompanying decision confidence, C_0 . The waiting decision variable (DV), x is created by setting $x_0 = \log \frac{\alpha C_0}{1 - \alpha C_0}$ where α is the fraction of non-probe trials. And x drifts with rate A toward a boundary Z . The agent is willing to wait for reward until x hits the bound Z , but gives up and moves immediately when the bound is hit. (A) Left panel shows example particles from a model in which the noise comes from variability in the drift rate. This noisy drift model produces the scale invariant property in which the ratio of the standard deviation of hitting times to the mean of hitting times is constant across all x_0 values. Traces are colored according to the binned normalized evidence favoring the choice, $|\frac{\Delta}{Z}| \times \text{correct}$. Second panel from left shows kernel density estimates of bound hitting times for each of the bins of normalized evidence favoring the choice with the same colormap. The third panel from the left shows the mean wait time as a function of the normalized evidence favoring the choice (black trace with points colored by bin). The second panel from the right shows the standard deviation of the bound hitting times as a function of the mean in each bin (black trace with points colored by bin). The generative relationship between standard deviation and mean is underlaid (pink trace). Any deviation reflects noise added by the psychometric decision process. The rightmost panel shows the coefficient of variation (CV) as a function of the mean wait time in each bin (black trace with points colored by bin). Again, the generative relationship is underlaid (pink trace). (B) Plots are as in A, but for a model with a diffusion noise process in which noise is added in every time step. The pink traces from A are maintained for comparison. (C) Plots as in A and B, but for a model in which noise comes from variability in the bound. The pink traces from A and B are maintained for comparison.

431 task is dominated by the same sources of variability as in interval timing tasks where
 432 the goal is to learn to respond when reward is most likely, rather than to persist until
 433 the moment that reward is sufficiently unlikely that it is worth giving up and moving
 434 on. It is also not trivial to separate noise in the timing decision from variability in the

435 confidence level that might result from any given percept.

436 We used the process model defined in the previous section to examine the patterns
437 of timing variability expected when different aspects of the process were corrupted by
438 noise. We considered three possible ways of adding timing noise to our process model.
439 In the first, we add noise to the drift rate, A (Fig 5a). This produces scale invariant
440 variability. For a given initial point, x_0 , the standard deviation in bound hitting times
441 that is proportional to the mean. In the second model, we added diffusion noise to the
442 position of x at each time step, which adds a random walk to the deterministic drift
443 (Fig 5b). In this model, the standard deviation of wait times is proportional to the square
444 root of the mean hitting times for a given x_0 , meaning that the standard deviation will
445 grow slower than for the scale invariant model. Finally, we considered a model with
446 a noisy bound (Fig 5c). In this model, the standard deviation is constant regardless
447 of the initial point x_0 . The noise parameters were chosen for each model to produce a
448 coefficient of variation (CV; ratio of standard deviation to the mean) of 0.3 when $x_0 = 0$.
449 For the scale invariant model, the CV is 0.3 for all values of x_0 , which is the level of noise
450 assumed in Lak et al.⁵.

451 While the patterns of timing variability produced by each of these models are simple
452 when the initial wait time decision variable, x_0 , is known, we don't have access to x_0 for
453 our rats. To understand the pattern of variability expected under each model when x_0 is
454 unknown, we generated simulated x_0 values for 50,000 trials. To do this, we supposed a
455 signal detection theory model of the decision process in which the stimulus is characterized
456 by the ratio between the click difference and the total number of clicks on each trial,
457 $s = \frac{\Delta}{\Sigma}$. For each trial, we generated a percept by adding noise with standard deviation
458 σ_s to the stimulus, $p = s + \xi$ where $\xi \sim \mathcal{N}(0, \sigma_s^2)$. Decisions we made by comparing the
459 stimulus to a decision boundary, $b = 0$. Confidence was then defined (beginning with
460 equation 3) as

$$\begin{aligned} C_0 &\equiv P(\text{correct} \mid p) \\ &= \int_s P(\text{correct} \mid p, s) P(p, s) ds \\ &= \int_s \mathbb{1}^{\text{sign}(p)=\text{sign}(s)} P(p \mid s) P(s) ds \end{aligned} \quad (22)$$

461 where we are integrating the probability of experiencing the percept p given the stimulus
462 s over all levels of s that would produce a correct choice (see Supplemental Information
463 for full equations). For the simulations, we assumed a uniform prior, $P(s)$, and the
464 probability of a percept given a stimulus is the Gaussian $P(p \mid s) = \mathcal{N}(s, \sigma_s^2)$. We used
465 a value of σ_s^2 that best fit an example rat (see Methods for details). We assumed an
466 accurate estimate of the non-probe trial frequency, α . Combining this with confidence,
467 we produced a sample x_0 for each trial using equation 19. We then generated a sample
468 willingness to wait on each trial by applying the drift dynamics (equation 18) until the
469 particle hit a bound Z (set to -3 in the simulations), chosen to produce a similar range
470 of wait times as observed in data.

471 To determine what patterns we would be able to see in our rat data, we analyzed the
472 simulated dataset for models with each source of timing noise as though we did not know
473 generative x_0 , but could only observe the stimulus, choice accuracy, and willingness to
474 wait on a given trial. We binned trials by the evidence supporting choice, $\frac{\Delta}{\Sigma} \times \text{correct}$.
475 We then plotted conditioned kernel density estimates of willingness to wait for each bin.
476 We also computed the mean, standard deviation, and coefficient of variation (standard
477 deviation divided by mean) in each bin.

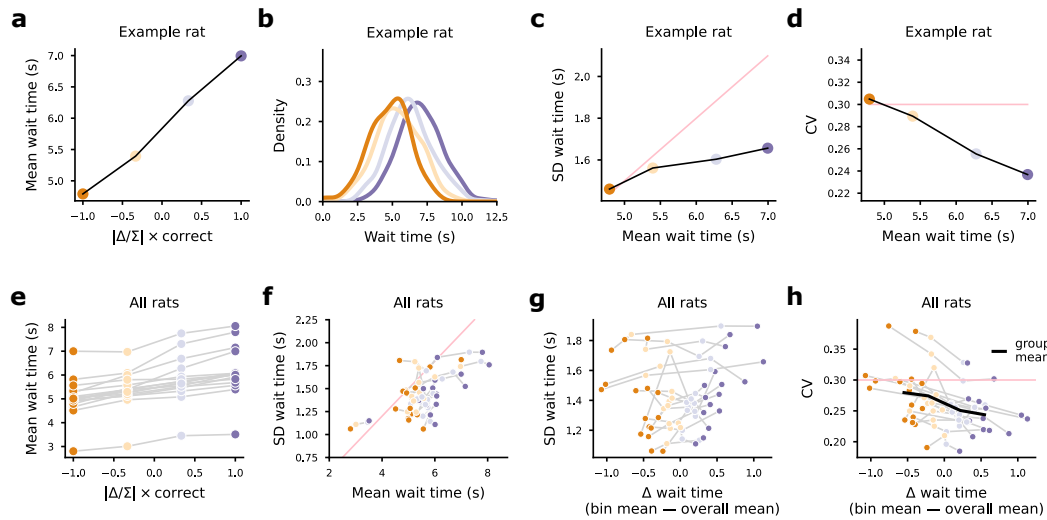


Figure 6: Scale invariant noise is not the dominant source of variability in the data (A) Average wait times for an example rat for correct probe trials and proportionally sampled error trials as a function of evidence for chosen option bin. (B) Kernel density estimate for the wait times in each bin from A (the same colormap is used to indicate evidence bin). (C) Standard deviation of the wait times in each bin (points connected by black trace) is overlaid on the predicted relationship between wait time standard deviation and mean under scale invariance with a coefficient of variation equal to .3 (pink trace). (D) Coefficient of variation (the standard deviation divided by the mean) in each bin is plotted as a function of mean wait time in each bin for the data (points) and for the scale invariant model shown in C (pink trace). (E) All rats' mean wait times plotted as a function of evidence for chosen option as in A. (F) All rats' wait time standard deviations plotted against mean wait time as in C. Again, overlaid on the prediction from scale invariance. (G) Standard deviation in each bin for all rats plotted against the difference between the mean wait time in the bin and the average of the bins for that rat. (H) Coefficient of variation plotted as in G for all rats with the mean across rats overlaid (black trace). If scale invariance was the dominant source of noise, each rat's trace should be flat.

478 All models achieved increasing mean willingness to wait as a function of evidence
 479 supporting choice. But, the models had different relationships between standard deviation
 480 and mean within each bin. The variable drift model produced a roughly proportional
 481 increase in standard deviation as the mean grew, corresponding to flat coefficient of
 482 variation, as expected under scale invariant noise (Fig 5a). But, there was additional
 483 noise across all bins that arose from variability in confidence (x_0) within each bin. The
 484 diffusion noise model produced a slightly sublinear increase in standard deviation as the
 485 mean grew, corresponding to a decreasing coefficient of variation (Fig 5b). Finally, the
 486 bound noise model produced a distinctly sublinear increase in standard deviation as the
 487 mean grew, with almost no increase between the last two evidence bins (Fig 5c). This
 488 corresponded to a dramatic decline in the coefficient of variation for the simulated data.

489 **Scalar variability is not the dominant source of noise in rat waiting data**

490 To test whether scale invariant timing noise was the dominant noise source affecting rats'
491 waiting time decisions, we analyzed the rat data using the analysis methods used to study
492 the simulated data. We analyzed the mean and standard deviation of probe trial wait
493 times in bins of normalized evidence strength favoring the rat's choice ($\Delta/\Sigma \times \text{correct}$).
494 The average wait time in each bin is shown for an example rat in Figure 6a. The
495 distribution of wait times in each of these bins is shown for the example rat in Figure 6b.
496 We compared the standard deviation of wait times in each of the bins to the mean (Fig 6c)
497 and computed the coefficient of variation in each bin (Fig 6d) for the example rat. The
498 pattern we observed was not consistent with the simulated data for the model with scale
499 invariant timing variability caused by noisy drift. Instead, the pattern we observed was
500 sublinear, with minimal increase in standard deviation between the last two bins.

501 We repeated this procedure for all rats. Average wait time as a function of binned
502 evidence for choice is shown for all rats in Figure 6e. The standard deviations as a
503 function of average wait time are plotted for all rats in Figure 6f (as in Figure 6c). To
504 compare across rats, we subtracted off the average wait time of the bins from each rat's
505 data (Fig 6g). We find a consistent pattern across rats that the relationship between
506 standard deviation and mean wait time is flatter than expected under the scale invariant
507 model. We then examined the coefficients of variation along this axis and plotted them
508 together for comparison (Fig 6h). Across rats, we see a consistent downward trend in the
509 coefficients of variation, inconsistent with scale invariant timing noise. This suggests that
510 other sources of noise dominate any scale invariant noise that exists in our rats' behavior.
511 Additionally, the standard deviation appears to increase more slowly than expected if
512 diffusion process was the dominant source of noise. Qualitatively, the variability in our
513 data appears most consistent with the model in which variability in the bound dominates.
514 This variability may stem from noise, but may also stem from continual learning of the
515 appropriate bound setting as a recency-weighted average of the reward rate history^{11,14}.

516 **Discussion**

517 We trained rats to perform a task requiring auditory evidence accumulation¹⁶ combined
518 with a post-decision temporal wager designed to assess their decision confidence⁵. The
519 time that animals are willing to wait for a delayed reward after making a decision has
520 become a popular proxy measurement of decision confidence, because we know that
521 optimal agents wait longer for rewards they are more confident they will receive^{5,7,10}.
522 However, willingness to wait in optimal agents is also influenced by the maximum possible
523 reward rate in the environment, which is in turn influenced by many environmental
524 statistics. These statistics determine the optimal overall average willingness to wait, but
525 they are not explicitly accounted for in previous studies of confidence-guided waiting.

526 Here, we developed an expression for the reward rate in the environment which made
527 all of the relevant environmental statistics explicit. Using this environmental reward rate,
528 we derived an expression for the conditions under which reward was maximized in the
529 environment. This generalized the marginal value theorem¹² into the case of stochastic
530 rewards¹⁵ with arbitrary initial expectations about the probability of eventual reward.
531 This work made it possible to test whether rats performing this task achieved overall
532 reward-rate-maximization, which we refer to as "optimal" behavior.

533 One of the key statistics that determines the optimal overall average willingness to
534 wait is the travel time incurred when deciding to move on from a given reward opportunity

535 and pursue the next. We observed that our animals were willing to wait longer than fully
536 optimal agents who minimized travel time and then maximized reward rate by finding the
537 best overall willingness to wait for that minimized travel time. Compared to these agents,
538 our animals “overharvested” reward on each trial, as has been seen in other studies of
539 foraging behaviors^{17,14}. Unlike fully optimal agents, our animals took longer to travel
540 between reward opportunities than was strictly necessary. We asked whether their waiting
541 behavior was optimal if we treated their travel times as constrained, meaning that they
542 had minimized travel time to the best of their ability and then optimized willingness to
543 wait given those travel times. We found that when we treated travel times as constrained,
544 the rats’ waiting behavior was near optimal.

545 One limitation of our study is that we don’t know for certain that our rats’ have
546 minimized travel time to the best of their ability. It is possible that rats could decrease
547 time between trials and achieve a higher reward rate. One way that future studies could
548 test this would be to impose a longer minimum intertrial interval and test whether rats
549 found a new willingness to wait that was optimal for the increased travel times.

550 To understand how our rats achieved near-optimal wait times, we developed a model
551 of the wait time decision process, which sought to capture the unfolding cognitive state
552 throughout the decision. Taking inspiration from the success of the drift diffusion model
553 in modeling two alternative decisions^{27,20,24}, we used the sequential probability ratio
554 test¹⁸ to develop an optimal update rule for a decision variable that can control the
555 port-leaving decision. This produced a continuously evolving cognitive process model
556 for controlling port-leaving time via the linear drift of a decision variable toward a fixed
557 bound. This process model provided a tractable algorithm for implementing optimal
558 waiting in the brain in which each of the separate parameters could be learned from
559 experience.

560 The model also allowed us to consider sources of variability that might contribute
561 to the decision process, drawing on extensive work on pacemaker accumulator models
562 of timing behavior²². To understand the mapping between willingness to wait and
563 confidence, it is useful to know what sources of variability are contributing to wait times.
564 Previous work has assumed that variability in the timing of willingness to wait would be
565 dominated by the scale invariant property⁵ in which the standard deviation of observed
566 wait times should be proportional to the animals’ desired wait time²¹. However, this
567 assumption had not been tested. We compared three models of timing variability in
568 the waiting decision process. The first was a noisy drift model, which produced scale
569 invariant timing noise. The second was a diffusion noise model, in which timing noise
570 grew like the square root of the interval to be timed. And finally, a noisy bound model,
571 in which timing noise was constant across desired wait times. We found that our data
572 was most consistent with the model dominated by bound variability.

573 While our model provides an improved description of the port-leaving decision process,
574 there are several avenues of possible improvement to the model that we should consider
575 in the future. First, there are well-documented aspects of the port choice decision process
576 that are not being accounted for here, including its evolution in time¹⁶, effects of trial
577 history²⁸, and change in the parameters of the decision process from trial to trial^{29,3}.
578 Second, there may be postprocessing of the stimulus following choice that leads to
579 evolution of confidence independently from that instructed by the environment³⁰. Finally,
580 just as the parameters governing the choice process may evolve from trial to trial, the
581 same may happen for the wait time decision either due to learning or changes in internal
582 state like increasing satiety or patience³¹. Indeed, we speculate that continual learning of
583 the bound controlling port-leaving may explain the variability we observed in our data.

584 Our model provides a tractable algorithm for solving this task, which can produce
585 optimal behavior. The model can also produce a variety of forms of variability and
586 deviation from optimality, which we have used to better understand the sources of
587 variability in confidence-guided waiting decisions. Future work investigating the neural
588 basis of confidence computations using the confidence-guided waiting paradigm should
589 seek to link neural activity and perturbations of brain regions to the parameters of a
590 dynamic model of the internal cognitive process for deciding when to give up and move
591 on, like the one developed here. Using such a model will increase the interpretive power
592 of experiments using this paradigm to understand how the brain computes confidence
593 estimates and uses them to guide subsequent behavior.

594 Methods

595 Subjects

596 Animal use procedures were approved by the Princeton University Institutional Animal
597 Care and Use Committee and carried out in accordance with NIH standards. All subjects
598 were adult male Long Evans rats bred either at Princeton Neuroscience Institute (VGAT-
599 ReaChR rats) or by one of the following vendors (wild type rats): Taconic, Hilltop
600 and Harlan, USA. Rats were pair-housed unless implanted with infusion cannulae at
601 which point they were single-housed. Rats were placed on a water restriction schedule to
602 motivate them to perform the task for water rewards.

603 Behavioral tasks

604 **Poisson Clicks** We trained rats on the Poisson Clicks task¹⁶ with a post-decision wait
605 time wager^{5,6} using an automated training protocol. Throughout training, rats were put
606 on a controlled water schedule where they received at least 3% of their weight every day.
607 Rats trained each day in training sessions of around 120 minutes.

608 In the final stage of training, each trial began with the illumination of a center nose
609 port by an LED light inside the port. This LED indicated that the rat could initiate a
610 trial by placing its nose into the center port. Rats were required to keep their nose in
611 the center port (“center fixation”) for a fixed duration until the LED turned off as a “go”
612 signal. During center fixation, two trains of randomly-timed auditory clicks were played
613 from speakers on either side of the center port after a variable delay. The duration of the
614 click trains was uniformly distributed. The two click trains were each associated with
615 one of two side ports and clicks in each click train were generated using different Poisson
616 rates. For a given rat, the two generative rates always summed to a fixed value (20 or 40
617 clicks s^{-1}).

618 After the “go” signal, rats made a port choice by poking their nose into one of the
619 two side ports. If they exited from the center port before the “go” signal, the trial was
620 considered a violation and they experienced a white noise stimulus followed by a short
621 time out. These trials did not yield decisions or wait times, but did contribute to travel
622 times.

623 Choices were considered correct, and potentially rewarded, if they corresponded to
624 the click train with the greater number of clicks, which corresponds to a noiseless ideal
625 observer’s estimate of the larger click rate.

626 **Confidence-guided waiting** Rewards were only delivered if the rat stayed at the side
627 port until a reward time t_r drawn from an exponential distribution between a minimum
628 $t_{r,\min} \in (.05s, .5s)$ and maximum $t_{r,\max} > 15s$ with time constant $\tau = 1.5s$. The resulting
629 mean reward delay was $\langle t_r \rangle = t_{r,\min} + \tau$. After errors, no feedback was delivered. Instead,
630 the animal had to eventually give up on waiting for reward and start a new trial.

631 With probability $\zeta \in (.05, .15)$, the trial was turned into a probe trial by setting
632 $t_r = 100s$. We did not allow multiple probe trials to occur consecutively. These probe
633 trials allowed us to observe port-leaving times on a subset of correct trials when they
634 might otherwise have been censored by reward delivery. Rats were given a grace period
635 between 500 and 1500ms for leaving and returning to the choice port. If they withdrew
636 from the reward port for longer than this grace period, reward was no longer available.
637 If the rat returned to the center port, during or after the grace period, a new trial was

638 immediately initiated. If they returned to the chosen side port after the grace period, or
639 entered the opposite side port at any time, the possibility of reward delivery was removed.
640 For analysis of uncensored wait times, we focused on trials where the rat initiated a new
641 trial by center poking within 2 seconds of leaving the side port.

642 Shaping

643 We shaped the animals by first training them to perform the Poisson Clicks task via a
644 standardized set of training stages. We then added the reward delay component. First,
645 fixed feedback delays were introduced on both correct and error trials and grew in each
646 trial until they reached $t_{r,\min}$. Then, the error feedback delay was incremented from trial
647 to trial until the rat never waited long enough to get the error feedback. At that point,
648 the error feedback delay was set to 100s. Next, the reward delays were randomized by
649 gradually increasing the exponential time constant τ and the maximum delay time $t_{r,\max}$.
650 When the $t_{r,\max}$ was larger than the rat's longest waiting times, we set it to 100s. When
651 τ reached its target value, we introduced probe trials. We did not allow multiple probe
652 trials to occur in a row.

653 Inclusion criteria

654 Rats trained on this task were included in this study if they had more than 30 sessions
655 that met the session inclusion criteria and if the fraction of unrewarded trials that ended
656 with a re-initiating center poke (as opposed to re-entry in the chosen side port or entry
657 into the opposite side port) met a minimum threshold of 55%. The session inclusion
658 criteria required that the rat perform at least 150 trials with an overall accuracy rate
659 exceeding 60%. In order to prevent the rats from developing biases towards particular
660 side ports, an anti-biasing algorithm detected biases and probabilistically generated trials
661 with the correct answer on the non-favored side.

662 Psychometric curves

663 Behavioral sensitivity was assessed using psychometric curves. The probability of choosing
664 the rightward port was computed as a function of the binned normalized click difference
665 ($\frac{\Delta}{\Sigma} \equiv \frac{\#R-\#L}{\#R+\#L}$). We fit psychometric curves with 2 parameters, a bias parameter b and a
666 noise parameter σ , for all rats as a function of the normalized click difference. We fit the
667 data by minimizing the negative log likelihood across trials where the probability of a
668 rightward choice on a given trial was given by

$$P(\text{go right}) = .5 \left(1 + \operatorname{erf} \left(-\frac{b - \frac{\Delta}{\Sigma}}{\sigma\sqrt{2}} \right) \right). \quad (23)$$

669 Wait time chronometric curves

670 Wait time modulation was assessed using error trials and correct probe trials to create
671 wait time chronometric curves, which relate mean wait time to the strength of evidence
672 supporting the chosen option. Strength of evidence supporting choice was computed as
673 $\frac{\Delta}{\Sigma} \times y$, where $y = \pm 1$ with positive values for correct port choices and negative values for
674 incorrect port choices. The trials with the most evidence supporting the chosen option
675 have large, positive values and the trials with the most evidence against the chosen option
676 have large, negative values. The most difficult trials, with the least evidence weighing on

677 the choice, have small magnitudes. We expect confidence to increase monotonically along
678 this axis. We fit a line to each rat’s wait times in the space of normalized click difference
679 supporting the choice.

680 **Optimality analysis**

681 To test whether rats’ waiting times maximized overall reward rate, we found the optimal
682 overall average wait time, t_w^* , by evaluating equation 14 for each rat. To do this, we
683 estimated the relevant terms contributing to equation 14 from the rats’ datasets: α was
684 the fraction of non-probe trials in the rat’s dataset, C_0 was the fraction of correct trials,
685 τ and $t_{r,\min}$ were estimated from the reward delays scheduled for the rat, and t_0 was
686 estimated from either the mean travel time achieved by the rat (after excluding the
687 longest 1% of travel times, because the rats occasionally fully disengaged from the task
688 for long periods of time), or the minimum travel time achieved by the rat. Because we
689 are only interested in average overall waiting time here, we don’t need to consider the
690 variations in wait time associated with confidence. Therefore, this agent was constrained
691 to wait the same amount of time on every trial, which allowed us to avoid making choices
692 about how to capture variations in confidence for this analysis. We used a root finding
693 algorithm to evaluate t_w^* for a given set of task statistics. We compared the willingness
694 to wait for each of these agents to the average waiting times for the corresponding rat
695 in correct probe trials and in a subset of error trials (subsampling to ensure that the
696 frequency of error trials in this comparison matched that in the overall dataset). We also
697 measured the optimal agent’s reward rate and the fraction of the agent’s reward rate
698 achieved by the rat.

699 **Process model simulations with candidate noise sources**

700 We used euler integration to simulate the wait time decision process for three candidate
701 noise models. In all simulations, we sampled 50,000 trial stimulus strengths, s , with
702 replacement from the dataset of an example rat. We then generated a percept, $p = s + \xi$,
703 for each trial, by adding Gaussian noise, $\xi \sim \mathcal{N}(0, \sigma_s^2)$, to the stimulus. The model made
704 a rightward choice if the resulting percept was greater than a decision boundary, which
705 we set to zero (i.e., $p > b$ for $b = 0$). Given this percept, we generated confidence levels
706 according to equation 22. We then produced a corresponding x_0 and updated it in 25
707 millisecond timesteps ($\Delta t = .025s$) according to equation 19. The drift was set to it’s
708 optimal setting $A = -\frac{1}{\tau}$ (per equation 18) and the bound was set to $Z = -3$ to produce
709 mean wait times across trials that roughly matched the example rat’s. To produce a
710 model with scale invariant timing noise, and specifically a coefficient of variation of 0.3, as
711 in Lak et al.⁵, we set the drift on each trial to be $A^{\text{trial}} = -\frac{1}{\hat{\tau}}$ where $\hat{\tau} \sim \mathcal{N}(\tau, 0.3\tau)$. To
712 produce a model with diffusion noise, we added Gaussian noise in each time step drawn
713 from $\mathcal{N}(0, c\sqrt{\Delta t})$ with $c = 0.3\sqrt{ZA}$ to produce an equivalent level of noise at $x_0 = 0$
714 as produced under the scale invariant model. To produce a model with constant noise,
715 we sampled a different bound on each trial $Z^{\text{trial}} \sim \mathcal{N}(Z, 0.3 \cdot |Z|)$. The magnitude of
716 noise was again chosen to produce the same noise level as the other models for $x_0 = 0$.
717 We recorded each models willingness to wait on each trial as the timestep in which the
718 particle x first crossed the bound Z .

719 **Analysis of variability in simulations and rat data**

720 We used our simulations to ask what patterns of variability would be expected as a
721 function of the stimulus. This was useful for analyzing rat data in which the confidence
722 level and x_0 level are unknown. To do this, we binned trials by the evidence supporting
723 the chosen option for both the simulated data and the rat data. Within these bins, we
724 computed kernel density estimates of the distribution, as well as computing the mean,
725 standard deviation, and coefficient of variation (ratio of standard deviation to mean) of the
726 wait times in each bin. These produced distinct patterns for each of the candidate models,
727 which we then compared qualitatively to the rat data. In particular, the assumption of
728 scale invariance predicted a flat coefficient of variation, which we did not observe in the
729 rat data. Instead, our data was most consistent with the constant bound noise in which
730 the standard deviation in each bin grows slowly as mean wait time increases.

731 **Acknowledgements**

732 We thank Athena Akrami, Adrian Bondy, Diksha Gupta, Thomas Luo, all other members
733 of the Brody lab, as well as Lukas Braun, Nathaniel Daw, Javier Masís, Stefano Sarao
734 Mannelli, and Pat Simen for helpful conversations and suggestions. TB acknowledges
735 support by NIH grant T32 MH 65214-16. This work was supported by a grant from the
736 Simons Foundation (Grant # 542953) awarded to CB, as well as NIH grant R01MH108358
737 awarded to CB.

738 **Author contributions**

739 T.B. and C.K. developed the rat training protocol. T.B. managed rat training and care.
740 T.B. and A.P. derived the equations and models. T.B. analyzed the data. T.B., A.P. and
741 C.B. wrote the manuscript. C.B. oversaw all aspects of the project.

742 **Competing interests statement**

743 The authors declare no competing interests

744 References

- 745 [1] Alexandre Pouget, Jan Drugowitsch, and Adam Kepecs. Confidence and certainty:
746 distinct probabilistic quantities for different goals. *Nat. Neurosci.*, 19(3):366–374,
747 March 2016.
- 748 [2] Balázs Hangya, Joshua I. Sanders, and Adam Kepecs. A mathematical framework
749 for statistical decision confidence. *Neural Computation*, 28(9):1840–1858, 09 2016.
- 750 [3] Jan Drugowitsch, André G. Mendonça, Zachary F. Mainen, and Alexandre Pouget.
751 Learning optimal decisions with confidence. *Proceedings of the National Academy of
752 Sciences*, 116(49):24872–24880, 2019. doi:10.1073/pnas.1906787116.
- 753 [4] Jan Drugowitsch, Rubén Moreno-Bote, Anne K. Churchland, Michael N. Shadlen, and
754 Alexandre Pouget. The cost of accumulating evidence in perceptual decision making.
755 *Journal of Neuroscience*, 32(11):3612–3628, 2012. doi:10.1523/JNEUROSCI.4010-
756 11.2012.
- 757 [5] Armin Lak, Gil M. Costa, Erin Romberg, Alexei A. Koulakov, Zachary F.
758 Mainen, and Adam Kepecs. Orbitofrontal cortex is required for opti-
759 mal waiting based on decision confidence. *Neuron*, 84(1):190–201, 2014.
760 doi:<https://doi.org/10.1016/j.neuron.2014.08.039>.
- 761 [6] Adam Kepecs, Naoshige Uchida, Hatim A. Zariwala, and Zachary F. Mainen. Neural
762 correlates, computation and behavioural impact of decision confidence. *Nature*, 455
763 (7210):227–231, August 2008. doi:10.1038/nature07200.
- 764 [7] Paul Masset, Torben Ott, Armin Lak, Junya Hirokawa, and Adam Kepecs. Behavior-
765 and Modality-General representation of confidence in orbitofrontal cortex. *Cell*, 182
766 (1):112–126.e18, July 2020.
- 767 [8] A Stolyarova, M Rakhshan, E E Hart, T J O’Dell, M A K Peters, H Lau, A Soltani,
768 and A Izquierdo. Contributions of anterior cingulate cortex and basolateral amygdala
769 to decision confidence and learning under uncertainty. *Nat. Commun.*, 10(1):4704,
770 October 2019.
- 771 [9] Hannah R. Joo, Hexin Liang, Jason E. Chung, Charlotte Geaghan-Breiner, Jiang Lan
772 Fan, Benjamin P. Nachman, Adam Kepecs, and Loren M. Frank. Rats use memory
773 confidence to guide decisions. *Current Biology*, 31(20):4571–4583, 2021/11/11 2021.
774 doi:10.1016/j.cub.2021.08.013.
- 775 [10] K. Schmack, M. Bosc, T. Ott, J. F. Sturgill, and A. Kepecs. Striatal dopamine
776 mediates hallucination-like perception in mice. *Science*, 372(6537), April 2021.
777 doi:10.1126/science.abf4740.
- 778 [11] Andrew Mah, Shannon S. Schiereck, Veronica Bossio, and Christine M. Constantino-
779 ple. Distinct value computations support rapid sequential decisions. *Nature Com-
780 munications*, 14(1):7573, November 2023. doi:10.1038/s41467-023-43250-x. Number:
781 1 Publisher: Nature Publishing Group.
- 782 [12] E L Charnov. Optimal foraging, the marginal value theorem. *Theor. Popul. Biol.*, 9
783 (2):129–136, April 1976.
- 784 [13] David W Stephens and John R Krebs. *Foraging Theory*, volume 1. Princeton
785 University Press, 1986. ISBN 9780691084411.
- 786 [14] Sara M. Constantino and Nathaniel D. Daw. Learning the opportunity cost of time

- 787 in a patch-foraging task. *Cognitive, Affective, and Behavioral Neuroscience*, 15(4):
788 837–853, April 2015. doi:10.3758/s13415-015-0350-y.
- 789 [15] John McNamara. Optimal patch use in a stochastic environment. *Theoretical*
790 *Population Biology*, 21(2):269–288, April 1982. doi:10.1016/0040-5809(82)90018-1.
- 791 [16] Bingni W. Brunton, Matthew M. Botvinick, and Carlos D. Brody. Rats and humans
792 can optimally accumulate evidence for decision-making. *Science*, 340(6128):95–98,
793 2013. doi:10.1126/science.1233912.
- 794 [17] Gary A Kane, Aaron M Bornstein, Amitai Shenhav, Robert C Wilson, Nathaniel D
795 Daw, and Jonathan D Cohen. Rats exhibit similar biases in foraging and intertempo-
796 ral choice tasks. *eLife*, 8:e48429, September 2019. doi:10.7554/eLife.48429. Publisher:
797 eLife Sciences Publications, Ltd.
- 798 [18] A Wald. Sequential tests of statistical hypotheses. *Ann. Math. Stat.*, 16(2):117–186,
799 1945.
- 800 [19] Roger Ratcliff and Gail McKoon. The diffusion decision model: theory and data
801 for two-choice decision tasks. *Neural Computation*, 20(4):873–922, April 2008.
802 doi:10.1162/neco.2008.12-06-420.
- 803 [20] Joshua I. Gold and Michael N. Shadlen. The neural basis of de-
804 cision making. *Annual Review of Neuroscience*, 30(1):535–574, 2007.
805 doi:10.1146/annurev.neuro.29.051605.113038.
- 806 [21] John Gibbon. Scalar expectancy theory and weber’s law in animal timing. *The*
807 *psychological review.*, 84(3), 1977.
- 808 [22] Patrick Simen, Francois Rivest, Elliot A. Ludvig, Fuat Balci, and Peter Killeen.
809 Timescale invariance in the pacemaker-accumulator family of timing models. *Timing*
810 *& Time Perception*, 1(2):159–188, January 2013. doi:10.1163/22134468-00002018.
811 Publisher: Brill.
- 812 [23] Patrick Simen, Ksenia Vlasov, and Samantha Papadakis. Scale (in)variance in a
813 unified diffusion model of decision making and timing. *Psychological Review*, 123(2):
814 151–181, 2016. doi:10.1037/rev0000014.
- 815 [24] Roger Ratcliff. A theory of memory retrieval. *Psychological Review*, 85(2):59–108,
816 March 1978. doi:10.1037/0033-295x.85.2.59.
- 817 [25] Jacob D. Davidson and Ahmed El Hady. Foraging as an evidence accu-
818 mulation process. *PLOS Computational Biology*, 15(7):e1007060, July 2019.
819 doi:10.1371/journal.pcbi.1007060.
- 820 [26] Benjamin Y Hayden, John M Pearson, and Michael L Platt. Neuronal basis of
821 sequential foraging decisions in a patchy environment. *Nature Neuroscience*, 14(7):
822 933–939, June 2011. doi:10.1038/nn.2856.
- 823 [27] Rafal Bogacz, Eric Brown, Jeff Moehlis, Philip Holmes, and Jonathan Cohen. The
824 physics of optimal decision making: A formal analysis of models of performance in
825 two-alternative forced-choice tasks. *Psychological Review*, 2006.
- 826 [28] Diksha Gupta, Brian DePasquale, Charles D. Kopec, and Carlos D. Brody. Trial-
827 history biases in evidence accumulation can give rise to apparent lapses in decision-
828 making. *Nature Communications*, 15(1), January 2024. doi:10.1038/s41467-024-
829 44880-5.

- 830 [29] Nicholas A. Roy, Ji Hyun Bak, Athena Akrami, Carlos D. Brody, and Jonathan W.
831 Pillow. Extracting the dynamics of behavior in sensory decision-making experiments.
832 *Neuron*, 109(4):597–610.e6, February 2021. doi:10.1016/j.neuron.2020.12.004.
- 833 [30] Joaquin Navajas, Bahador Bahrami, and Peter E Latham. Post-decisional accounts
834 of biases in confidence. *Current Opinion in Behavioral Sciences*, 11:55–60, October
835 2016. doi:10.1016/j.cobeha.2016.05.005.
- 836 [31] Michael Bukwich, Malcolm G. Campbell, David Zoltowski, Lyle Kingsbury, Mom-
837 chil S. Tomov, Joshua Stern, HyungGoo R. Kim, Jan Drugowitsch, Scott W. Linder-
838 man, and Naoshige Uchida. Competitive integration of time and reward explains
839 value-sensitive foraging decisions and frontal cortex ramping dynamics. *bioRxiv*,
840 September 2023. doi:10.1101/2023.09.05.556267.

Supplementary Information

A cognitive process model captures near-optimal confidence-guided waiting in rats

J Tyler Boyd-Meredith, Alex T Piet, Chuck D Kopec, and Carlos D Brody

Contents

1	Derivation of reward-rate-maximizing behavior	2
1.1	Expected reward per trial	2
1.2	Expected time per trial	3
1.2.1	Reward maximization doesn't depend on consumption time . . .	4
1.3	Derivation of posterior belief that reward will be delivered	5
1.4	Derivation of optimal willingness to wait	5

1 Derivation of reward-rate-maximizing behavior

The total reward rate in the task is defined as the expected reward per trial, $g(t_w)$, divided by the expected time spent in each trial, T_{total} :

$$RR_{\text{total}} = \frac{g(t_w)}{T_{\text{total}}(t_w)}. \quad (\text{S1})$$

To maximize the reward rate, we find the condition such that its derivative is zero, $\frac{\partial RR_{\text{total}}}{\partial t_w} = 0$. We compute the derivative using the quotient rule

$$\frac{\partial RR_{\text{total}}}{\partial t_w} = \frac{\frac{\partial g(t_w)}{\partial t_w} T_{\text{total}}(t_w) - \frac{\partial T_{\text{total}}(t_w)}{\partial t_w} g(t_w)}{T_{\text{total}}^2(t_w)},$$

which is equal to zero when

$$\begin{aligned} \frac{\partial g(t_w^*)}{\partial t_w} T_{\text{total}}(t_w^*) &= \frac{\partial T_{\text{total}}(t_w^*)}{\partial t_w} g(t_w^*) \\ \frac{\partial g(t_w^*)}{\partial t_w} \frac{\partial t_w}{\partial T_{\text{total}}(t_w^*)} &= \frac{g(t_w^*)}{T_{\text{total}}(t_w^*)} \\ \frac{\partial g}{\partial T_{\text{total}}}(t_w^*) &= RR_{\text{total}}^* \end{aligned} \quad (\text{S2})$$

where t_w^* is the reward-maximizing willingness to wait and RR_{total}^* is the reward achieved at t_w^* .

To find a solution for t_w^* for a given set of task statistics from equation S2, we need expressions for $g(t_w)$ and $T_{\text{total}}(t_w)$. To do so, we will use the notation introduced in the main text to simplify these expressions. We will use

$$r_w \equiv (t_w, t_w + \delta t) \quad (\text{S3})$$

to indicate whether reward is set to be delivered in some infinitesimal timestep δt beginning at time t_w . Then, we can indicate whether reward is set to be delivered before time t_w using the sum

$$R_w \equiv \sum_{i=0}^{w-1} r_i. \quad (\text{S4})$$

We will use the negation, $\neg R_w$, to indicate when no reward is delivered by time t_w .

1.1 Expected reward per trial

Because at most 1 reward is delivered per trial and it always has the same magnitude, we can set the reward magnitude to 1 and make the expected reward per trial equivalent to the probability of reward in a trial

$$g(t_w) \equiv P(R_w). \quad (\text{S5})$$

Our reward distribution is exponential, meaning that, given that reward is coming, the delivery times, t_r , are distributed according to

$$P(t_w | R_\infty) = \frac{1}{\tau} e^{-(t_w - t_{r,\text{min}})/\tau}$$

where τ is an experimenter-specified time constant and $t_{r,\min}$ is an experimenter-specified minimum reward time (equation 1 in the main text). The probability of receiving a reward before time t_w depends on this distribution and on the probability that reward will be delivered on this trial, $P(R_\infty)$. The probability that reward will be delivered on this trial is estimated based on the decision confidence, C_0 , and the probability that a trial is not a probe a trial, α :

$$P(R_\infty) = \alpha C_0$$

(equation 4 in the main text). The expected reward per trial as the probability of receiving a reward before time t_w is the cumulative density function for the exponential given that the reward is coming times the prior probability that the reward is coming, $P(R_\infty)$:

$$\begin{aligned} g(t_w) &= P(R_w | R_\infty)P(R_\infty) \\ &= \alpha C_0(1 - e^{-(t_w - t_{r,\min})/\tau}) \end{aligned} \quad (\text{S6})$$

where we have used the cumulative density for an exponential to compute $P(R_w | R_\infty)$.

1.2 Expected time per trial

The expected time per trial can be broken into three epochs: the time between leaving the reward port on the previous trial and entering a reward port on the current trial, t_0 , the time spent at the port on the current trial, t_{port} , and the time spent consuming reward on the current trial, t_{drink} . Adding together the expected duration of each epoch, we get:

$$T_{\text{total}} = t_0 + \mathbb{E}[t_{\text{port}} | t_w] + \mathbb{E}[t_{\text{drink}} | t_w]. \quad (\text{S7})$$

The first quantity is referred to as the “travel time” and, for the reward-maximizing agent, does not depend on t_w . The other two quantities depend on whether reward is set to be delivered and how long the agent is willing to wait. The consumption time, t_{drink} , is either 0, if no reward is received, or a constant, if reward is delivered. Its expectation can be written

$$\mathbb{E}[t_{\text{drink}} | t_w] = t_{\text{drink}}P(R_w) = t_{\text{drink}}g(t_w). \quad (\text{S8})$$

We will show that $\mathbb{E}[t_{\text{drink}} | t_w]$ can be ignored for the reward maximization process.

Expected time at the port To compute the expected time at the port, $\mathbb{E}[t_{\text{port}} | t_w]$, we will separately consider trials in which reward is not set to be delivered and trials in which reward will be delivered if the agent waits long enough. Marginalizing over these possibilities gives us

$$\mathbb{E}[t_{\text{port}} | t_w] = \mathbb{E}[t_{\text{port}} | t_w, \neg R_\infty]P(\neg R_\infty) + \mathbb{E}[t_{\text{port}} | t_w, R_\infty]P(R_\infty) \quad (\text{S9})$$

When no reward is set to be delivered, the agent always gives up and moves on at the time t_w :

$$\mathbb{E}[t_{\text{port}} | t_w, \neg R_\infty] = t_w. \quad (\text{S10})$$

In trials where reward is set to be delivered at some time t_r , the time at the port can take one of two values. If the agent is not willing to wait long enough to get the reward, the agent will give up before the reward is delivered and we will again observe time t_w spent at the port:

$$\mathbb{E}[t_{\text{port}} \mid t_w, R_\infty, t_r > t_w] = t_w \quad (\text{S11})$$

If the agent is willing to wait long enough to get the reward, the reward delivery will censor the agent's willingness to wait and we will observe t_r time spent at the port. To compute the expected port time for this trial type, we need to marginalize over the possible values that t_r can take, as follows:

$$\begin{aligned} \mathbb{E}[t_{\text{port}} \mid t_w, R_\infty, t_r \leq t_w] &= \int_{t_{r,\min}}^{t_w} t \cdot P(t_r = t \mid R_\infty, t_r \leq t_w) dt \\ &= \int_{t_{r,\min}}^{t_w} t \cdot \frac{P(t_r \leq t_w \mid R_\infty, t_r = t) P(t_r = t \mid R_\infty)}{P(t_r \leq t_w \mid R_\infty)} dt \\ &= \frac{1}{P(t_r \leq t_w \mid R_\infty)} \int_{t_{r,\min}}^{t_w} t P(t_r = t \mid R_\infty) dt \\ &= \frac{1}{P(t_r \leq t_w \mid R_\infty)} \int_0^{t_w} \frac{t}{\tau} e^{-(t-t_{r,\min})/\tau} dt \\ &= \frac{t_{r,\min} + \tau(1 - e^{-(t_w-t_{r,\min})/\tau}) - t_w e^{-(t_w-t_{r,\min})/\tau}}{P(t_r \leq t_w \mid R_\infty)} \quad (\text{S12}) \end{aligned}$$

Combining equations S11 and S12 and multiplying each by their probabilities, we can compute the expected time at the port on trials where reward is set to be delivered eventually:

$$\begin{aligned} \mathbb{E}[t_{\text{port}} \mid t_w, R_\infty] &= \mathbb{E}[t_{\text{port}} \mid t_w, R_\infty, t_r > t_w] P(t_r > t_w \mid t_w, R_\infty) + \\ &\quad \mathbb{E}[t_{\text{port}} \mid t_w, R_\infty, t_r \leq t_w] P(t_r \leq t_w \mid t_w, R_\infty) \\ &= t_w e^{-(t_w-t_{r,\min})/\tau} + t_{r,\min} + \tau(1 - e^{-(t_w-t_{r,\min})/\tau}) - t_w e^{-(t_w-t_{r,\min})/\tau} \\ &= t_{r,\min} + \tau(1 - e^{-(t_w-t_{r,\min})/\tau}) \quad (\text{S13}) \end{aligned}$$

Finally, we can combine the expected port time in trials where no reward is baited (equation S10) and trials where reward is set to be delivered if the agent waits long enough (equation S13) to get the expected time at the port overall:

$$\begin{aligned} \mathbb{E}[t_{\text{port}} \mid t_w] &= \mathbb{E}[t_{\text{port}} \mid t_w, -R_\infty] P(-R_\infty) + \mathbb{E}[t_{\text{port}} \mid t_w, R_\infty] P(R_\infty) \\ &= (1 - \alpha C_0) t_w + \alpha C_0 \left(t_{r,\min} + \tau \left(1 - e^{-(t_w-t_{r,\min})/\tau} \right) \right) \quad (\text{S14}) \end{aligned}$$

1.2.1 Reward maximization doesn't depend on consumption time

As mentioned above, we can ignore the consumption time in the reward maximization process, which simplifies equation S2. We will use $T(t_w) \equiv t_0 + \mathbb{E}[t_{\text{port}} \mid t_w]$, to represent the expected time spent searching for, but not consuming reward. We will use $RR \equiv \frac{g(t_w)}{T(t_w)}$ to represent the reward rate per time spent searching for reward. We can rewrite

equation S2 with the consumption times made explicit and show that it can be ignored

$$\begin{aligned} \frac{\partial g}{\partial T_{\text{total}}}(t_w^*) &= RR_{\text{total}}^* \\ \frac{\frac{\partial}{\partial t_w} g(t_w^*)}{\frac{\partial}{\partial t_w} T(t_w^*) + t_{\text{drink}} \frac{\partial}{\partial t_w} g(t_w^*)} &= \frac{g(t_w^*)}{T(t_w^*) + t_{\text{drink}} g(t_w^*)} \\ \frac{\frac{\partial}{\partial t_w} T(t_w^*)}{\frac{\partial}{\partial t_w} g(t_w^*)} + t_{\text{drink}} &= \frac{T(t_w^*)}{g(t_w^*)} + t_{\text{drink}} \\ \frac{\frac{\partial}{\partial t_w} g(t_w^*)}{\frac{\partial}{\partial t_w} T(t_w^*)} &= \frac{g(t_w^*)}{T(t_w^*)} \\ \frac{\partial g}{\partial T}(t_w^*) &= RR^*. \end{aligned}$$

We will use equation 8 to find the optimal waiting behavior.

1.3 Derivation of posterior belief that reward will be delivered

From equation 11 in the main text, we know that the posterior belief that reward will be delivered on a given trial after waiting for time t_w without receiving reward is

$$P(R_{\infty} | \neg R_w) = \frac{P(\neg R_w | R_{\infty})P(R_{\infty})}{P(\neg R_w)}.$$

The first term in the numerator is the probability that reward is not delivered by time t_w given that it will be delivered eventually, which is the survivor function of the exponential distribution (or 1 minus the CDF)

$$P(\neg R_w | R_{\infty}) = e^{-(t_w - t_{r,\min})/\tau}. \quad (\text{S15})$$

The denominator can be expressed as

$$\begin{aligned} P(\neg R_w) &= P(\neg R_w | \neg R_{\infty})P(\neg R_{\infty}) + P(\neg R_w | R_{\infty})P(R_{\infty}) \\ &= 1 - \alpha C_0 + \alpha C_0 e^{-(t_w - t_{r,\min})/\tau} \end{aligned} \quad (\text{S16})$$

where we have used equation S15 and the fact that $P(\neg R_{\infty}) = 1 - P(R_{\infty})$. Combining these expressions with the definition of $P(R_{\infty})$ (equation 4), we get:

$$P(R_{\infty} | \neg R_w) = \frac{\alpha C_0 e^{-(t_w - t_{r,\min})/\tau}}{1 - \alpha C_0 + \alpha C_0 e^{-(t_w - t_{r,\min})/\tau}}.$$

1.4 Derivation of optimal willingness to wait

We rearrange the terms of the optimality condition from equation 8 and use the expression we derived for the instantaneous reward expectation (equation 13) to find the

optimal willingness to wait

$$\begin{aligned}
 \frac{\partial g}{\partial T}(t_w^*) &= RR^* \\
 P(r_w \mid \neg R_w) &= RR^* \\
 \frac{1}{\tau} \cdot \frac{\alpha C_0 e^{-(t_w^* - t_{r,\min})/\tau}}{1 - \alpha C_0 + \alpha C_0 e^{-(t_w^* - t_{r,\min})/\tau}} &= RR^* \\
 \frac{\alpha C_0}{(1 - \alpha C_0) e^{(t_w^* - t_{r,\min})/\tau} + \alpha C_0} &= RR^* \tau \\
 (1 - \alpha C_0) e^{(t_w^* - t_{r,\min})/\tau} &= \frac{\alpha C_0}{RR^* \tau} - \alpha C_0 \\
 e^{(t_w^* - t_{r,\min})/\tau} &= \frac{\alpha C_0 - \alpha C_0 RR^* \tau}{(1 - \alpha C_0) RR^* \tau} \\
 e^{(t_w^* - t_{r,\min})/\tau} &= \frac{\alpha C_0 (1 - RR^* \tau)}{(1 - \alpha C_0) RR^* \tau} \\
 t_w^* &= t_{r,\min} + \tau \log \left(\frac{\alpha C_0}{1 - \alpha C_0} \frac{1 - RR^* \tau}{RR^* \tau} \right) \\
 t_w^* &= t_{r,\min} + \tau \left(\log \frac{\alpha C_0}{1 - \alpha C_0} - \log \frac{RR^* \tau}{1 - RR^* \tau} \right).
 \end{aligned}$$