

Optimal Response Vigor and Choice Under Non-stationary Outcome Values

Amir Dezfouli¹

School of Psychology, UNSW, Sydney, Australia

Bernard W. Balleine

School of Psychology, UNSW, Sydney, Australia

Richard Nock

Data61, the Australian National University and the University of Sydney, Australia

Abstract

¹Corresponding author. Email: a.dezfouli@unsw.edu.au. Address: School of Psychology, UNSW, Sydney, Australia

Within a rational framework a decision-maker selects actions based on the reward-maximization principle which stipulates that they acquire outcomes with the highest value at the lowest cost. Action selection can be divided into two dimensions: selecting an action from various alternatives, and choosing its vigor, i.e., how fast the selected action should be executed. Both of these dimensions depend on the values of outcomes, which are often affected as more outcomes are consumed together with their associated actions. Despite this, previous research has only addressed the computational substrate of optimal actions in the specific condition that the values of outcomes are constant. It is not known what actions are optimal when the values of outcomes are non-stationary. Here, based on an optimal control framework, we derive a computational model for optimal actions when outcome values are non-stationary. The results imply that, even when the values of outcomes are changing, the optimal response rate is constant rather than decreasing. This finding shows that, in contrast to previous theories, commonly observed changes in action rate cannot be attributed solely to changes in outcome value. We then prove that this observation can be explained based on uncertainty about temporal horizons; e.g., the session duration. We further show that, when multiple outcomes are available, the model explains probability matching as well as maximization strategies. The model provides, therefore, a quantitative analysis of optimal action and explicit predictions for future testing.

Keywords: choice, response vigor, reward learning, optimal actions

Introduction

According to normative theories of decision-making, actions made by humans and animals are chosen with the aim of earning the maximum amount of future reward whilst incurring the lowest cost (Marshall, 1890; von Neumann & Morgenstern, 1947). Within such theories individuals optimize their actions by learning about their surrounding environment so as to satisfy their long-term objectives. The problem of finding the optimal action is, however, argued to have two aspects: (1) choice, i.e., deciding which action to select from several alternatives; and (2) vigor, i.e., deciding how fast the selected action should be executed. For a rat in a Skinner box, for example, the problem of finding the optimal action involves selecting a lever (choice) and deciding at what rate to respond on that lever (vigor). High response rates can have high costs (e.g., in terms of energy consumption), whereas a low response rate could have an opportunity cost if the experimental session ends before the animal has earned sufficient reward. Optimal actions provide the right balance between these two factors and, based on the reinforcement-learning framework and methods from optimal control theory, the characteristics of optimal actions and their consistency with various experimental studies have been previously elaborated (Dayan, 2012; Niv, Daw, Joel, & Dayan, 2007; Niyogi, Shizgal, & Dayan, 2014; Salimpour & Shadmehr, 2014).

These previous models have assumed, however, that outcome values are stationary and do not change on-line over the course of a decision-making session. To see the limitations of such an assumption, imagine the rat is in a Skinner box and has started to earn outcomes (e.g., food pellets) by taking actions. One can assume that, as a result of consuming rewards, the motivation of the animal to earn more food outcomes will decrease (e.g., because of satiety) and, therefore, over time, the outcomes earned will have a lower value. Such changes in value should affect both optimal choice and vigor (Killeen, 1995) but have largely been ignored in previous models. Nevertheless, in most experimental and real-world scenarios, outcome values are affected by the history of outcome consumption, a phenomenon known as the “law of diminishing marginal utility”² in the economics literature, and as “drive reduction theory” in psychological accounts of motivation, which suppose that the drive for earning an outcome decreases as the consequence of its prior consumption (Hull, 1943; Keramati & Gutkin, 2014).

²Also known as “First Law of Gossen” named for Hermann Heinrich Gossen (1810 – 1858).

Here, building on previous work, we introduce the concept of a *reward field*, which captures non-stationary outcome values. Using this concept and methods from optimal control theory, we derive the optimal response vigor and choice strategy without assuming that outcome values are stationary. In particular, the results indicate that, even when the values of outcomes are changing, the optimal response rate in a free-operant procedure³ is a constant response rate. This finding rules out previous suggestions that the commonly observed decrease in within-session response rates is due to decreases in outcome value (Killeen, 1995). Instead, we show that decreases in within-session response rates can be explained by uncertainty regarding session duration. This later analysis is made possible by explicitly representing session duration in the current model, which is another dimension in which the current model extends previous work. The framework is then extended to choice situations and specific predictions are made concerning the conditions under which the optimal strategy involves maximization or probability matching.

Model Specification

The outcome space

We define the *outcome space* as a coordinate space with n dimensions, where n is the number of outcomes in the environment. For example, in a free-operant procedure in which the outcomes are water and food pellets, the outcome space will have two dimensions corresponding to water and food pellets. Within the outcome space, the state of the decision-maker at time t is defined by two factors: (i) the amount of *earned outcome* up to time t , which is denoted by \mathbf{x}_t and can be thought of as the position of the decision-maker in outcome space; e.g., in the above example, $\mathbf{x}_t = [1, 2]$ would indicate that one unit of water and two units of food pellet have been gained up to time t ; and (ii) the *outcome rate* at time t , denoted by \mathbf{v}_t , which can be considered the velocity of the decision-maker in the outcome space ($\mathbf{v}_t = d\mathbf{x}_t/dt$); e.g., if a rat is earning two units of water and one unit of food pellet per unit of time, then $\mathbf{v}_t = [2, 1]$. In general, we assume that the outcome rate cannot be negative ($\mathbf{v} \geq 0$), which means that the cumulative number of earned outcomes cannot decrease with time.

³In a free-operant procedure an animal is free to make responses continuously and repeatedly to earn outcomes.

The reward

We assume that there exists an n -dimensional *reward field*, denoted by $A_{\mathbf{x},t}$, where each element of $A_{\mathbf{x},t}$ represents the per unit value of each of the outcomes. For example, the element of $A_{\mathbf{x},t}$ corresponding to food pellets represents the value of one unit of food pellet at time t , given that \mathbf{x} units of outcome have been previously consumed. As such, $A_{\mathbf{x},t}$ is a function of both time and the amount of outcome earned. This represents the fact that (i) the reward value of an outcome can change value as a result of consuming previous outcomes, e.g., due to satiety (depending on \mathbf{x}) and (ii) the reward value of an outcome can change purely with the passage of time; e.g., an animal can get hungrier over time causing the reward value of food pellets to increase (depending on t).

In general, we assume that $A_{\mathbf{x},t}$ has two properties:

$$\frac{\partial A_{\mathbf{x},t}}{\partial \mathbf{x}} \leq \mathbf{0}, \frac{\partial A_{\mathbf{x},t}}{\partial t} \geq \mathbf{0}, \quad (1)$$

which entail that (i) the outcome values decrease (or remain constant) as more outcomes are earned, and (ii) that outcome values do not decrease with the passage of time. The latter assumption for example entails that, a rat experiences a higher amount of reward from consuming food pellets as it gets hungrier over time (even if no action is taken) due to the baseline metabolic rate at which the subject turns calories to energy.

Cost

Within the context of free-operant experiments, previous studies have expressed the cost of earning outcomes as a function of the delay between consecutive responses made to earn outcomes. For example, if a rat is required to make several lever presses to earn outcomes, then the cost will be higher if the delay between lever presses is short. More precisely, if the previous response has occurred τ time steps ago, then the cost of the current lever press will be:

$$C_\tau = \frac{a}{\tau} + b, \quad (2)$$

where a and b are constants (Dayan, 2012; Niv et al., 2007). b is the constant cost of each lever press, which is independent of the delay between lever presses whereas the factor a controls the

rate-dependent component of the cost. Previous research has established that predictions derived from this definition of cost are consistent with experimental data (Dayan, 2012; Niv et al., 2007). Note that costs such as basal metabolic rate and the cost of operating the brain, although consuming a high portion of energy produced by the body, are not included in the above definition because they are constant and independent of response rate and, therefore, are not directly related to the analysis of response vigor and choice.

Here, we express cost as a function of rate of earning outcomes rather than the rate of action execution⁴. We define the cost function K_v as the cost paid at each time step for earning outcomes at rate v . In the specific case that the outcome space has one dimension (there is only one outcome), and under ratio schedules of reinforcement (fixed-ratio, variable-ratio, random-ratio) in which the decision-maker is required to perform either precisely or on average k responses to earn one unit of outcome, the cost defined in equation 2 will be equivalent to:

$$K_v = ak^2v^2 + kbv. \quad (3)$$

See Appendix A for the proof. The cost is composed of two terms: a linear term (kbv), and a quadratic term (ak^2v^2). The linear term is coming from the constant cost of lever presses, i.e., for earning v amount of outcome, kv responses are required each at cost b (k is the average number of responses required for earning one unit of the outcome) and therefore the total cost will be kbv . The quadratic term comes from the rate-dependent component of the cost. That is, earning outcomes at rate v implies that kv responses were made at one unit of time, and therefore the delay between responses will be $1/kv$. The cost of each response is inversely proportional to the delay between responses, and therefore the cost of each response will be akv . Since kv responses are required to earn one unit of the outcome, the total cost will be $akv \times kv = ak^2v^2$, which is the quadratic term in equation 3. Such a quadratic form, independent of its connections to equation 2, is further motivated by the fact that quadratic forms are typically used to represent motor costs across optimal control studies (e.g., Berniker, O'Brien, Kording, & Ahmed, 2013; Salimpour & Shadmehr, 2014; Uno, Kawato, & Suzuki, 1989), which is partially due to the

⁴Note that the rate of earning outcomes is a function of the rate of action execution. For example, if k is the average number of responses required for earning one unit of the outcome, then the outcome rate is $1/k$ times the rate of action execution.

its simplicity while providing a reasonable approximation to more complex cost functions.

This definition of cost implies that it is only a function of outcome rate and is time-independent ($\partial K_v / \partial t = 0$). Although, in general, it may seem reasonable to assume that, as time passes within a session, the cost of taking actions will increase because of factors such as effector fatigue, here we made a time-independence assumption based on previous studies showing that factors such as effector fatigue have a negligible effect on response rate (McSweeney, Hinson, & Cannon, 1996). In general, we assume that at least one response is required to earn an outcome ($k > 0$), and the cost of earning outcomes is non-zero ($a > 0$).

Value

The reward earned in each time step can be calculated as the reward produced by one unit of each outcome ($A_{\mathbf{x},t}$) multiplied by the amount earned from each outcome, which will be $\mathbf{v} \cdot A_{\mathbf{x},t}$. The cost of earning this amount of reward is K_v , and therefore the net amount of reward earned (in dt time step) will be:

$$L_{\mathbf{x},\mathbf{v},t} = \mathbf{v} \cdot A_{\mathbf{x},t} - K_v. \quad (4)$$

A decision-making session starts at $t = 0$ and the total duration of that session is T . We denote the total reward gained in this period as $S_{0,T}$, which is the sum of the net rewards earned at each point in time:

$$S_{0,T} = \int_0^T L_{\mathbf{x},\mathbf{v},t} dt. \quad (5)$$

The quantity $S_{0,T}$ is called the *value* function, and the goal of the decision-maker is to find the optimal rate of earning outcomes that yields the highest value. The optimal rates that maximize $S_{0,T}$ can be found using different variational calculus methods, such as the Euler-Lagrange equation or the Hamilton-Jacobi-Bellman equation (Liberzon, 2011). The results presented in the next sections are derived using the Euler-Lagrange equation (see Appendix A for details of the value function in non-deterministic schedules).

Results

Optimal response vigor

In this section we use the model presented above to analyze optimal response vigor when there is one outcome and one response available in the environment. The analysis is divided into two sections. In the first section, we assume that the decision-maker is certain about session duration, i.e., that the session will continue for T time units, and we will extend this analysis in the next section to a condition in which the decision-maker assumes a probabilistic distribution of session lengths.

Response vigor when the session duration is known. We maintain the following theorem:

Theorem 1 *If the duration of the session is fixed and the time-dependent change in the reward field is zero ($\partial A_{x,t}/\partial t = 0$), then the optimal outcome rate is constant ($dv/dt = 0$). If the time-dependent change in the reward field is positive ($\partial A_{x,t}/\partial t > 0$), then the optimal outcome rate will be accelerating ($dv/dt > 0$).*

See Appendix B,C for a proof of this theorem. Note that the assumption $\partial A_{x,t}/\partial t = 0$ implies that the passage of time has no significant effect on the reward value of the outcome; e.g., a rat is not getting hungrier during an instrumental conditioning session⁵, which is a reasonable assumption given the short duration of such experiments (typically less than an hour). Within this condition, the above theorem states that the optimal response rate is constant throughout the session, even under conditions in which the reward value of the outcome decreases within the session as a result of earning outcomes, e.g., because of satiety. As an intuitive explanation for why a constant rate is optimal, imagine a decision-maker who chooses a non-constant outcome rate that results in a total of x_T outcomes during the session. If, instead of the non-constant rate, the decision-maker selects a constant rate $v = x_T/T$, then the total outcomes earned will be the same as before; however, the cost will be lower because it is a quadratic function of the outcome rate and, therefore, it is better to earn outcomes at a constant rate (Figure 1). Nevertheless, although this prediction is clear enough, it is not consistent with the experimental results,

⁵In an instrumental conditioning experiment an animal learns to perform specific actions on which the delivery of valued outcomes are contingent

described next.

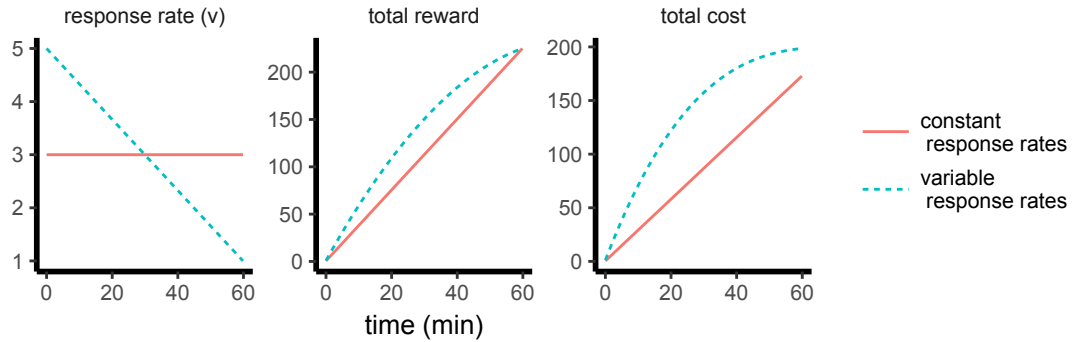


Figure 1. Total amount of reward and total cost paid during the session in two different conditions. **Left panel:** In the first condition (variable response rates), response rates are initially high at the beginning of the session, and then gradually decrease toward the end of the session. In the second condition, (constant response rates), response rates stay the same throughout the session. The unit of the y-axis is responses per minute. **Middle panel:** Total reward since the beginning of the session in each condition. In both conditions, the total amount of reward during the session is the same. The unit of the y-axis is arbitrary. **Right panel:** Total cost paid since the beginning of the session in each condition. The cost paid in the variable response rates condition is higher than the cost in the constant response rates condition, despite the fact that the amount of reward in both conditions at the end of the session is the same. The unit of the y-axis is arbitrary.

Within-session pattern of responses. It has been established that in various schedules of reinforcement, including variable-ratio (McSweeney, Roll, & Weatherly, 1994) and fixed-ratio (Bouton, Todd, Miles, León, & Epstein, 2013) schedules, the rate of responding within a session adopts a particular pattern: the response rate reaches its maximum a short time after the session starts (warm-up period), and then gradually decreases toward the end of the session (Figure 2:left panel). Killeen (1994) proposed a mathematical description of this phenomenon, which can be expressed as follows (Killeen & Sitomer, 2003):

$$\beta = \frac{r}{\delta r + 1/\alpha}, \quad (6)$$

where β is the response rate, δ is the minimum delay between responses, r is the resulting outcome rate, and α is called *specific activation*⁶. The model suggests that as the decision-maker earns outcomes during the session, the value of α gradually declines due to satiety, which will

⁶Note that in the original notation in Killeen and Sitomer (2003), α is denoted by a and β is denoted by b .

cause a decrease in response rate⁷. Although this model has been shown to provide a quantitative match to the experimental data, it is not consistent with Theorem 1 which posits that, even under conditions in which outcome values are changing within a session, the optimal response rate should not decrease during the session. As a consequence, the present model suggests that the cause of any decrease in the within-session response rate cannot be due purely to a change in outcome value.

Note, however, that the optimal response rate advocated by Theorem 1 is not consistent with reports of decreasing response rates across a session, which implies that some of the assumptions made to develop the model may not be accurate. Although the form of the cost and reward functions is reasonably general, we assumed that the duration of the session, T , is fixed and known by the decision-maker. In the next section we show that relaxing this assumption such that the duration of the session is unknown results in much closer concordance between predictions from the model and the experimental data.

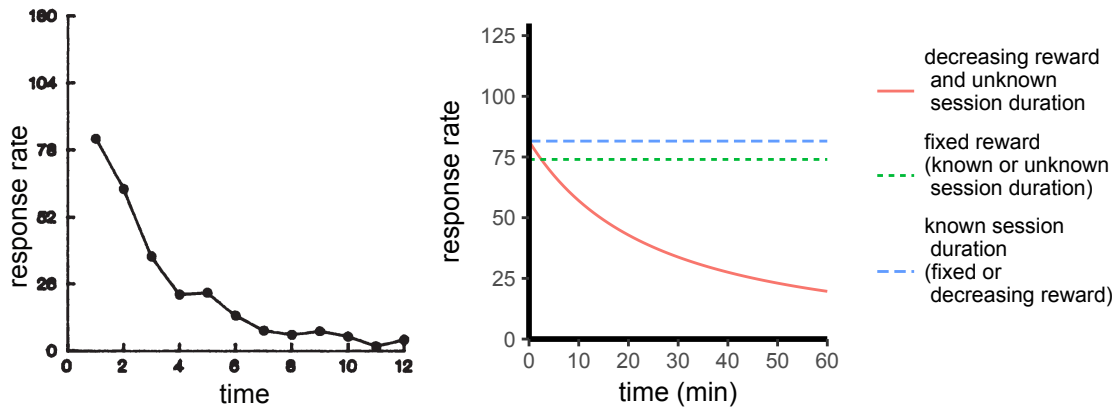


Figure 2. The pattern of within-session response rates (responses per minute). **Left panel:** Experimental data. The rate of responding per minute during successive intervals (each interval is 5 minutes) in a variable-ratio (VR) schedule ($k = 15$; VR15). The figure is adopted from McSweeney et al. (1994). **Right panel:** The theoretical pattern of within-session responses predicted by the model in different conditions. See text for details of each condition.

Response vigor when session duration is unknown. In this section we assume that the decision-maker is uncertain about the session duration, which can be either because the session duration is in fact non-deterministic, or because of the inherent inaccuracies in interval timing in animals (Gallistel & Gibbon, 2000; Gibbon, 1977). Since the session length is unknown, the

⁷Here satiety refers to both post-ingestive factors (such as blood glucose level; Killeen, 1995) and/or pre-ingestive factors (for example sensory specific satiety; McSweeney, 2004).

decision-maker assumes that the session can end at any point in time (T) with a probability distribution function $p(T)$. In this condition, a plausible way to calculate optimal response rates is to use $p(T)$ to set an expectation as to how long the session will last and to calculate the optimal response rate based on that expectation. Based on this, if t' time step has passed since the beginning of the session, then the expected session duration is $\mathbb{E}_{T \sim p(T)}[T|T > t']$ and therefore the value of the rest of the session will be $S_{t', \mathbb{E}[T|T > t']}$. The following theorem maintains that the optimal rate of outcome delivery that maximizes the value function is a decreasing function of the current time in the session t' , and therefore the optimal response rates will decrease throughout the session.

Theorem 2 *Assuming $S_{t', \mathbb{E}[T|T > t']}$ is the value function and that (i) the time dependent change in the reward field is zero ($\partial A_{x,t}/\partial t = 0$), (ii) the probability that the session ends at each point in time is non-zero ($p(T) > 0$), (iii) values of outcomes decrease as more are consumed ($\partial A_{x,t}/\partial x < 0$), then the optimal rate of outcome delivery is a decreasing function of t' :*

$$\frac{dv_{t'}^*}{dt'} < 0. \quad (7)$$

Furthermore, if conditions (i) and (ii) hold and the values of outcomes are constant ($\partial A_{x,t}/\partial x = 0$), then the optimal outcome rate is constant ($dv/dt = 0$).

See Appendix B,D for the proof of this theorem. Theorem 2 stems from two bases: (i) the optimal rate of outcome delivery is a decreasing function of session length, i.e., when the session duration is long the decision-maker can afford to earn outcomes more slowly, and (ii) when the session duration is unknown, expected session duration should increase with the passage of time (Figure 3). This phenomenon, which has been elaborated within the context of delayed gratification (McGuire & Kable, 2013; Rachlin, 2000), is more significant if the decision-maker assumes a heavy-tail distribution over T . Putting (i) and (ii) together implies that the optimal response rate will decrease with the passage of time. Importantly, this suggests, from a normative perspective, that uncertainty about the session duration *and* a decrease in the value of the outcomes are both necessary to explain within-session decreases in response rates.

For simulation of the model we characterized the session duration using a Generalized Pareto distribution following McGuire and Kable (2013). Simulations of the model are depicted

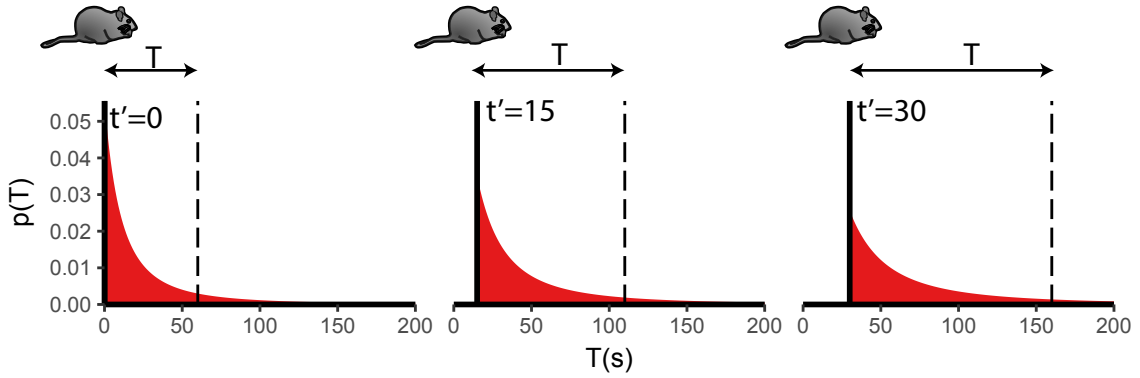


Figure 3. The expected length of the session changes as time passes within the session. The red areas in the panels show the probability distribution function of the length of the session ($p(T)$). The vertical dashed-lines represent the expected length of the session ($\mathbb{E}_{T \sim p(T)}[T | T > t']$), and the vertical solid-lines represent the current time in the session (t'). **Left panel:** at the beginning of the session ($t' = 0$), the animal expects the session to last for 60 minutes. After 15 minutes have passed since the beginning of the session (**middle panel**; $t' = 15$), the expected duration of the session becomes 110 minutes. As more time passes (**right panel**; $t' = 30$), the expected duration of the session increases to 160 minutes. The unit of the x-axes in the panels is minutes.

in Figure 2:right panel. The simulations were obtained using analytical equations derived from Theorem 2 and trial-by-trial updates of the expected session length (see Appendix D.2 for details). Simulations show three different conditions. In condition (i) the session duration is known and, as the figure shows, irrespective of whether the value of outcomes is decreasing or fixed, the optimal response rate is constant. In condition (ii) session duration is unknown, but the value of outcomes is constant. Again in this condition the optimal response rate is constant. In condition (iii) session duration is unknown *and* the reward value decreases as more outcomes are consumed. Only in this condition, consistent with experimental data and Theorem 2, response rates decrease as time passes. Therefore, the simulations confirm that a decrease in outcome value alone is not sufficient to explain within-session response rates and that uncertainty about session duration is also required to reproduce a pattern of responses consistent with the experimental data. Note that a similar pattern can also be obtained using any other distribution that assigns a non-zero probability to positive values of T .

Relationship to temporal discounting. There are, however, alternative explanations available based on changes in outcome value. One candidate explanation is based on the temporal discounting of outcome value according to which the value of future rewards is discounted

compared to more immediate rewards. Typically, the discount value due to delay is assumed to be a function of the interaction of delay and outcome value. If, at the beginning of the session, outcome values are large (e.g., because a rat is hungrier), then any discount produced by selecting a slow response rate will be larger at that point than later in the session when the value of the outcome is reduced (e.g., due to satiety) and so a delay will have less impact. It could be argued, therefore, that it is less punitive to maintain a high response rate at the beginning of the session to avoid delaying outcomes and then to decrease response rate as time passes within the session. As such, temporal discounting predicts decreases in within-session response rates consistent both with experimental observations and with the argument that outcome value decreases within the session (e.g., the satiety effect).

Prediction. Although plausible, such explanations make very different predictions compared to the model. The most direct prediction from the model is that introducing uncertainty into the session duration without altering the average duration should nevertheless lead to a sharper decline in response rate within the session; e.g., if for one group of subjects the session lasts exactly 30 minutes whereas for another group the session length is uncertain and can end at any time (but ends on average after 30 minutes), then the model predicts that the response rate in the second group will be higher at the start and decrease more sharply than in the first group. This effect is not anticipated by the temporal discounting account of the effect.

Another prediction of the model is with regard to the effect of training on within-session response rates. By experiencing more training sessions, subjects should be able to build a more accurate representation of the session length. This implies that, for this case, the expected length of the session will remain relatively unchanged as time passes within a session and, therefore, the decrement in within-session response rates should be predicted to grow smaller with more training. Consistent with this prediction, some experimental results indicate that the gap between the highest and the lowest response rates within a session does decrease with more training (McSweeney & Hinson, 1992, Figure 11)⁸ while other studies show that the gap becomes larger as training proceeds. It is worth noting that in the former study the shaping sessions were excluded

⁸Note that in these experiments, animals were trained on a variable-interval schedule. In a variable-interval schedule of reinforcement, it is the time period since the last outcome delivery that determines whether the next response will be rewarded, which is in contrast to ratio schedules where outcome delivery depends on the number of responses. The current model and theorems apply to ratio schedules, and therefore, this prediction can be tested more accurately using a ratio procedure.

when comparing early and late training sessions, while in the latter study they were not. Based on this, further analysis and experimental studies are required to test this prediction accurately.

Effect of experimental parameters. Optimal response rates predicted by the model are affected by experimental parameters (e.g., reward magnitude), which can be compared against experimental data. In general, in an instrumental conditioning experiment, the duration of the session can be divided into three sections: (i) outcome handling/consumption time, which refers to the time that an animal spends consuming the outcome, (ii) *post-reinforcer pause*, which refers to the pause that occurs after consuming the outcome and before starting to make the next response (e.g., lever press), something consistently reported in studies using a fixed-ratio schedule, and (iii) *run time*, which refers to the time spent making responses (e.g., lever pressing). Experimental manipulations have been shown to have different effects on the duration of these three sections of the session (see below), and decisions about whether each of these sections is included when calculating response rates can affect the results. The predictions of the current model are with regard to response rate; whether manipulating experimental parameters should be expected to change the two other measures (outcome handling and post-reinforcer pause) cannot be predicted by the current model. In the following sections, we briefly review the currently available data from instrumental conditioning experiments and their relationship to predictions of the model. Simulations are obtained using analytical equations derived in Theorem 1 (see Appendix D.3 for details)⁹.

The effect of response cost (a and b). Experimental studies in rats working on a fixed-ratio schedule (Alling & Poling, 1995) indicate that increasing the force required to make responses causes increases in both inter-response time and the post-reinforcement pause. The same trend has been reported in Squirrel monkeys (Adair & Wright, 1976). Consistent with this observation the present model predicts that increases in the cost of responding, for example by increasing the effort required to press the lever (increases in a and/or b), lead to a lower rate of earned outcomes and a lower rate of responding (Figure 4). The reason for this is that, by increasing the cost, the response rate for each outcome should slow in order to compensate for the increase in the cost and so maintain a reasonable gap between the reward and the cost of each outcome.

⁹Note that, for simplicity, the simulations in this section are made under the assumption that the session duration is fixed.

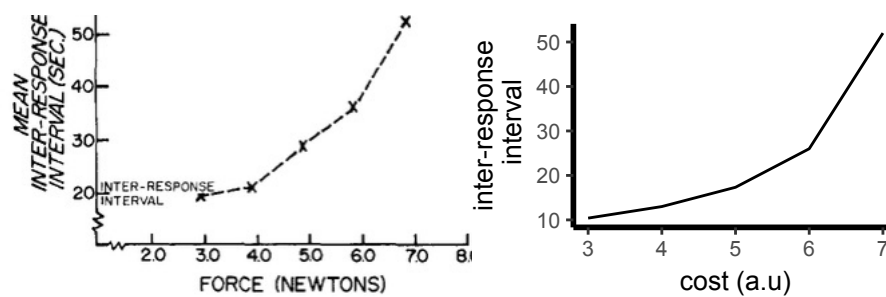


Figure 4. Effect of response cost on response rates. **Left panel:** Empirical data. Inter-response intervals (in seconds) when the force required to make a response is manipulated. Figure is adopted from Adair and Wright (1976). **Right panel:** Model prediction. Inter-response interval (in seconds; equal to the inverse of response rates) as a function of cost of responses (b).

The effect of ratio-requirement (k). Experimental studies mainly suggest that the rate of earned outcomes decreases with increases in the ratio-requirement (Aberman & Salamone, 1999; Barofsky & Hurwitz, 1968), which is consistent with the general trend in the optimal rate of outcome delivery implied by the present model (see below).

Experimental studies on the rate of responding on fixed-ratio schedules indicate that the post-reinforcement pause increases with increases in the ratio-requirement (Ferster & Skinner, 1957, Figure 23)(Felton & Lyon, 1966; Powell, 1968; Premack, Schaeffer, & Hundt, 1964). In terms of overall response rates, some experiments report that response rates increase with increases in the ratio-requirement up to a point beyond which response rates will start to decline, in rats (Barofsky & Hurwitz, 1968; Kelsey & Allison, 1976; Mazur, 1982), pigeons (Baum, 1993) and mice (Greenwood, Quartermain, Johnson, Cruce, & Hirsch, 1974), although other studies have reported inconsistent results in pigeons (Powell, 1968), or a decreasing trend in response rate with increases in the ratio-requirement (Felton & Lyon, 1966; Foster, Blackman, & Temple, 1997). The inconsistency is partly due to the way in which response rates are calculated in the different studies; for example in some studies outcome handling and consumption time are not excluded when calculating response rates (Barofsky & Hurwitz, 1968), in contrast to other studies (Foster et al., 1997). As a consequence, the experimental data regarding the relationship between response rate and the ratio-requirement is inconclusive.

With regard to this issue, the present model predicts that the relationship between response rate and the ratio-requirement is an inverted U-shaped pattern (Figure 5: left panel), which is consistent with the studies mentioned previously, depending on the value of other ex-

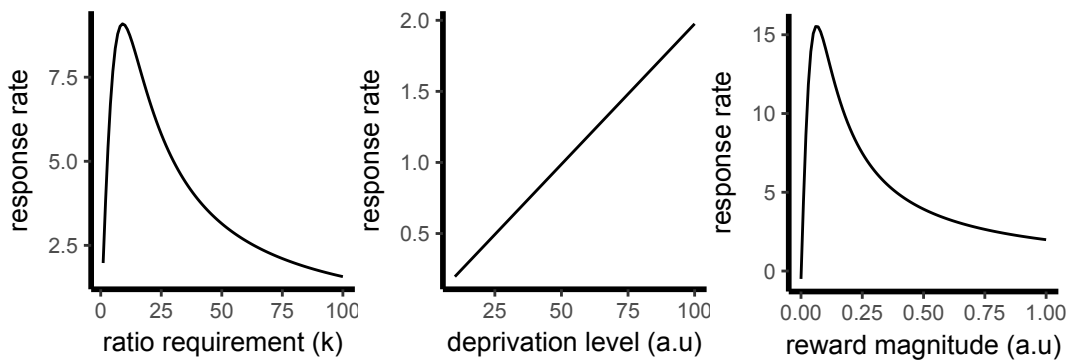


Figure 5. Left panel: The effect of ratio-requirement on the response rate (responses per minute). **Middle panel:** The effect of deprivation level on response rates. **Right panel:** The effect of the reward magnitude on response rates.

perimental parameters. The model makes an inverted U-shaped prediction because, under a low ratio-requirement, the cost is generally low and, therefore, as the ratio-requirement increases, the decision-maker will make more responses to compensate for the drop in the amount of reward. In contrast, when the ratio-requirement is high, the cost of earning outcomes is high and the margin between the cost and the reward of each outcome becomes significantly tighter as the ratio-requirement increases. The margin can, however, be loosened by decreasing the response rate (see Appendix D.2 for the exact source of this effect in the model).

The Effect of deprivation level. Experimental studies that have used fixed-ratio schedules suggest that response rates increase with increases in deprivation (Chapter 4, Ferster & Skinner, 1957; Sidman & Stebbins, 1954). However, such increases are mainly due to decreases in the post-reinforcement pause, and not due to the increases in the actual rate of responding after the pause (see Pear, 2001, Page 289 for a review of other reinforcer schedules; see for example Eldar, Morris, & Niv, 2011 for the case of variable-interval schedules). The model predicts that, with increases in deprivation, the rate of responding and the rate of earned outcomes will increase linearly (Figure 5: middle panel). The rate of increase is, however, negligible when the outcomes are small and the generated satiety after earning each outcome is insignificant. This provides a potential reason why the effect of deprivation on response rate has not previously been observed in experimental data. Similarly, when the session duration is long, even under high deprivation, sufficient time is available to earn enough reward and become satiated, and therefore the effect of deprivation levels on response rate will be minor.

The effect of reward magnitude. Some studies show that post-reinforcement pauses increase as the magnitude of the reward increases (Powell, 1969), whereas other studies suggest that the post-reinforcement pause decreases (Lowe, Davey, & Harzem, 1974); however, in this later study the magnitude of reward was manipulated within-session and a follow-up study found that, at a steady state, the post-reinforcement pause increases with increases in the magnitude of the reward (Meunier & Starratt, 1979). Reward magnitude does not, however, have a reliable effect on overall response rate (Keesey & Kling, 1961; Lowe et al., 1974; Powell, 1969). Regarding predictions from the model, the effect of reward magnitude on earned outcome and response rates is, again, predicted to take an inverted U-shaped relationship (Figure 5: right panel), and, therefore, depending on the value of the parameters, the predictions of the model are consistent with the experimental data. The model makes a U-shaped prediction because, when the reward magnitude is large then, given high response rates, the animals will become satiated quickly and, therefore, the reward value of future outcomes will decrease substantially if the animal maintains this high response rate. As a consequence, under a high reward magnitude condition, an increase in reward will cause response rates to decrease. Under a low reward magnitude condition, however, satiety has a negligible effect and a high response rate ensures that sufficient reward can be earned before the session ends.

Summary. Table 1 shows the summary of the predictions of the model presented here and also the predictions of the model in Niv (2007) with regard to the effect of experimental parameters. The predictions of the models are different with respect to the effect of reward magnitude on response rates. The previous work predicts that higher reward magnitudes lead to higher response rates, whereas the study here predicts an inverted U-shaped relationship between them, i.e., further increases in reward magnitude when it is already high, will lead to lower response rates. The reason is that according to the current study, high reward magnitudes cause satiety and thus diminish outcome values, which can support lower response rates. This effect of satiety (within a session) is not explicitly modelled in the previous work and thus the predictions of the two frameworks differ.

Experimental parameters	Current model	Niv (2007)
increase in response cost	lower response rates (Figure 4: right panel)	lower response rates (page 59; Niv, 2007)
increase in ratio-requirement	inverted U-shaped (Figure 5: left panel)	lower response rates or inverted U-shaped (Figure 2.10a; Niv, 2007)
increase in deprivation levels*	higher response rates (Figure 5: middle panel)	higher response rates (Figure 2.10d; Niv, 2007)
increase in reward magnitude	inverted U-shaped (Figure 5: right panel)	higher response rates (Figure 2.10d; Niv, 2007)

Table 1

Summary of the predictions of the current model with regard to the experimental parameters. The table also presents the predictions of Niv (2007).

**Increases in the deprivation levels are assumed to increase the reward magnitude.*

Optimal choice and response vigor

In this section we address the choice problem, i.e., the case where there are multiple outcomes available in the environment and the decision-maker needs to make a decision about the response rate along each outcome dimension. An example of this situation is a concurrent free-operant procedure in which two levers are available and pressing each lever produces an outcome on a ratio schedule. Unlike the case with single outcome environments, the optimal rate of earning outcomes is not necessarily constant and can take different forms depending on whether the reward field is a *conservative field* or a *non-conservative field*, and whether the costs of responses along the outcome dimensions are independent of each other. Below, we derive the optimal choice strategy in each condition.

Conservative reward field. A reward field is conservative if the amount of reward experienced by consuming different outcomes does not depend on the *order* of consumption and depends only on the number of each outcome earned by the end of the session. This property holds in two conditions (i) when the value of each outcome is unrelated to the prior consumption of other outcomes; and (ii) the consumption of an outcome affects the value of other outcomes and this effect is symmetrical. As an example of condition (i), imagine an environment with two outcomes in which one of the outcomes only satisfies thirst and the other only satisfies hunger¹⁰. Here, consumption of one of the outcomes will not affect the amount of reward that will be ex-

¹⁰In this example we assumed that hunger and thirst are independent motivational drives.

perienced by consuming the other outcome and, therefore, the total reward during the session does not depend on the order of choosing the outcomes. As an example of condition (ii), imagine an environment with two outcomes in which both outcomes satisfy hunger and, therefore, consuming one of the outcomes reduces the amount of future reward produced by the other outcome. Here, if the effect of the outcomes on each other is symmetrical, i.e., consuming outcome O_1 reduces the reward elicited by outcome O_2 by the same amount that consuming outcome O_2 reduces the reward elicited by outcome O_1 , then it will not matter which outcome is consumed first and the total reward during the session will be independent of the chosen order. As such, the reward field will be conservative.

Under the conditions that a reward field is conservative, the optimal response rate will be constant for each outcome relative to the other. Intuitively, this is because, in terms of the total rewards per session, the only thing that matters is the final number of earned outcomes and, therefore, there is no reason why the relative allocation of responses to outcomes should vary within the session. The actual response rate for each outcome will, however, depend on whether the costs of the outcomes are independent, a point elaborated in the next section.

Conservative reward field and independent response cost. The costs of various outcomes are independent if the decision-maker can increase their work for one outcome without affecting the cost of other outcomes. As an example, imagine a decision-maker that is using their left hand to make responses that earn one outcome and their right-hand to make responses that earn a second outcome. In this case, the independence assumption entails that the cost of responding with one or other hand is determined by the response rate on that hand; e.g., the decision-maker can increase or decrease rate of responding on the left hand without affecting the cost of responses on the right hand. More precisely, the independence assumption entails that the Hessian matrix of K_v is diagonal:

$$\frac{\partial^2 K_v}{\partial v_i \partial v_j} = 0, i \neq j. \quad (8)$$

In this situation even if some of the outcomes have a lower reward or a higher cost (inferior outcomes) compared to other outcomes (superior outcomes), it is still optimal to allocate a portion of responses to the inferior outcomes. This is because responding for inferior outcomes has no effect on the net reward earned from superior outcomes and, therefore, as long as the

response rate for inferior outcomes is sufficiently low that the reward earned from them is more than the cost paid, responding for them is justified. The portion of responses allocated to each outcome depends, however, on the cost and reward of each outcome. We maintain the following theorem:

Theorem 3 *If (i) the reward field is conservative, (ii) the time-dependent term of the reward field is zero ($\partial A_{\mathbf{x},t}/\partial t = \mathbf{0}$), and (iii) the cost function satisfies equation 8, then the optimal rate of earning outcome \mathbf{v}^* will be constant ($d\mathbf{v}/dt = \mathbf{0}$) and satisfies the following equation:*

$$\frac{\partial K_{\mathbf{v}^*}}{\partial \mathbf{v}^*} = A_{T\mathbf{v}^*,T}. \quad (9)$$

See Appendix E,E.1 for the proof and for the specification of optimal responses. As an example, imagine a concurrent fixed-ratio schedule in which a subject is required to make k responses with the left hand to earn O_1 , and lk responses with the right hand to earn O_2 , and both outcomes have the same reward properties. According to Theorem 3, the optimal response rate for O_1 (the outcome with the lower ratio-requirement) will be l times greater than the response rate for the second outcome O_2 . Figure 6:left panel independent cost condition shows the simulation of the model and the optimal trajectory in the outcome space. As the figure shows, the rate of earning O_1 is higher than O_2 , however, the proportion of each outcome of the total remains the same throughout the session.

Relationship to probability matching. The above results are generally in line with the *probability matching* notion, which states that a decision-maker allocates responses to outcomes based on the ratio of responses required for each outcome (Bitterman, 1965; Estes, 1950). Probability matching is often argued to violate rational choice theory because, within this theory, it is expected that a decision-maker works exclusively for the outcome with the higher probability (i.e., the lower ratio-requirement). However, based on the model proposed here, probability matching is the optimal strategy when the cost of actions is independent, and therefore consistent with rational decision-making.

Relationship to matching law. The *matching law* refers to the observation that the rate of responses for different actions is proportional to the rate of rewards obtained from the corresponding actions (Herrnstein, 1961). For example, if ν_1 and ν_2 are the response rates for two

different actions, and z_1 and z_2 refer to the rate of rewards obtained from each action, then the matching law implies that,

$$\frac{v_1}{v_2} = \frac{z_1}{z_2}. \quad (10)$$

In contrast to the matching law which is about rewards obtained from each action, in probability matching the responses are allocated to actions according to the probability of rewards being *available* for each action. In this respect these two behavioral phenomena are different. For example, although maximization (exclusively selecting the action with the higher reward probability/lower ratio-requirement) is inconsistent with probability matching, it is indeed consistent with the matching law (because in maximization 100% of responses are made on one of the actions and 100% of rewards are obtained from that action). The results that we obtained in the previous sections are related to the rate at which outcomes are available on each action, and therefore, they are not directly related to the matching law. Furthermore, the matching law mostly applies to the case of variable-interval schedules¹¹, and is not particularly informative in the case of ratio schedules, which are the focus of the current analysis. This is because in ratio schedules, the rate of earning rewards from actions is directly related to the rate of responding on the corresponding actions no matter how the decision-maker distributes responses over actions.

The relationship to Kubanek (2017). Typically, in computational models of the matching law and probability matching, the effect of effort, i.e., the cost of obtaining rewards, is not explicitly modelled (e.g., Iigaya & Fusi, 2013; Loewenstein, Prelec, & Seung, 2009; Sakai & Fukai, 2008). An exception can be found in the study of Kubanek (2017) in which the matching law is regarded as a consequence of the diminishing returns associated with variable-interval schedules of reinforcement. In such schedules, outcome rate grows almost proportional to response rate when response rates are low, whereas outcome rate saturates when response rates are high (because in these schedules a certain period of time has to pass before the next outcome can be earned) and, therefore, to produce a slight increase in outcome rate will require a significant increase in response rate. Based on this, outcomes are more expensive to earn at high response rates, which justifies allocating a portion of responses to inferior actions, on which the outcomes are not yet saturated and are still (relatively) cheap. As such, in variable-interval schedules we

¹¹In variable-interval schedules, the subject needs to wait a certain amount of time (according to a probability distribution) before being able to obtain the next reward by selecting actions.

would expect animals to match rather than respond exclusively on the superior action, and indeed, Kubanek (2017) showed that the matching law is the optimal strategy when faced with these schedules.

This prediction for variable-interval schedules is essentially the same as the prediction generated in the current study for ratio schedules (and independent response costs) even though, unlike variable interval schedules, the outcome rates are non-saturating. This is because, although on ratio schedules outcome rates are non-saturating and proportional to response rates, the cost of earning outcomes increases as response rates increase due to the quadratic cost of responses (as implied in equation 3), meaning that it is better to limit response rates even on superior actions. As such, although the model proposed here is focused on ratio schedules and the one in Kubanek (2017) on variable-interval schedules, both approach optimal decisions based on the fact that the outcomes are more expensive when response rates are high; and whereas in the former it is due to the quadratic cost function, in the latter it is due to the properties of interval schedules, and in this respect the two studies are complementary. In addition, the model proposed here extends previous work by addressing the role of changes in outcome value on choice, in addition to the role that the cost of earning outcomes plays.

Conservative reward field and dependent response cost. In this section we assume that the cost of responses for an outcome is affected by the rate of responding required to earn other outcomes. In other words, what determines the cost is the delay between subsequent responses either for the same or for a different outcome; i.e., the cost is proportional to the rate of earning all of the outcomes. In concurrent free-operant procedures this assumption entails, for example, that the cost of pressing, say, the right lever is determined by the time that has passed since the last press on either the right or a left lever. In this condition the optimal strategy is *maximization*; i.e., to take the action with the higher reward (or lower ratio-requirement) and to stop taking the other action (see Appendix G). The reason is because, unlike the case with independent costs, allocating more responses to earn an inferior outcome will increase the cost of earning superior outcomes and, therefore, it is better to pay the cost for the superior outcome only, which requires fewer responses per unit of outcome.

For example, under a concurrent fixed-ratio schedule in which an animal needs to make k responses on the left lever to earn O_1 , and lk responses on the right lever to earn O_2 (O_1 and O_2

have the same reward properties), the optimal response rate will be a constant response rate on the left lever and a zero response rate on the right lever. Figure 6:left panel dependent cost condition shows a simulation of the model and the optimal trajectory in outcome space, which shows that the subject only earns O_1 . Note the difference between this example, and the example mentioned in the previous section is that, here the costs of earning outcomes are not independent, while in the previous section we assumed that the costs of earning O_1 and O_2 are independent of each other.

Prediction. One way of testing the above explanation for maximization and matching strategies is to compare the pattern of responses when two different effectors are used to make responses for the outcomes vs. when a single effector is being used to earn both outcomes. In the first condition the costs of the two outcomes are independent of each other whereas in the second condition they are dependent on each other. As a consequence, the theory predicts that, in the first condition, response rates will reflect probability matching whereas in the second condition they will reflect the maximization strategy.

Probability matching and maximization. As such, whether the outcome rate reflects a probability matching or a maximization strategy depends on the cost function and, in concurrent free-operant procedures, the cost that reflects the maximization strategy is more readily applicable. Regarding the experimental data, evidence from concurrent instrumental conditioning experiments in pigeons tested using variable-ratio schedules (Herrnstein & Loveland, 1975) is in-line with the maximization strategy and consistent with predictions from the model.

Within the wider scope of decision-making tasks, some studies are consistent with probability matching, whereas other studies provide evidence in-line with a maximization strategy (see Vulkan, 2000 for a review). However, many of these latter studies use discrete-trial tasks in which, unlike free-operant tasks (which are the focus of the current analysis), actions are typically disjoint and therefore unlikely to convey a rate-dependent cost. Even within the domain of free-operant tasks, for the cost of actions to be independent of each other the decision-maker should be able to respond using effectors independent of each other (e.g., left hand and right hand), otherwise, as argued, probability matching will no longer be the optimal strategy. In spite of this, some evidence suggests that probability matching occurs even in settings where the task is discrete-trial or when responses are not independent. In these settings, observed probability

matching will be unrelated to the current analysis and might stem from other factors such as cognitive efforts and limitations (e.g., Schulze & Newell, 2016), tendency of the subjects to find patterns in random sequences (e.g., Gaissmaier & Schooler, 2008), or it could be the effect of competition in certain environments (Schulze, van Ravenzwaaij, & Newell, 2015).

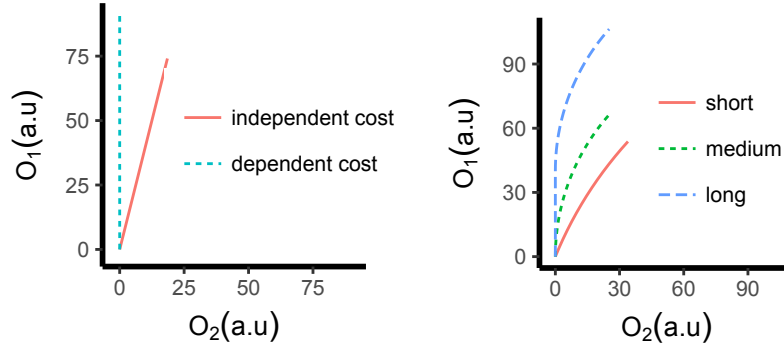


Figure 6. Left panel: Optimal trajectory in a conservative reward field. Earning O_1 requires k responses and earning O_2 requires lk responses. Initially the amount of earned outcome is zero (starting point is at point $[0,0]$), and the graph shows the trajectories that the decision-maker takes in two different conditions corresponding to when the costs of outcomes are independent, and when the costs are dependent on each other. **Right panel:** The optimal trajectories in the outcome space when the reward field is non-conservative. The graph shows the optimal trajectory in the conditions that the session duration is short ($T = 7$), medium ($T = 15.75$) and long ($T = 23$). O_1 generates more reward than O_2 .

Non-conservative reward field. A reward field is non-conservative if the total amount of reward experienced during the session depends on the order of the consumption of the outcomes. Imagine an environment with two outcomes say O_1 and O_2 , where both outcomes have the same motivational properties, e.g., consumption of one unit of either O_1 or O_2 decreases hunger by one unit, however, they generate different amounts of rewards, e.g., one unit of O_1 generates more reward than one unit of O_2 . As an example, let's denote the amount of earned O_1 and O_2 by x_1 and x_2 respectively. Based on this, the current food deprivation level will be $H - x_1 - x_2$, where H is the initial deprivation level. Here, although both outcomes have the same effect on reducing the deprivation level, in a non-conservative reward field, one of the outcomes (O_1 in this example) generates more reward than the other:

$$A_{\mathbf{x},t} = \left[\underbrace{l(H - x_1 - x_2)}_{O_1}, \underbrace{H - x_1 - x_2}_{O_2} \right], \quad (11)$$

which implies that the reward generated by both outcomes is proportional to the current food deprivation level, *and* the reward of O_1 is l times greater than the reward generated by O_2 . Within such an environment, the total amount of reward experienced depends on the order of consuming outcomes. This is because if hunger is high then consuming O_1 generates significantly more reward than O_2 and, therefore, early in the session it is better to allocate more responses to O_1 ; whereas later in the session when hunger is presumably lower and the difference in the value of the outcomes is small, responses for O_2 can gradually increase. If we reverse this order, i.e., first O_2 is consumed and then O_1 , then early consumption of O_2 will cause satiety and the decision-maker will lose the chance to experience high reward from O_1 when hungry. As such, the amount of experienced reward depends on the order of consuming the outcomes and, based on the above explanation, a larger amount of reward can be earned during the session if more responses are allocated to the outcome with the higher reward at the beginning of the session (see Appendix H). Figure 6:right panel shows the simulations of the model under different session durations (simulations are obtained using analytical solutions). In each simulation, at the beginning of the session the initially earned outcomes are zero and each line in the figure shows the trajectory of the amount earned from each outcome during the session. As the figure shows, in all conditions the rate of earning O_1 is higher than O_2 and this difference is more apparent under long session durations, in which a large amount of reward can be gained during the session and it makes a significant difference to earn O_1 first.

Prediction. A test of the above prediction would involve an experiment in which the subject is responding for two food outcomes containing an equal number of calories (and therefore having the same impact on motivation) but with different levels of the desirability (e.g., having different flavors) and, therefore, having a different reward effect. Theorem A3 predicts that, if the session duration is long enough, early in the session the response rate for the outcome with the greater desirability will be higher whereas, later in the session, responses for the other outcome will increase.

Relationship to motivational drives. Formally, a reward field is conservative if there exists a scalar field, denoted by $D_{\mathbf{x}}$, such that:

$$A_{\mathbf{x},t} = -\frac{\partial D_{\mathbf{x}}}{\partial \mathbf{x}}. \quad (12)$$

Keramati and Gutkin (2014) conceptualized D_x as the motivational drive for different outcomes and provided a definition of motivational drives as deviations of the *internal state* of a decision-maker from their homeostatic set-points. Based on this definition, according to equation 12, rewards are generated as a consequence of reductions in drive and, more precisely, the reward field is the amount of change in the motivational drive that is due to the consumption of one unit of each outcome. It can be shown that if a reward field satisfies equation 12 then the amount of reward experienced in a session depends on the total number of earned outcomes and, therefore, it is conservative. For the case of non-conservative reward fields, the drive for earning an outcome not only depends on the number of earned outcomes, but also on the order in which they were earned. However, D_x only depends on the number of earned outcomes (dependency on x) and not on their order, and because of this it cannot be defined in non-conservative reward fields. In this respect, the current study extends the model proposed by Keramati and Gutkin (2014) to cases where rewards do not correspond to any underlying motivational drive.

Conclusion and Discussion

Computational models of action selection are essential for understanding decision-making processes in humans and animals, and here we extended these models by providing a general analytical solution to the problem of response vigor and choice. Table 2 shows the summary of the results obtained for different conditions. The results provide (i) a normative basis for commonly observed decrements in within-session response rates, and (ii) a normative explanation for probability matching and reward maximization, as two commonly observed choice strategies.

Relationship to previous models of response vigor. There are two significant differences between the model proposed here and previous models of response vigor (Dayan, 2012; Niv et al., 2007). Firstly, although the effect of between-session changes in outcome values on response vigor was addressed in previous models (Niv, Joel, & Dayan, 2006), the effects of on-line changes in outcome values within a session were not addressed. On the other hand, the effect of changes in outcome value on the choice between actions has been addressed in some previous models (Keramati & Gutkin, 2014), however their role in determining response vigor has not been investigated. We address such limitations directly in this model.

	n	T	reward field	cost function	response rates	thrm	pre. works
1	n=1	known	constant	-	constant	1	constant
2	n=1	known	non-constant	-	constant/increasing	1	-
3	n=1	unknown	constant	-	constant	2	constant
4	n=1	unknown	non-constant	-	decreasing	2	-
5	n>1	known	conservative	independent	prob. matching	3	-/prob. matching
6	n>1	known	conservative	dependent	maximization	A2	-/maximization
7	n>1	known	non-conserv	independent	see text	A3	-

Table 2

Summary of the results. 'n' refers to the number of dimensions of the outcome space and 'response rates' refers to the optimal response rates obtained by the corresponding theorem. 'constant' reward field implies that the values of the outcomes do not change as more outcomes are consumed. 'non-constant' reward field implies that the values of the outcomes decrease as more outcomes are consumed. 'known' session length implies that the decision-maker is certain about the session length (T). 'unknown' session length implies that the decision-maker is uncertain about the session length. 'non-conserv' refers to non-conservative reward field. 'prob. matching' refers to probability matching. 'pre. works' refers to 'previous works'. 'thrm' refers to the corresponding theory in each case.

Secondly, previous work conceptualized the structure of the task as a semi-Markov decision process in which taking an action leads to outcomes after a delay. Based on that, the optimal actions and the delay between them that maximize the average reward per unit of time (average reward) were derived. Here, we used a variational analysis to calculate the optimal actions that maximize the reward earned within the session. One benefit of the approach taken in the previous works is that it extends naturally to a wide range of instrumental conditioning schedules such as interval schedules, whereas the extension of the model proposed here to the case of interval schedules is not trivial. Optimizing the average reward (as adopted in previous work) is equivalent to the maximization of reward in an infinite-horizon time scale; i.e., the session duration is unlimited. In contrast, the model used here explicitly represents the duration of the session which, as we showed, plays an important role in the pattern of responses.

In addition to the predictions of the current model, Table 2 shows the predictions of previous models of response vigor in each condition. The cases that involve non-constant reward fields are not addressed in previous work and, therefore, their predictions are not mentioned in the table. In the case of environments in which one outcome type is available ($n = 1$), and the reward field is constant, the prediction of the previous works is that the response rates will be constant, which is the same as the prediction of the current model (Table 2 rows #1). In the case

of environments with more than one outcome type ($n = 2$), and constant reward fields, we expect the prediction from previous research to be optimal in both ‘dependent’ and ‘independent’ cost conditions (Table 2 row #6,#7). This is because, in these conditions, the optimal response rates are constant within a session, and therefore the previous models should be able to learn them, in which case their predictions will be the same as the predictions from the current model.

Relationship to principle of least action. A basic assumption that we made here is that the decision-maker takes actions that yield the highest amount of reward. This reward maximization assumption has a parallel in the physics literature known as *the principle of least action*, which implies that among all trajectories that a system can take, the true trajectories are the ones that minimize the action. Here action has a different meaning from that used in the psychology literature, and it refers to the integral of the Lagrangian (L) along the trajectory. In the case of a charged particle with charge q and mass m in a magnetic field B , the Lagrangian will be:

$$L = \frac{1}{2} m \mathbf{v}^2 + q \mathbf{v} \cdot \mathbf{A}, \quad (13)$$

where \mathbf{A} is the vector potential ($B = \nabla \times \mathbf{A}$). By comparing equation 13 with equations 4 and 5, we can see that the reward field $A_{\mathbf{x},t}$ corresponds to the vector potential, and the term K_v corresponds to $\frac{1}{2} m \mathbf{v}^2$ by assuming $m = 2ak^2$, and $b = 0$. This parallel can provide some insights into the properties of the response rates. For example, it can be shown that when the Lagrangian is not explicitly dependent on time (time-translational invariance), which here implies that $\partial A_{\mathbf{x},t} / \partial t = 0$, then the Hamiltonian (\mathcal{H} , or energy) of the system is conserved. The Hamiltonian in the case of the system defined in equation 4 (with single outcome) is:

$$\begin{aligned} \mathcal{H} &= K_v - \frac{\partial K_v}{\partial v} v \\ &= -ak^2 v^2 \quad (\text{using equation 3}). \end{aligned}$$

Conservation of the Hamiltonian implies that $ak^2 v^2$ (and therefore v) is constant (response rate is constant), as stated by Theorem 1. Further exploration of this parallel can be an interesting future direction.

Disclosures and Acknowledgments

A.D was support by the CSIRO and by grants DP150104878, FL0992409 from the Australian Research Council. B.W.B was supported by a Senior Principal Research Fellowship from the National Health & Medical Research Council of Australia, GNT1079561. The author declares that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest. We are grateful to Hadi Lookzadeh and Peter Dayan for helpful discussions.

References

- Aberman, J. E., & Salamone, J. D. (1999). Nucleus accumbens dopamine depletions make rats more sensitive to high ratio requirements but do not impair primary food reinforcement. *Neuroscience*, 92(2), 545–52.
- Adair, E. R., & Wright, B. A. (1976). Behavioral thermoregulation in the squirrel monkey when response effort is varied. *Journal of Comparative and Physiological Psychology*, 90(2), 179.
- Alling, K., & Poling, A. (1995, may). The effects of differing response-force requirements on fixed-ratio responding of rats. *Journal of the experimental analysis of behavior*, 63(3), 331–46.
- Barofsky, I., & Hurwitz, D. (1968). Within ratio responding during fixed ratio performance. *Psychonomic Science*, 11(7), 263–264.
- Baum, W. M. (1993). Performances on ratio and interval schedules of reinforcement: Data and theory. *Journal of the Experimental Analysis of Behavior*, 59(2), 245.
- Berniker, M., O'Brien, M. K., Kording, K. P., & Ahmed, A. A. (2013). An Examination of the Generalizability of Motor Costs. *PLoS ONE*, 8(1).
- Bitterman, M. E. (1965). Phyletic differences in learning. *American Psychologist*, 20(6), 396.
- Bouton, M. E., Todd, T. P., Miles, O. W., León, S. P., & Epstein, L. H. (2013). Within- and between-session variety effects in a food-seeking habituation paradigm. *Appetite*, 66, 10–19.
- Dayan, P. (2012, apr). Instrumental vigour in punishment and reward. *The European Journal of Neuroscience*, 35(7), 1152–68.
- Eldar, E., Morris, G., & Niv, Y. (2011, sep). The effects of motivation on response rate: a hidden semi-Markov model analysis of behavioral dynamics. *Journal of neuroscience methods*, 201(1), 251–61.
- Estes, W. K. (1950). Toward a statistical theory of learning. *Psychological review*, 57(2), 94.
- Felton, M., & Lyon, D. O. (1966). The post-reinforcement pause. *Journal of the Experimental Analysis of Behavior*, 9(2), 131–134.
- Ferster, C. B., & Skinner, B. F. (1957). *Schedules of reinforcement*. Prentice-hall inc.

- Foster, M., Blackman, K., & Temple, W. (1997). Open versus closed economies: performance of domestic hens under fixed ratio schedules. *Journal of the Experimental Analysis of Behavior*, 67(1), 67.
- Gaissmaier, W., & Schooler, L. J. (2008). The smart potential behind probability matching. *Cognition*, 109(3), 416–422.
- Gallistel, C. R., & Gibbon, J. (2000). Time, rate, and conditioning. *Psychological review*, 107(2), 289.
- Gibbon, J. (1977). Scalar expectancy theory and Weber's law in animal timing. *Psychological review*, 84(3), 279.
- Greenwood, M. R., Quartermain, D., Johnson, P. R., Cruce, J. A., & Hirsch, J. (1974, nov). Food motivated behavior in genetically obese and hypothalamic-hyperphagic rats and mice. *Physiology & behavior*, 13(5), 687–92.
- Herrnstein, R. J. (1961). Relative and absolute strength of response as a function of frequency of reinforcement. *Journal of the experimental analysis of behavior*, 4(3), 267–272.
- Herrnstein, R. J., & Loveland, D. H. (1975). Maximizing and matching on concurrent ratio schedules. *Journal of the experimental analysis of behavior*, 24(1), 107.
- Hull, C. L. (1943). *Principles of Behavior*. New York: Appleton-Century.
- Iigaya, K., & Fusi, S. (2013). Dynamical regimes in neural network models of matching behavior. *Neural computation*, 25(12), 3093–3112.
- Keesey, R. E., & Kling, J. W. (1961). Amount of reinforcement and free-operant responding. *Journal of the Experimental analysis of behavior*, 4(2), 125–132.
- Kelsey, J. E., & Allison, J. (1976). Fixed-ratio lever pressing by VMH rats: Work vs accessibility of sucrose reward. *Physiology & behavior*, 17(5), 749–754.
- Keramati, M., & Gutkin, B. S. (2014). Homeostatic reinforcement learning for integrating reward collection and physiological stability. *eLife*, 3.
- Killeen, P. R. (1994, feb). Mathematical principles of reinforcement. *Behavioral and Brain Sciences*, 17, 105–172.
- Killeen, P. R. (1995, nov). Economics, ecologies, and mechanics: The dynamics of responding under conditions of varying motivation. *Journal of the Experimental Analysis of Behavior*, 64(3), 405–431.
- Killeen, P. R., & Sitomer, M. T. (2003, apr). MPR. *Behavioural Processes*, 62(1-3), 49–64.
- Kubaneck, J. (2017). Optimal decision making and matching are tied through diminishing returns. *Proceedings of the National Academy of Sciences*, 114(32), 8499–8504.
- Liberzon, D. (2011). *Calculus of Variations and Optimal Control Theory: A Concise Introduction*. Princeton University Press.
- Loewenstein, Y., Prelec, D., & Seung, H. S. (2009, oct). Operant matching as a Nash equilibrium of an intertemporal game. *Neural computation*, 21(10), 2755–73.

- Lowe, C. F., Davey, G. C. L., & Harzem, P. (1974). Effects of reinforcement magnitude on interval and ratio schedules. *Journal of the Experimental analysis of behavior*, 22(3), 553–560.
- Marshall, A. (1890). *Principles of Economics*. London: Macmillan and Co., Ltd.
- Mazur, J. E. (1982). A molecular approach to ratio schedule performance. In M. L. Commons, R. J. Herrnstein, & H. Rachlin (Eds.), *Quantitative analyses of behavior. vol. 2: Matching and maximizing accounts*. Ballinger.
- McGuire, J. T., & Kable, J. W. (2013, apr). Rational temporal predictions can underlie apparent failures to delay gratification. *Psychological review*, 120(2), 395–410.
- McSweeney, F. K. (2004). Dynamic changes in reinforcer effectiveness: Satiation and habituation have different implications for theory and practice. *The Behavior Analyst*, 27(2), 171–188.
- McSweeney, F. K., & Hinson, J. M. (1992, jul). Patterns of responding within sessions. *Journal of the Experimental Analysis of Behavior*, 58(1), 19–36.
- McSweeney, F. K., Hinson, J. M., & Cannon, C. B. (1996). Sensitization–habituation may occur during operant conditioning. *Psychological Bulletin*, 120(2), 256.
- McSweeney, F. K., Roll, J. M., & Weatherly, J. N. (1994, jul). Within-session changes in responding during several simple schedules. *Journal of the Experimental Analysis of Behavior*, 62(1), 109–132.
- Meunier, G. F., & Starratt, C. (1979). On the magnitude of reinforcement and fixed-ratio behavior. *Bulletin of the Psychonomic Society*, 13(6), 355–356.
- Niv, Y. (2007). *The effects of motivation on habitual instrumental behavior* (Ph.D. Thesis). Hebrew University.
- Niv, Y., Daw, N. D., Joel, D., & Dayan, P. (2007, apr). Tonic dopamine: opportunity costs and the control of response vigor. *Psychopharmacology*, 191(3), 507–20.
- Niv, Y., Joel, D., & Dayan, P. (2006, aug). A normative perspective on motivation. *Trends in cognitive sciences*, 10(8), 375–81.
- Niyogi, R. K., Shizgal, P., & Dayan, P. (2014). Some Work and Some Play: Microscopic and Macroscopic Approaches to Labor and Leisure. *PLoS Computational Biology*, 10(12).
- Pear, J. (2001). *The science of learning*. Psychology Press.
- Powell, R. W. (1968). The effect of small sequential changes in fixed-ratio size upon the post-reinforcement pause. *Journal of the Experimental Analysis of Behavior*, 11(5), 589–593.
- Powell, R. W. (1969). The effect of reinforcement magnitude upon responding under fixed-ratio schedules. *Journal of the Experimental Analysis of Behavior*, 12(4), 605–608.
- Premack, D., Schaeffer, R. W., & Hundt, A. (1964). Reinforcement of drinking by running: effect of fixed ratio and reinforcement time. *Journal of the experimental analysis of behavior*, 7(1), 91–96.
- Rachlin, H. (2000). *The Science of Self-Control*. Harvard University Press.

- Sakai, Y., & Fukai, T. (2008). The Actor-Critic Learning Is Behind the Matching Law: Matching Versus Optimal Behaviors. *Neural Computation*, 20(1), 227–251.
- Salimpour, Y., & Shadmehr, R. (2014, jan). Motor costs and the coordination of the two arms. *The Journal of neuroscience*, 34(5), 1806–18.
- Schulze, C., & Newell, B. R. (2016). Taking the easy way out? Increasing implementation effort reduces probability maximizing under cognitive load. *Memory & cognition*, 44(5), 806–818.
- Schulze, C., van Ravenzwaaij, D., & Newell, B. R. (2015). Of matchers and maximizers: How competition shapes choice under risk and uncertainty. *Cognitive Psychology*, 78, 78–98.
- Sidman, M., & Stebbins, W. C. (1954). Satiation effects under fixed-ratio schedules of reinforcement. *Journal of Comparative and Physiological Psychology*, 47(2), 114.
- Uno, Y., Kawato, M., & Suzuki, R. (1989). Formation and control of optimal trajectory in human multijoint arm movement. Minimum torque-change model. *Biological cybernetics*, 61(2), 89–101.
- von Neumann, J., & Morgenstern, O. (1947). *Theory of games and economic behavior*. Princeton university press Princeton.
- Vulkan, N. (2000). An economist's perspective on probability matching. *Journal of economic surveys*, 14(1), 101–118.

Appendix

Contents		
A	Value in non-deterministic schedules	32
B	Optimal actions in one-dimensional outcome space	36
C	Theorem 1: Proof	38
C.1	Simulation details of Figure 1	38
D	Theorem 2: Proof and simulation details	38
D.1	Proof of Theorem 2	38
D.2	Simulation details	40
D.3	Simulation details of Figures 4, 5	42
E	Optimal actions in multi-dimensional outcome space	42
F	Theorem 3: Proof and simulation details	44
F.1	Proof of Theorem 3	44
F.2	Simulation details	45
G	Theorem A2: Definition, proof and simulation details	46
G.1	Proof of Theorem A2	46
G.2	Simulation details	47
H	Theorem A3: Definition, proof and simulation details	48
H.1	Proof of Theorem A3	48
H.2	Simulation details	51

A Value in non-deterministic schedules

The value of a trajectory in the outcome space is the sum of the net amount of rewards that can be earned during a session. However, the amount of reward earned during a session can be non-deterministic, as for example in the case of variable-ratio and random-ratio schedules of reinforcement, actions lead to outcomes probabilistically. Similarly, the cost of earning outcomes will also be non-deterministic, since the number of responses required to earn outcomes is non-deterministic. Let's denote the cost of earning outcomes under such non-deterministic schedules by K'_V . Using this, we define the value function as the sum of the *expected* net amount

of rewards that will be earned during a session:

$$S_{0,T} = \int_0^T \mathbb{E}[\mathbf{v} \cdot A_{\mathbf{x},t} - K'_{\mathbf{v}}] dt = \int_0^T L_{\mathbf{x},\mathbf{v},t} dt, \quad (\text{A.1})$$

where the expectation is w.r.t the number of earned outcomes along each outcome dimension during dt time step. Following the above definition, we have:

$$L_{\mathbf{x},\mathbf{v},t} = \mathbb{E}[\mathbf{v} \cdot A_{\mathbf{x},t} - K'_{\mathbf{v}}], \quad (\text{A.2})$$

where $L_{\mathbf{x},\mathbf{v},t}$ is the expected net reward earned in dt time step. In the main text and in the following sections, $\mathbb{E}[\mathbf{v}]$ is denoted by \mathbf{v} for simplicity of notation. By replacing \mathbf{v} by $\mathbb{E}[\mathbf{v}]$ in equation 4 we get:

$$L_{\mathbf{x},\mathbf{v},t} = \mathbb{E}[\mathbf{v}] \cdot A_{\mathbf{x},t} - K_{\mathbb{E}[\mathbf{v}]}. \quad (\text{A.3})$$

In the main text, equation A.3 (equation 4 in the main text) was used instead of equation A.2, and the aim of this section is to show that equation A.3 and equation A.2 are equivalent. We first consider environments with one-dimensional outcome spaces, and then we extend it to the case of environments with multi-dimensional outcome spaces. We maintain the following theorem:

Theorem A1 *Assume that the cost of one response, given that the delay since the last response is τ , is as follows:*

$$C_{\tau} = a/\tau + b. \quad (\text{A.4})$$

Furthermore, assume that on average, or exactly, k responses are required to earn one unit of the outcome, and r is the number of outcomes earned. Then we have:

$$L_{x,v,t} = \mathbb{E}_r[v] A_{x,t} - K_{\mathbb{E}_r[v]}, \quad (\text{A.5})$$

where

$$K_v = vk(kav + b). \quad (\text{A.6})$$

Proof. We first provide an intuitive explanation for why the cost defined in equation A.4 is the same as the cost defined in equation A.6 in the case of fixed-ratio schedules of reinforcement (i.e., exactly k responses are required to earn an outcome). Earning the outcome at rate v implies

that the time it takes to earn the outcome is $1/\nu$, and since k responses have been executed in this period, the delay between responses is:

$$\tau = \frac{1}{k\nu}, \quad (\text{A.7})$$

and therefore using equation A.4 (equation 2 in the main text), the cost of making one response will be $ak\nu + b$. Since k responses are required for earning each outcome, the total cost of earning one unit of the outcome will be k times the cost of one response, which will be $k(ak\nu + b)$. Since the total amount of outcome earned is νdt , the total cost per unit of time will be:

$$\begin{aligned} K_\nu &= \frac{k(ak\nu + b)\nu dt}{dt} \\ &= \nu k(ak\nu + b), \end{aligned} \quad (\text{A.8})$$

which is equivalent to equation 3 and A.6.

We now show that equation A.5 and equation A.2 are equivalent in order to prove Theorem A1. Equation A.2 has two terms. As for the first term, $A_{x,t}$ can be assumed to be constant in dt time step, and therefore we have:

$$\mathbb{E}_r[\nu A_{x,t}] = \mathbb{E}_r[\nu] A_{x,t}. \quad (\text{A.9})$$

As for the second term we maintain that:

$$\mathbb{E}_r[K'_\nu] = K_{\mathbb{E}_r[\nu]}. \quad (\text{A.10})$$

To show the above relation, assume that r is the number of outcomes earned after making one response. Since according to the definition of random-ratio and variable-ratio schedules, out of N responses on average N/k will be rewarded, we have $\mathbb{E}_r[r] = 1/k$ and the expected rate of outcome earning will be:

$$\mathbb{E}_r[\nu] = \mathbb{E}_r\left[\frac{r}{\tau}\right] = \frac{1}{k\tau}. \quad (\text{A.11})$$

Furthermore, according to equation A.4 the cost of one response is $a/\tau + b$, and therefore, the cost per unit of time will be:

$$K'_v = \frac{a/\tau + b}{\tau}. \quad (\text{A.12})$$

Therefore:

$$\begin{aligned} \mathbb{E}_r[K'_v] &= \frac{a/\tau + b}{\tau} \\ &= \mathbb{E}_r[v] k(ak\mathbb{E}_r[v] + b) \text{ (using equation A.11)} \\ &= K_{\mathbb{E}_r[v]} \quad \text{(using equation A.6),} \end{aligned}$$

which proves equation A.10. Substituting equations A.10, A.9 in equation A.2 yields equation A.5, which proves the theorem.

We now turn to the case of multi-dimensional outcome spaces. The aim is to show equation A.2 is equivalent to equation A.3. To show this, we first maintain that:

$$\mathbb{E}[\mathbf{v} \cdot A_{\mathbf{x},t}] = \mathbb{E}[\mathbf{v}] \cdot A_{\mathbf{x},t}, \quad (\text{A.13})$$

which holds since $A_{\mathbf{x},t}$ can be assumed to be constant during dt time step. Next, we show that:

$$\mathbb{E}[K'_v] = K_{\mathbb{E}[\mathbf{v}]}, \quad (\text{A.14})$$

which states that $\mathbb{E}[K'_v]$ can be represented as a function of $\mathbb{E}[\mathbf{v}]$. To show this, assume r_i is the number of outcomes earned after making one response for outcome i , and τ_i is the delay between responses for outcome i . We have:

$$\mathbb{E}[v_i] = \mathbb{E}\left[\frac{r_i}{\tau_i}\right] = \frac{\mathbb{E}[r_i]}{\tau_i}, \quad (\text{A.15})$$

and therefore τ_i can be expressed as a function of $\mathbb{E}[v_i]$. Next, assume that $[C_\tau]_i$ is the cost of making one response for outcome i with delay τ_i between the responses, and τ is a vector containing the delay between responses for each outcome ($\tau = [\tau_1 \dots \tau_n]$). In dt time step, dt/τ_i responses for outcome i are made, and therefore the total cost in dt time period will be $[C_\tau]_i dt/\tau_i$, which implies that the cost for outcome i per unit of time is $[C_\tau]_i/\tau_i$. Given this, the total cost

paid for all the outcomes per unit of time will be:

$$\begin{aligned}\mathbb{E}[K'_v] &= \sum_i \frac{[C_\tau]_i}{\tau_i} \\ &= \sum_i [C_\tau]_i \frac{\mathbb{E}[v_i]}{\mathbb{E}[r_i]} \text{ (using equation A.15)} \\ &= K_{\mathbb{E}[v]},\end{aligned}$$

where we used the fact that τ in C_τ can be expressed using $\mathbb{E}[v]$ (using equation A.15), and therefore $\mathbb{E}[K'_v]$ can be expressed as a function of $\mathbb{E}[v]$, which is denoted by $K_{\mathbb{E}[v]}$ (as noted in equation A.14). Substituting equation A.14, A.13 in equation A.2 yields equation A.3.

B Optimal actions in one-dimensional outcome space

The aim is to derive optimal actions when the outcome space has one dimension. Given the reward field $A_{x,t}$, the reward of gaining dx units of outcome will be $A_{x,t}dx$, and since $dx = vdt$, the reward earned in each time step is $vA_{x,t}$. Given that K_v is the cost function (the cost paid in each time step), the net reward in each time step can be written as:

$$L_{x,v,t} = vA_{x,t} - K_v, \quad (\text{B.1})$$

and based on this, the value, which is the sum of net rewards in each time step, will be:

$$S_{0,T} = \int_0^T L_{x,v,t} dt. \quad (\text{B.2})$$

The optimal rates that maximize $S_{0,T}$ can be found using different variational calculus methods such as the Euler-Lagrange equation, or the Hamilton-Jacobi-Bellman equation (Liberzon, 2011). Here we use the Euler-Lagrange form, which sets a necessary condition for $\delta S = 0$:

$$\frac{d}{dt} \frac{\partial L}{\partial v} = \frac{\partial L}{\partial x}. \quad (\text{B.3})$$

Furthermore, since the end-point of the trajectory is free (the amount of outcomes that can be gained during a session is not limited, but the duration of the session is limited to T), the optimal

trajectory will also satisfy the transversality conditions:

$$\left. \frac{\partial L}{\partial v} \right|_{t=T} = 0, \quad (\text{B.4})$$

which implies:

$$\left. \frac{\partial K_v}{\partial v} \right|_{t=T} = A_{x,t}|_{t=T}, \quad (\text{B.5})$$

where as mentioned T is the total session duration.

By substituting equation B.1 in equation B.3 we will have:

$$\frac{d}{dt} \left(-\frac{\partial K_v}{\partial v} + A_{x,t} \right) = v \frac{dA_{x,t}}{dx}. \quad (\text{B.6})$$

The term $dA_{x,t}/dt$ has two components: the first component is the change in $A_{x,t}$ due to the change in x and the second component is due to the time-dependent changes in $A_{x,t}$:

$$\begin{aligned} \frac{dA_{x,t}}{dt} &= \frac{dx}{dt} \frac{\partial A_{x,t}}{\partial x} + \frac{\partial A_{x,t}}{\partial t} \\ &= v \frac{\partial A_{x,t}}{\partial x} + \frac{\partial A_{x,t}}{\partial t}. \end{aligned} \quad (\text{B.7})$$

Furthermore we have:

$$\frac{d}{dt} \left(\frac{\partial K_v}{\partial v} \right) = \frac{dv}{dt} \frac{\partial^2 K_v}{\partial v^2}. \quad (\text{B.8})$$

Substituting equations B.7, B.8 in equation B.6 yields:

$$\frac{dv}{dt} \left(\frac{\partial^2 K_v}{\partial v^2} \right) = \frac{\partial A_{x,t}}{\partial t}. \quad (\text{B.9})$$

In the condition that the rate of outcome earning is constant ($dv/dt = 0$), we have $x_T = vT$, which by substituting in equation B.5 yields:

$$\frac{\partial K_{v^*}}{\partial v^*} = A_{Tv^*,T}. \quad (\text{B.10})$$

The above equation will be used in order to calculate the optimal rate of outcome earning.

C Theorem 1: Proof

The cost function K_v defined in equation 3 satisfies the following relation:

$$\frac{\partial^2 K_v}{\partial v^2} > 0, \quad (\text{C.1})$$

which holds as long as at least one response is required to earn an outcome ($k > 0$), and the cost of earning outcomes is non-zero ($a > 0$).

Assuming that $\partial A_{x,t}/\partial t = 0$, and given equation C.1, the only admissible solution to equation B.9 is:

$$\frac{dv}{dt} = 0. \quad (\text{C.2})$$

Furthermore, assuming $\partial A_{x,t}/\partial t > 0$, and given equation C.1, the only admissible solution to equation B.9 is:

$$\frac{dv}{dt} > 0, \quad (\text{C.3})$$

which proves Theorem 1.

C.1 Simulation details of Figure 1. For the illustration depicted in Figure 1, following parameters were used: $a = 0.32$, $b = 0$, $k = 4$, $A_{x,t} = H = 5$. The cost and the reward were calculated at each time-step using the response rates shown in the figure. The cost were calculated using equation 3 and the reward field was assumed to be constant throughout the session.

D Theorem 2: Proof and simulation details

D.1 Proof of Theorem 2. In order to prove the theorem, we first provide a lemma. Assuming that the total session duration (T) has the probability density function $f(T)$ and that $f(T) > 0$, here we show that the expectation of the total session duration never decreases as time passes throughout the session.

Lemma 1 *Let's denote the expectation of the session duration at time t' with T'*

$$T' = \mathbb{E}[T | T > t'], \quad (\text{D.1})$$

and assume T has the following probability density function:

$$T \sim f(T), f(T) > 0. \quad (\text{D.2})$$

Then:

$$\frac{\partial T'}{\partial t'} > 0. \quad (\text{D.3})$$

Proof. We have:

$$\begin{aligned} \frac{\partial T'}{\partial t'} &= \frac{\partial \mathbb{E}[T|T > t']}{\partial t'} \\ &= \frac{\partial}{\partial t'} \left[\int_{t'}^{\infty} \frac{T f(T)}{1 - F(t')} dT \right] \\ &= \frac{\partial}{\partial t'} \left[\frac{1}{1 - F(t')} \int_{t'}^{\infty} T f(T) dT \right] \\ &= \frac{f(t')}{[1 - F(t')]^2} \int_{t'}^{\infty} T f(T) dT - \frac{t' f(t')}{1 - F(t')} \\ &= \frac{f(t')}{1 - F(t')} \underbrace{\int_{t'}^{\infty} \frac{T f(T)}{1 - F(t')} dT}_{\mathbb{E}[T|T > t']} - \frac{t' f(t')}{1 - F(t')} \\ &= \frac{f(t')}{1 - F(t')} (\mathbb{E}[T|T > t'] - t') > 0, \end{aligned} \quad (\text{D.4})$$

where $F(T)$ is the cumulative distribution function of T .

Based on the above lemma, we show that the optimal response rate is a decreasing function of t' . Since in Theorem 2 the value is calculated under the assumption that the section length is T' , and based on equation B.5, the optimal response rate satisfies the following equation:

$$\left. \frac{\partial K_v}{\partial v} \right|_{t=T'} = A_{x,t} \Big|_{t=T'}. \quad (\text{D.5})$$

Taking the derivative w.r.t to the current time in the session, i.e., t' we get:

$$\frac{dv}{dt'} \left(\frac{\partial^2 K_v}{\partial v^2} \right) = \frac{\partial T'}{\partial t'} \left(v \frac{\partial A_{x,t}}{\partial x} + \frac{\partial A_{x,t}}{\partial T'} \right). \quad (\text{D.6})$$

Theorem 2 assumes that $\partial A_{x,t}/\partial x < 0$ and $\partial A_{x,t}/\partial T' = 0$, which given equations D.3, C.1, and that $v > 0$ yields:

$$\frac{dv^*}{dt'} < 0, \quad (\text{D.7})$$

which implies that the rate of earning outcomes decreases as time passes within a session. The second part of Theorem 2 assumes that $\partial A_{x,t}/\partial x = 0$ and $\partial A_{x,t}/\partial T' = 0$, which given equation D.6 implies $dv/dt' = 0$, and therefore the optimal rate of outcome earning is not changing by the current time in the session, i.e., it is constant.

D.2 Simulation details. The simulation of the model depicted in Figure 2:right panel requires defining (i) the reward field, (ii) the cost function, and (iii) a probability distribution over the session duration. As for the probability distribution of the session duration, following McGuire et al (McGuire & Kable, 2013), we assumed that T follows a Generalized Pareto distribution:

$$F(T) = 1 - \left(1 + \frac{kT}{\sigma}\right)^{-1/k}, \quad (\text{D.8})$$

where k is a shape parameter (note that k is not the ratio-requirement here) and σ is a scale parameter, and the third parameter (location μ) was assumed to be zero. Furthermore we have:

$$F(T|T > t') = 1 - \left(1 + \frac{kT}{\sigma + kt'}\right)^{-1/k}, \quad (\text{D.9})$$

which has the following expected value:

$$\mathbb{E}[T|T > t'] = \frac{\sigma + kt'}{1 - k} + t', \quad (\text{D.10})$$

which as we expect is an increasing function of t' . Note that at point t' the expected *remaining* time until the end of the session is $\frac{\sigma + kt'}{1 - k}$.

For the simulation of the model we assumed that $k = 0.7$ and $\sigma = 18$, which represents that the initial expectation for the session duration is 60 minutes.

For the cost function, in all the simulations the cost defined in equation 3 was used, which is equivalent to the cost function used in the previous works (Dayan, 2012; Niv et al., 2007).

For the definition of the reward field, we used the framework provided by Keramati and Gutkin (2014), which provides a computational model for how the values of outcomes change

with the consumptions of the outcomes. They suggested that the dependency of the reward field on the amount of outcome earned is indirect and it is through *the motivational drive*. They conceptualized the motivational drive as the deviations of the *internal states* of a decision-maker from their homeostatic set-points. For example, let's assume that there is only one internal state, say hunger, where H denotes its homeostatic set-point (which corresponds to the deprivation level, assuming that initial value of x is zero), and there is an outcome which consuming each unit of it satisfies l units of the internal state. In this condition, the motivational drive at point x , denoted by D_x , will be:

$$D_x = \frac{1}{2}(H - lx)^2. \quad (\text{D.11})$$

Keramati and Gutkin (2014) showed that such a definition of the motivational drive has implications that are consistent with the behavioral evidence. According to the framework, the reward generated by earning δx units of the outcome is proportional to the change in the motivational drive, which can be expressed as:

$$A_{x,t} = -\frac{\partial D_x}{\partial x} = l(H - lx). \quad (\text{D.12})$$

As equation D.12 suggests, with earning more outcomes (increase in x) $A_{x,t}$ decreases. Given the above reward field, the optimal response rate of outcome earning, obtained by equation B.10, will be:

$$v^* = \frac{Hl - bk}{Tl^2 + 2ak^2}. \quad (\text{D.13})$$

Equation D.13 was used in the simulations of the “decreasing reward and unknown session duration” condition in Figure 2:right panel. The simulation of this condition was done using parameters $k = 15$, $l = 1.0$, $a = 0.05$, $b = 0.1$, $H = 450$. As for T , in each time t' within the session, the expected session duration ($E[T|T > t']$) was calculated using equation D.10, and was used as T in equation D.13.

For the “known session duration (fixed or decreasing reward)” condition in Figure 2:right panel, the same parameters as the previous condition were used, but the session duration was fixed to $T = 60$. For the “fixed reward (known or unknown session duration)” condition, we

assumed that the reward field is independent of the amount of reward earned:

$$A_{x,t} = lH. \quad (\text{D.14})$$

Given the above reward field, the optimal rate of outcome earning is:

$$v^* = \frac{Hl - bk}{2ak^2}. \quad (\text{D.15})$$

The simulation of this condition was done using parameters $k = 15$, $l = 1.0$, $a = 0.05$, $b = 0.1$, $H = 450/4$. Note that in this condition the response rate was independent of the session duration. The response rates in all the conditions were obtained by multiplying the outcome rates by k (since k responses are required to earn one unit of outcome).

D.3 Simulation details of Figures 4, 5. The simulation depicted in Figure 4 and Figure 5 are using equation D.13 with the following parameters (note that the optimal response rates were obtained by multiplying v^* by k). For Figure 4:right panel simulation parameters are $T = 50$, $k = 1$, $l = 1$, $a = 1$, $H = 8$. Parameter b is varied between 3 to 7 in order to generate the plot.

In Figure 5:left panel simulation parameters are $T = 50$, $l = 1$, $a = 0.3$, $b = 0.05$, $H = 100$. Parameter k was varied between 1 to 100 in order to generate the plot.

In Figure 5:middle panel simulation parameters are $T = 50$, $k = 1$, $l = 1$, $a = 0.3$, $b = 0.05$. Parameter H was varied between 10 to 100 in order to generate the plot.

In Figure 5:right panel simulation parameters are $T = 50$, $k = 1$, $a = 0.1$, $b = 0.1$, $H = 100$. Parameter l was varied between 0 to 1 in order to generate the plot.

E Optimal actions in multi-dimensional outcome space

The aim of this section is to derive the optimal actions in the condition that the outcome space is multi-dimensional. Optimal trajectory will satisfy the Euler-Lagrange equation along each outcome dimension:

$$\frac{d}{dt} \frac{\partial L}{\partial \mathbf{v}} = \frac{\partial L}{\partial \mathbf{x}}, \quad (\text{E.1})$$

where:

$$L_{\mathbf{x},\mathbf{v},t} = A_{\mathbf{x},t} \cdot \mathbf{v} - K_{\mathbf{v}}. \quad (\text{E.2})$$

Furthermore since the end point of the trajectory is free (the total amount of outcomes is not fixed) we have:

$$\left. \frac{\partial L}{\partial \mathbf{v}} \right|_{t=T} = \mathbf{0}. \quad (\text{E.3})$$

Using equation E.1, E.2 we have:

$$\frac{d}{dt} \left(\frac{d}{d\mathbf{v}} (-K_{\mathbf{v}} + \mathbf{v} \cdot A_{\mathbf{x},t}) \right) = \frac{d(\mathbf{v} \cdot A_{\mathbf{x},t})}{d\mathbf{x}}. \quad (\text{E.4})$$

For the right hand side of the above equation we have:

$$\frac{d(\mathbf{v} \cdot A_{\mathbf{x},t})}{d\mathbf{x}} = \mathbf{v}^T \frac{\partial A_{\mathbf{x},t}}{\partial \mathbf{x}}. \quad (\text{E.5})$$

We also have:

$$\frac{dA_{\mathbf{x},t}}{dt} = \frac{\partial A_{\mathbf{x},t}}{\partial t} + \frac{\partial A_{\mathbf{x},t}}{\partial \mathbf{x}} \mathbf{v}, \quad (\text{E.6})$$

which by substitution into equation E.4 yields:

$$\frac{d}{dt} \frac{\partial K_{\mathbf{v}}}{\partial \mathbf{v}} = \frac{\partial A_{\mathbf{x},t}}{\partial t} + \left(\frac{\partial A_{\mathbf{x},t}}{\partial \mathbf{x}} - \frac{\partial A_{\mathbf{x},t}^T}{\partial \mathbf{x}} \right) \mathbf{v}. \quad (\text{E.7})$$

We now provide two lemmas, which will be used in the proof of the following theorems.

Lemma 2 Assume that \mathbf{H} is the Hessian matrix of $K_{\mathbf{v}}$, i.e.,

$$[\mathbf{H}]_{i,j} = \frac{K_{\mathbf{v}}^2}{\partial v_i \partial v_j}, \quad (\text{E.8})$$

and furthermore assume that the cost of earning outcomes along each dimension is independent of the outcome rate on the other dimensions, i.e.,

$$\mathbf{H}_{i,j} = 0, i \neq j. \quad (\text{E.9})$$

Then:

$$\frac{d}{dt} \frac{\partial K_{\mathbf{v}}}{\partial \mathbf{v}} = \frac{d\mathbf{v}}{dt} \odot \left(\frac{\partial^2 K_{\mathbf{v}}}{\partial \mathbf{v}^2} \right), \quad (\text{E.10})$$

where $\partial^2 K_{\mathbf{v}} / \partial \mathbf{v}^2$ represents the diagonal terms of \mathbf{H} , and \odot is entrywise Hadamard product.

Proof. Using equation E.9 we have:

$$\begin{aligned} \frac{d}{dt} \frac{\partial K_{\mathbf{v}}}{\partial \mathbf{v}} &= \frac{d\mathbf{v}}{dt} \mathbf{H} \\ &= \frac{d\mathbf{v}}{dt} \odot \left(\frac{\partial^2 K_{\mathbf{v}}}{\partial \mathbf{v}^2} \right), \end{aligned} \quad (\text{E.11})$$

where the last equation comes from the fact that \mathbf{H} is a diagonal matrix.

Lemma 3 *Assuming that the reward field is conservative, i.e.,*

$$A_{\mathbf{x},t} = -\frac{\partial D_{\mathbf{x}}}{\partial \mathbf{x}}, \quad (\text{E.12})$$

then:

$$\mathbf{M} = \frac{\partial A_{\mathbf{x},t}}{\partial \mathbf{x}} - \frac{\partial A_{\mathbf{x},t}^\top}{\partial \mathbf{x}} = \mathbf{0}. \quad (\text{E.13})$$

Proof. Using equation E.12 we get:

$$\begin{aligned} [\mathbf{M}]_{i,j} &= \frac{\partial [A_{\mathbf{x},t}]_i}{\partial x_j} - \frac{\partial [A_{\mathbf{x},t}]_j}{\partial x_i} \\ &= -\frac{\partial^2 D_{\mathbf{x}}}{\partial x_j \partial x_i} + \frac{\partial^2 D_{\mathbf{x}}}{\partial x_i \partial x_j} \\ &= -\frac{\partial^2 D_{\mathbf{x}}}{\partial x_i \partial x_j} + \frac{\partial^2 D_{\mathbf{x}}}{\partial x_i \partial x_j} \text{ (using Schwarz's theorem)} \\ &= 0. \end{aligned}$$

Note that the use of Schwarz's theorem is based on the assumption that $D_{\mathbf{x}}$ is twice differentiable, which holds in the circumstances that we consider here.

F Theorem 3: Proof and simulation details

E.1 Proof of Theorem 3. Theorem 3 assumes that (i) the costs of earning outcomes are independent (equation E.9), (ii) the reward field is conservative (equation E.12), and (iii) the reward field is independent of time ($\partial A_{\mathbf{x},t}/\partial t = \mathbf{0}$). Based on Lemma 2, Lemma 3 and equation E.7 we have:

$$\frac{d\mathbf{v}}{dt} \odot \left(\frac{\partial^2 K_{\mathbf{v}}}{\partial \mathbf{v}^2} \right) = \mathbf{0}. \quad (\text{E.1})$$

Given that equation C.1 holds along each outcome dimension ($\partial^2 K_{\mathbf{v}} / \partial \mathbf{v}^2 > \mathbf{0}$), the only admissible solution to equation E.1 is $d\mathbf{v}/dt = 0$, which shows that the optimal rate of earning outcomes is constant. Since the optimal rate is constant, we have $\mathbf{x}_T = T\mathbf{v}^*$, which by substituting in boundary conditions implied by equation E.3 yields equation 9:

$$\frac{\partial K_{\mathbf{v}^*}}{\partial \mathbf{v}^*} = A_{T\mathbf{v}^*, T}, \quad (\text{E2})$$

which completes the proof the theorem.

E2 Simulation details. For the simulation of the model in Figure 6:left panel “independent cost” condition, it is assumed that the two outcomes have the same reward effect, but earning the second outcome requires l times more responses. Following Keramati et al (Keramati & Gutkin, 2014), since the two outcomes have the same reward properties we defined the motivational drive as follows:

$$D_{\mathbf{x}} = \frac{1}{2}(H - x_1 - x_2)^2, \quad (\text{E3})$$

where as mentioned $D_{\mathbf{x}}$ is the motivational drive and it represents the deviations of the internal state of the decision-maker from its homeostatic set-point (H). x_1 is the amount of \mathbf{O}_1 earned and x_2 is the amount of \mathbf{O}_2 earned, and the current motivational drive for earning outcomes depends on the difference between the total amount of earned outcomes ($x_1 + x_2$) and the homeostatic set-point (H).

Given the motivational drive, the amount of reward generated by consuming each outcome will be equal to the amount of change in the motivational drive due to the consumption of the outcomes (equation 12), and therefore, we have:

$$A_{\mathbf{x}, t} = -\frac{\partial D_{\mathbf{x}}}{\partial \mathbf{x}} = [H - x_1 - x_2, H - x_1 - x_2]. \quad (\text{E4})$$

The above equation was used as the reward field in the simulations. As for the cost function, earning one unit of \mathbf{O}_1 requires k responses on the left hand, and earning one unit of \mathbf{O}_2 requires lk responses on the right hand. Based on this and using equation 3, the cost function will be:

$$K_{\mathbf{v}} = v_1[ak^2v_1 + kb] + v_2[ak^2l^2v_2 + klb], \quad (\text{E5})$$

where v_1 is the rate of earning \mathbf{O}_1 and v_2 is the rate of earning \mathbf{O}_2 .

Using Theorem 3, the optimal response rate will be (assuming $b = 0$):

$$\text{response rate} = \left[\overbrace{\frac{kl^2 H}{Tl^2 + 2ak^2 l^2 + T}}^{\text{for left hand}}, \overbrace{\frac{kl H}{Tl^2 + 2ak^2 l^2 + T}}^{\text{for right hand}} \right], \quad (\text{E.6})$$

where as mentioned in the main text “left hand” is the response that should be taken for earning \mathbf{O}_1 , and “right hand” is the response that should be taken for earning \mathbf{O}_2 . Parameters used for simulations are $k = 1$, $l = 2$, $a = 1$, $b = 0$, $H = 100$, and $T = 20$. Note that for obtaining the response rates, the outcome rate for \mathbf{O}_1 was multiplied by k , and the outcome rate for \mathbf{O}_2 was multiplied by kl .

G Theorem A2: Definition, proof and simulation details

G.1 Proof of Theorem A2. The aim of this section is to derive optimal actions in the conditions that the costs of earning outcomes are dependent on each other. In this condition, one can assume what determines the cost is the delay between subsequent responses, either for the same or for a different outcome, i.e., the cost is proportional to the rate of earning all of the outcomes. In particular, if for earning \mathbf{O}_1 , k responses are required and for earning \mathbf{O}_2 , lk responses are required ($l \neq 1$), then the delay between subsequent responses (τ) will be $1/(kv_1 + lkv_2)$. Given equation 2, the cost of earning one unit of \mathbf{O}_1 will be $k[a(kv_1 + lkv_2) + b]$, and the cost of earning one unit of \mathbf{O}_2 will be $kl[a(kv_1 + lkv_2) + b]$. Such a cost function can be achieved by defining the cost as follows:

$$K_v = v_1[ak(kv_1 + lkv_2) + kb] + v_2[akl(kv_1 + lkv_2) + klb]. \quad (\text{G.1})$$

In the following theorem, we maintain that given the above cost function, the optimal actions are to make no response for \mathbf{O}_2 , and to make responses for \mathbf{O}_1 at a constant rate.

Theorem A2 *Given the cost function defined in equation G.1 and assuming that the two outcomes have the same reward properties, i.e.,:*

$$[A_{\mathbf{x},t}]_1 = [A_{\mathbf{x},t}]_2. \quad (\text{G.2})$$

Then the optimal actions satisfy the following equations:

$$\begin{aligned}\frac{dv_1}{dt} &= 0, \\ v_2 &= 0.\end{aligned}\tag{G.3}$$

Proof. By substituting equation G.1 in equation E.2 we have:

$$\begin{aligned}L = & -v_1 [ak(kv_1 + lk v_2) + kb] - v_2 [akl(kv_1 + lk v_2) + klb] + \\ & v_1 [A_{x,t}]_1 + v_2 [A_{x,t}]_2.\end{aligned}\tag{G.4}$$

Using the boundary condition mentioned in equation E.3 we have:

$$\begin{aligned}[A_{\mathbf{x}_T, T}]_1 - 2ak^2lv_2 - 2ak^2v_1 - bk &= 0, \\ [A_{\mathbf{x}_T, T}]_2 - 2ak^2l^2v_2 - 2ak^2lv_1 - bkl &= 0.\end{aligned}\tag{G.5}$$

Using equation G.2 we get:

$$v_1 = -lv_2 - \frac{b}{2ak},\tag{G.6}$$

which is not admissible given constraints $v_1 \geq 0$ and $v_2 \geq 0$, and therefore we assume either v_1 or v_2 will be equal to zero. The trajectory will have a higher value by setting v_2 to zero since \mathbf{O}_2 has a higher cost, and therefore the optimal solution will be $v_2 = 0$. Since $v_2 = 0$ the problem degenerates to a one-dimensional problem, in which according to Theorem 1 the optimal response rate is constant, and therefore the rate of responding for \mathbf{O}_1 will be constant, which proves the theorem.

G.2 Simulation details. For the simulation of the model in Figure 6:left panel “dependent cost” condition, it is assumed that k responses on the left lever are required to earn \mathbf{O}_1 and lk response are required on the right lever to earn \mathbf{O}_2 . Similar to the “independent cost” condition mentioned in the previous section, the reward field was assumed as follows:

$$A_{\mathbf{x},t} = -\frac{\partial D_{\mathbf{x}}}{\partial \mathbf{x}} = [H - x_1 - x_2, (H - x_1 - x_2)].\tag{G.7}$$

Since the response rate for one of the outcomes will be zero (according to Theorem A2), the problem degenerates to an environment with one action and one outcome. Using Theorem 1,

and equation B.10 the optimal response rate will be:

$$\text{response rate} = \left[\underbrace{k \frac{H - bk}{T + 2ak^2}}_{\text{for left lever}}, \underbrace{0}_{\text{for right lever}} \right]. \quad (\text{G.8})$$

Parameters used for simulations are $k = 1$, $a = 1$, $b = 0$, $H = 100$, and $T = 20$.

H Theorem A3: Definition, proof and simulation details

H.1 Proof of Theorem A3. The aim of Theorem A3 is to derive optimal actions when the reward field is non-conservative and the costs of actions are independent. An example of a non-conservative reward field is when the amount of reward that consuming an outcome produces is greater or smaller than the change in the motivational drive. For example, assume that there are two outcomes available, and the consumption of both outcomes has a similar effect on the motivational drive:

$$D_{\mathbf{x}} = \frac{1}{2}(H - x_1 - x_2)^2, \quad (\text{H.1})$$

but the reward that the second outcome generates is l times larger ($l \neq 1$) than the change it creates in the motivational drive:

$$A_{\mathbf{x},t} = \left[-l \frac{\partial D_{\mathbf{x}}}{\partial x_1}, -\frac{\partial D_{\mathbf{x}}}{\partial x_2} \right] = [l(H - x_1 - x_2), H - x_1 - x_2]. \quad (\text{H.2})$$

In this condition, $\partial[A_{\mathbf{x},t}]_1/\partial x_2 = -l$ and $\partial[A_{\mathbf{x},t}]_2/\partial x_1 = -1$, and therefore the reward of the second outcome due to the consumption of the first outcome decreases more sharply than the reward of the first outcome would, due to the consumption of the second outcome. We have:

$$\mathbf{M} = \frac{\partial A_{\mathbf{x},t}}{\partial \mathbf{x}} - \frac{\partial A_{\mathbf{x},t}^\top}{\partial \mathbf{x}} = \begin{bmatrix} 0 & 1 - l \\ l - 1 & 0 \end{bmatrix}, \quad (\text{H.3})$$

and as long as $l \neq 1$ then $\mathbf{M} \neq \mathbf{0}$, and therefore the reward field is non-conservative, because if it was conservative then according to Lemma 3 we should have $\mathbf{M} = \mathbf{0}$.

If the reward field is non-conservative, i.e., there does not exist a scalar field $D_{\mathbf{x}}$ such that $A_{\mathbf{x},t}$ satisfies equation 12, then the optimal response rates are as follows: early in the session

the decision-maker exclusively works for the outcome with the higher reward value (O_1) and, when the time remaining in the session is less than the threshold (T_c), the decision-maker then gradually starts working for the outcome with the lower reward value (O_2). More precisely we maintain the following theorem:

Theorem A3 *If the reward field follows equation H.2, $\partial A_{\mathbf{x},t}/\partial t = \mathbf{0}$, and the cost is as follows:*

$$K_{\mathbf{v}} = \frac{1}{2}mv_1^2 + \frac{1}{2}mv_2^2, \quad (\text{H.4})$$

then the optimal trajectory in the outcome space will be:

$$[v_1, v_2] = \begin{cases} \left[\frac{H(l-1)}{Tl-T_c}, 0 \right], & T-t > T_c \\ \text{arc of a circle} & T-t \leq T_c \end{cases}, \quad (\text{H.5})$$

where

$$T_c = m \frac{\arctan(1/l)}{l-1}, \quad (\text{H.6})$$

$$m = 2ak^2.$$

Proof. We have:

$$\frac{\partial A_{\mathbf{x},t}}{\partial \mathbf{x}} - \frac{\partial A_{\mathbf{x},t}^\top}{\partial \mathbf{x}} = \begin{bmatrix} 0 & 1-l \\ l-1 & 0 \end{bmatrix}, \quad (\text{H.7})$$

and based on equations E.9, H.4, E.7 we get:

$$\begin{aligned} \frac{dv_1}{dt} &= \frac{1-l}{m}v_2, \\ \frac{dv_2}{dt} &= \frac{l-1}{m}v_1. \end{aligned} \quad (\text{H.8})$$

Defining $w = (l-1)/m$, the solution to the above set of differential equations has the form:

$$\mathbf{x} = [q_1 + r/w \sin(wt + \alpha), q_2 + r/w \cos(wt + \alpha)], \quad (\text{H.9})$$

which is an arc of a circle centered at $[q_1, q_2]$, and r and α are free parameters. The parameters can be determined using the boundary condition imposed by equation E.3, and also assuming

that the initial position is $\mathbf{x} = 0$. The boundary condition in equation E.3 implies:

$$m\mathbf{v} = A_{\mathbf{x},t}|_{t=T} = \left[l\sqrt{2D_{\mathbf{x}}}, \sqrt{2D_{\mathbf{x}}} \right], \quad (\text{H.10})$$

which implies that at the end of the trajectory the rate of earning the second outcome is l times larger than the first outcome. Therefore, the general form of the trajectory will be an arc starting from the origin and ending along the above direction. Given the constraint that $\mathbf{v} \geq 0$ only the solutions in which $q_2 \leq 0$ are acceptable ones (i.e., the center of the circle is below the x -axis). Solving equation H.9 for $q_2 \leq 0$ we get:

$$T \leq T_c, \quad (\text{H.11})$$

where

$$T_c = m \frac{\arctan(1/l)}{l-1}, \quad (\text{H.12})$$

and therefore T_c is independent of H (the initial motivational drive). As such if $T \leq T_c$ (equation H.11) then the optimal trajectory will be an arc of a circle starting from the origin. Otherwise, if $T > T_c$, the optimal trajectory will be composed of two segments. In the first segment, v_2 will take the boundary condition $v_2 = 0$ and the decision-maker earns only the first outcome (the outcome with the higher reward effect). The first segment continues until the remaining time in the session satisfies equation H.11 (the remaining time is less than T_c), after which the second segment starts, which is an arc of a circle defined by equation H.9. The rate of earning the first outcome, v_1 , in the first segment of the trajectory (when $v_2 = 0$) can be obtained by calculating the rates at the beginning of the circular segment. The initial rate at the start of the circular segment is as follows:

$$r = \frac{H(l-1)}{Tl - T_c}, \quad (\text{H.13})$$

which implies that at the first segment of the trajectory we have:

$$[v_1, v_2] = \left[\frac{H(l-1)}{Tl - T_c}, 0 \right], \quad (\text{H.14})$$

which completes the proof of Theorem A3.

It is interesting to mention that there is a parallel between the trajectory that a decision-

maker takes in the outcome space, and the motion of a charged particle in a magnetic field. In the case that the outcome space is three dimensional, using equation E.7 the optimal path in the outcome space satisfies the following properties:

$$\begin{aligned} m \frac{d\mathbf{v}}{dt} &= \left(\frac{\partial A_{\mathbf{x},t}}{\partial \mathbf{x}} - \frac{\partial A_{\mathbf{x},t}^\top}{\partial \mathbf{x}} \right) \mathbf{v} \\ &= -\mathbf{v} \times \mathbf{B}, \end{aligned} \tag{H.15}$$

where \times is the cross product, \mathbf{B} is the curl of the reward field ($\mathbf{B} = \text{curl} A_{\mathbf{x},t}$), and $m = 2ak^2$. The equation H.15 in fact lays out the motion of a unit charged particle (negatively charged) with mass m in a magnetic field with magnitude \mathbf{B} .

H.2 Simulation details. Simulations shown in Figure 6:right panel are based on Theorem A3, and the parameters used are $k = 1$, $l = 1.1$, $a = 1$, $b = 0$, $H = 100$, $m = 2ak^2$.