FOR PEER REVIEW - CONFIDENTIAL

# Homeostatic reinforcement learning for integrating reward collection and physiological stability

Tracking no: 17-09-2014-RA-eLife-04811R1

Mehdi Keramati (Gatsby Computational Neuroscience Unit) and Boris Gutkin (Ecole Normale Supérieure)

**Abstract:**
Efficient regulation of internal homeostasis and defending it against perturbations requires adaptive behavioral strategies. However, the computational principles mediating the interaction between homeostatic and associative learning processes remain undefined. Here we use a definition of primary rewards, as outcomes fulfilling physiological needs, to build a normative theory showing how learning motivated behaviors may be modulated by internal states. Within this framework, we mathematically prove that seeking rewards is equivalent to the fundamental objective of physiological stability, defining the notion of physiological rationality of behavior. We further suggest a formal basis for temporal discounting of rewards by showing that discounting motivates animals to follow the shortest path in the space of physiological variables toward the desired setpoint. We also explain how animals learn to act predictively to preclude prospective homeostatic challenges, and several other behavioral patterns. Finally, we suggest a computational role for interaction between hypothalamus and the brain reward system.

**Impact statement:** We propose a normative theoretical framework for how the brain's reward learning and homeostatic regulation processes interact.

**Competing interests:** No competing interests declared

**Author contributions:**
Mehdi Keramati: Doing simulations, Deriving analytical proofs; Conception and design; Analysis and interpretation of data; Drafting or revising the article Boris Gutkin: Discussing the results; Drafting or revising the article

**Datasets:**
N/A

**Ethics:**
Human Subjects: No Animal Subjects: No

**Author Affiliation:**
Mehdi Keramati(University College London,Gatsby Computational Neuroscience Unit,United Kingdom;Departément des Etudes Cognitives, Ecole Normale Supérieure, Ecole Nor,Group for Neural Theory, INSERM U960,France) Boris Gutkin(Group for Neural Theory, LNC INSERM U960, Institute for the Study of Cognition,Ecole Normale Supérieure,France;Center for Cognition and Decision Making,National Research University Higher School of Economics,Russia)

**Dual-use research:** No

**Permissions:** Have you reproduced or modified any part of an article that has been previously published or submitted to another journal? Yes The general idea proposed in this manuscript was accepted to a machine learning conference, but none of the figures in the current manuscript were used in that conference article: Keramati, M. and Gutkin, B.S., A Reinforcement Learning Theory for Homeostatic Regulation, NIPS, 2011. ------------------------------------------- Also, an early draft of the current manuscript has been uploaded to the bioRxiv website: http://biorxiv.org/content/early/2014/06/05/005140 Published Under CCAL: Yes

1    **Homeostatic reinforcement learning for integrating reward collection and**

2    **physiological stability**

3    **Authors:**  Mehdi Keramati[1,2,*], Boris Gutkin[1,3,*]

4    **Affiliations:**

5    [1] Group for Neural Theory, INSERM U960, Département des Etudes Cognitives, Ecole Normale

6    Supérieure, Ecole Normale Supérieure - PSL Research University, 75005 Paris, France.

7    [2] Gatsby Computational Neuroscience Unit, University College London, London, UK.

8    [3] National Research University Higher School of Economics, Center for Cognition and Decision

9    Making, Moscow, Russia.

10    *Correspondence to: Mehdi@gatsby.ucl.ac.uk or Boris.gutkin@ens.fr

11

12    **Abstract**: Efficient regulation of internal homeostasis and defending it against perturbations

13    requires adaptive behavioral strategies. However, the computational principles mediating the

14    interaction between homeostatic and associative learning processes remain undefined. Here we

15    use a definition of primary rewards, as outcomes fulfilling physiological needs, to build a

16    normative theory showing how learning motivated behaviors may be modulated by internal

17    states. Within this framework, we mathematically prove that seeking rewards is equivalent to the

18    fundamental objective of physiological stability, defining the notion of physiological rationality

19    of behavior. We further suggest a formal basis for temporal discounting of rewards by showing

20    that discounting motivates animals to follow the shortest path in the space of physiological

21    variables toward the desired setpoint. We also explain how animals learn to act predictively to

22 preclude prospective homeostatic challenges, and several other behavioral patterns. Finally, we

23 suggest a computational role for interaction between hypothalamus and the brain reward system.

24

## Introduction

26 Survival requires living organisms to maintain their physiological integrity within the

27 environment. In other words, they must preserve homeostasis (e.g. body temperature, glucose

28 level, etc.). Yet, how might an animal learn to structure its behavioral strategies to obtain the

29 outcomes necessary to fulfill and even preclude homeostatic challenges? Such, efficient

30 behavioral decisions surely should depend on two brain circuits working in concert: the

31 hypothalamic homeostatic regulation (HR) system, and the cortico-basal ganglia reinforcement

32 learning (RL) mechanism. However, the computational mechanisms underlying this obvious

33 coupling remain poorly understood.

34 The previously developed classical negative feedback models of HR have tried to explain the

35 hypothalamic function in behavioral sensitivity to the "internal" state, by axiomatizing that

36 animals minimize the deviation of some key physiological variables from their hypothetical

37 setpoints (Marieb & Hoehn, 2012). To this end, a direct corrective response is triggered when a

38 deviation from setpoint is sensed or anticipated (Sibly & McFarland, 1974; Sterling, 2012). A

39 key lacuna in these models is how a simple corrective action (e.g. "go eat") in response to a

40 homeostatic deficit might be translated into a complex behavioral strategy for interacting with

41 the dynamic and uncertain external world.

42 On the other hand, the computational theory of RL has proposed a viable computational account

43 for the role of the cortico-basal ganglia system in behavioral adaptation to the "external"

44 environment, by exploiting experienced environmental contingencies and reward history

45 (Rangel, Camerer, & Montague, 2008; Sutton & Barto, 1998). Critically, this theory is built upon

46 one major axiom, namely, that the objective of behavior is to maximize reward acquisition. Yet,

47 this suite of theoretical models does not resolve how the brain constructs the reward itself, and

48 how the variability of the internal state impacts overt behavior.

49 Accumulating neurobiological evidence indicates intricate intercommunication between the

50 hypothalamus and the reward-learning circuitry (Palmiter, 2007; Rangel, 2013; Yeo & Heisler,

51 2012). The integration of the two systems is also behaviorally manifest in the classical

52 behavioral pattern of anticipatory responding in which, animals learn to act predictively to

53 preclude prospective homeostatic challenges. Moreover, the "good regulator" theoretical

54 principle implies that "every good regulator of a system must be a model of that system" (Conant

55 & Ashby, 1970), accentuating the necessity of learning a model (either explicit or implicit) of the

56 environment in order to regulate internal variables, and thus, the necessity of associative learning

57 processes being involved in homeostatic regulation.

58 Given the apparent coupling of homeostatic and learning processes, here, we propose a formal

59 hypothesis for the computations, at an algorithmic level, that may be performed in this biological

60 integration of the two systems. More precisely, inspired by previous descriptive hypotheses on

61 the interaction between motivation and learning (Hull, 1943; Mowrer, 1960; Spence, 1956), we

62 suggest a principled model for how the rewarding value of outcomes is computed as a function

63 of the animal's internal state, and of the approximated need-reduction ability of the outcome. The

64 computed reward is then made available to RL systems that learn over a state-space including

65 both internal and external states), resulting in approximate reinforcement of instrumental

66 associations that reduce or prevent homeostatic imbalance.

67  The paper is structured as follows: After giving a heuristic sketch of the theory, we show several

68  analytical, behavioral, and neurobiological results. On the basis of the proposed computational

69  integration of the two systems, we prove analytically that reward-seeking and physiological

70  stability are two sides of the same coin, and also provide a normative explanation for temporal

71  discounting of reward. Behaviorally, the theory gives a plausible unified account for anticipatory

72  responding and the rise-fall pattern of the response rate. We show that the interaction between

73  the two systems is critical in these behavioral phenomena and thus, neither classical RL nor

74  classical HR theories can account for them. Neurobiologically, we show that our model can shed

75  light on recent findings on the interaction between the hypothalamus and the reward-learning

76  circuitry, namely, the modulation of dopaminergic activity by hypothalamic signals.

77  Furthermore, we show how orosensory information can be integrated with internal signals in a

78  principled way, resulting in accounting for experimental results on consummatory behaviors, as

79  well as the pathological condition of over-eating induced by hyperpalatability. Finally, we

80  discuss limitations of the theory, compare it with other theoretical accounts of motivation and

81  internal state regulation, and outline testable predictions and future directions.

82  **Results**

83  **Theory sketch.** A self-organizing system (i.e. an organism) can be defined as a system that

84  opposes the second law of thermodynamics (Friston, 2010). In other words, biological systems

85  actively resist the natural tendency to disorder by regulating their physiological state to fall

86  within narrow bounds. This general process, known as homeostasis (Bernard, 1957; Cannon,

87  1929), includes adaptive behavioral strategies for counteracting and preventing self-entropy in

88  the face of constantly changing environments. In this sense, one would expect organisms to

89  reinforce responses that mitigate deviation of the internal state from desired "setpoints". This is

90    reminiscent of the drive-reduction theory (Hull, 1943; Mowrer, 1960; Spence, 1956) according

91    to which, one of the major mechanisms underlying reward is the usefulness of the corresponding

92    outcome in fulfilling the homeostatic needs of the organism (Cabanac, 1971). Inspired by these

93    considerations (i.e. preservation of self-order and reduction of deviations), we propose a formal

94    definition of primary reward (equivalently: reinforcer, economic utility) as the approximated

95    ability of an outcome to restore the internal equilibrium of the physiological state. We then

96    demonstrate that our formal homeostatic reinforcement learning framework accounts for some

97    phenomena that classical drive-reduction was unable to explain.

98     We first define "homeostatic space" as a multidimensional metric space in which each

99    dimension represents one physiologically-regulated variable (the horizontal plane in Figure 1).

100   The physiological state of the animal at each time $t$ can be represented as a point in this space,

101   denoted by $H_t = (h_{1,t}, h_{2,t}, .., h_{N,t})$, where $h_{i,t}$ indicates the state of the $i$-th physiological

102   variable. For example, $h_{i,t}$ can refer to the animal's glucose level, body temperature, plasma

103   osmolality, etc. The homeostatic setpoint, as the ideal internal state, can be denoted by $H^* =$

104   $(h_1^*, h_2^*, .., h_N^*)$. As a mapping from the physiological to the motivational state, we define the

105   "drive" as the distance of the internal state from the setpoint (the three-dimensional surface in

106   Figure 1):

$$D(H_t) = \sqrt[m]{\sum_{i=1}^{N} \left| h_i^* - h_{i,t} \right|^n} \qquad (1$$

107   $m$ and $n$ are free parameters that induce important nonlinear effects on the mapping between

108   homeostatic deviations and their motivational consequences. Note that for the simple case of

109   $m = n = 2$, the drive function reduces to Euclidian distance. We will later consider more

110   general nonlinear mappings in terms of classical utility theory. We will also discuss that the drive

111     function can be viewed as equivalent to the information-theoretic notion of *surprise*, defined as

112     the negative log-probability of finding an organism in a certain state $(D(H_t) = -\ln p(H_t))$.

113     Having defined drive, we can now provide a formal definition for primary reward. Let's assume

114     that as the result of an action, the animal receives an outcome $o_t$ at time $t$. The impact of this

115     outcome on different dimensions of the animal's internal state can be denoted by $K_t =$

116     $(k_{1,t}, k_{2,t}, .., k_{N,t})$. For example, $k_{i,t}$ can be the quantity of glucose received as a result of

117     outcome $o_t$. Hence, the outcome results in a transition of the physiological state from $H_t$ to

118     $H_{t+1} = H_t + K_t$ (see Figure 1) and thus, a transition of the drive state from $D(H_t)$ to $D(H_{t+1}) =$

119     $D(H_t + K_t)$. Accordingly, the rewarding value of this outcome can be defined as the consequent

120     reduction of drive:

$$
\begin{aligned}
r(H_t, K_t) \quad &= D(H_t) - D(H_{t+1}) \\
&= D(H_t) - D(H_t + K_t)
\end{aligned}
\tag{2}
$$

121     Intuitively, the rewarding value of an outcome depends on the ability of its constituting elements

122     to reduce the homeostatic distance from the setpoint or equivalently, to counteract self-entropy.

123     As discussed later, the additive effect $(K_t)$ of these constituting elements on the internal state can

124     be approximated by the orosensory properties of outcomes. We will also discuss how erroneous

125     estimation of drive reduction can potentially be a cause for maladaptive consumptive behaviors.

126     We hypothesize in this paper that the primary reward constructed as proposed in Equation 2 is

127     used by the brain's reward learning machinery to structure behavior. Incorporating this

128     physiological reward definition in a normative RL theory allows us to derive one major result of

129     our theory, which is that the rationality of behavioral patterns is geared toward maintaining

130     physiological stability.

131    **Rationality of the theory.** Here we show that our definition of reward reconciles the RL and HR

132    theories in terms of their normative assumptions: reward acquisition and physiological stability

133    are mathematically equivalent behavioral objectives. More precisely, given the proposed

134    definition of reward and given that animals discount future rewards (Chung & Herrnstein, 1967),

135    any behavioral policy, $\pi$, that maximizes the sum of discounted rewards ($SDR$) also minimizes

136    the sum of discounted deviations from the setpoint, and vice versa. In fact, starting from an

137    initial internal state $H_0$, the sum of discounted deviations ($SDD$) for a certain behavioral policy $\pi$

138    that causes the internal state to move in the homeostatic space along the trajectory $p(\pi)$, can be

139    defined as:

$$SDD_\pi(H_0) = \int_{p(\pi)} \gamma^t . D(H_t) . dt \tag{3}$$

140    Similarly, the sum of discounted rewards (SDR) for a policy $\pi$ can be defined as:

$$SDR_\pi(H_0) = \int_{p(\pi)} \gamma^t . r_t . dt = \int_{p(\pi)} \gamma^t . \big(D(H_t) - D(H_{t+dt})\big) . dt \tag{4}$$

141    It is then rather straightforward to show that for any initial state $H_0$, we will have (see Materials

142    and Methods for the proof):

$$if \ \gamma < 1: \quad \underset{\pi}{\mathrm{argmin}} \, SDD_\pi(H_0) = \underset{\pi}{\mathrm{argmax}} \, SDR_\pi(H_0) \tag{5}$$

143    where $\gamma$ is the discount factor. In other words, the same behavioral policy satisfies optimal

144    reward-seeking as well as optimal homeostatic maintenance. In this respect, reward acquisition

145    sought by the RL system is an efficient means to guide an animal's behavior toward fulfilling the

146    basic objective of defending homeostasis. Thus, our theory suggests a physiological basis for the

147    rationality of reward seeking.

148  **Normative role of temporal discounting.** In the domain of animal behavior, one fundamental

149  question is why animals should discount rewards the further they are in the future. Our theory

150  indicates that reward seeking without discounting (i.e., if $\gamma = 1$) would not lead, and may even

151  be detrimental, to physiological stability (see Materials and Methods). Intuitively, this is because

152  a future-discounting agent would always tend to expedite bigger rewards and postpone

153  punishments. Such an agent, therefore, tries to reduce homeostatic deviations (which is

154  rewarding) as soon as possible, and thus, tries to find the shortest path toward the setpoint. A

155  non-discounting agent, in contrast, can always compensate for a deviation-induced punishment

156  by reducing that deviation any time in the future.

157  While the formal proof of the necessity of discounting is given in the Materials and Methods, let

158  us give an intuitive explanation. Imagine you had to plan a one-hour hill walk from a drop-point

159  toward a pickup point, during which you wanted to minimize the height (equivalent to drive)

160  summed over the path you take. In this summation, if you give higher weights to your height in

161  the near future as compared to later times, the optimum path would be to descend the hill and

162  spend as long as possible at the bottom (i.e. homeostatic setpoint) before returning to the pickup

163  point. Equation 5 shows that this optimization is equivalent to optimizing the total discounted

164  rewards along the path, given that descending and ascending steps are defined as being

165  rewarding and punishing, respectively (equation 2).

166  In contrast, if at all points in time you give equal weights to your height, then the summed height

167  over path only depends on the drop and pickup points, since every ascend can be compensated

168  with a descend at any time. In other words, in the absence of discounting, the rewarding value of

169  a behavioral policy that changes the internal state only depends on the initial and final internal

170  states, regardless of its trajectory in the homeostatic space. Thus, when $\gamma = 1$, the values of any

171   two behavioral policies with equal net shifts of the internal state are equal, even if one policy

172   moves the internal state along the shortest path, whereas the other policy results in large

173   deviations of the internal state from the setpoint and threatens survival. These results hold for

174   any form of temporal discounting (e.g., exponential, hyperbolic). In this respect, our theory

175   provides a normative explanation for the necessity of temporal discounting of reward: to

176   maintain internal stability, it is necessary to discount future rewards.

177   **A normative account of anticipatory responding.** A paradigmatic example of behaviors

178   governed by the internal state is the anticipatory responses geared to preclude perturbations in

179   regulated variables even before any physiological depletion (negative feedback) is detectable.

180   Anticipatory eating and drinking that occur before any discernible homeostatic deviation (S C

181   Woods & Seeley, 2002), anticipatory shivering in response to a cue that predicts the cold

182   (Hjeresen, Reed, & Woods, 1986; Mansfield, Benedict, & Woods, 1983), and insulin secretion

183   prior to meal initiation (S C Woods, 1991), are only a few examples of anticipatory responding.

184   One clear example of a conditioned homeostatic response is animals' progressive tolerance to

185   ethanol-induced hypothermia. Experiments show that when ethanol injections are preceded (i.e.,

186   are predictable) by a distinctive cue, the ethanol-induced drop of the body core temperature of

187   animals diminishes along the trials (Mansfield & Cunningham, 1980). Figure 2 shows that when

188   the temperature was measured 30, 60, 90, and 120 minutes after daily injections, the drop of

189   temperature below the baseline was significant on the first day, but gradually disappeared over

190   eight days. Interestingly, in the first extinction trial on the 9[th] day where the ethanol was omitted,

191   the animal's temperature exhibited a significant increase above normal after cue presentation.

192   This indicates that the enhanced tolerance response to ethanol is triggered by the cue, and results

193   in an increase of temperature in order to compensate for the forthcoming ethanol-induced

194    hypothermia. Thus, this tolerance response is mediated by associative learning processes, and is

195    aimed at regulating temperature. Here we demonstrate that the integration of HR and RL

196    processes accounts for this phenomenon.

197    We simulate the model in an artificial environment where on every trial, the agent can choose

198    between initiating a tolerance response and doing nothing, upon observing a cue (Figure 3a). The

199    cue is then followed by a forced drop of temperature, simulating the effect of ethanol (Figure

200    3b). We also assume that in the absence of injection, the temperature does not change. However,

201    if the agent chooses to initiate the tolerance response in this condition, the temperature increases

202    gradually (Figure 3d). Thus, if ethanol injection is preceded by cue-triggered tolerance response,

203    the combined effect (Figure 3f, as superposition of Figure 3b and d) will have less deviation

204    from the setpoint as compared to when no response is taken (Figure 3b). As punishment (as the

205    opposite of reward) in our model is defined by the extent to which the deviation from the

206    setpoint increases, the 'null' response will have a bigger punishing value than the 'tolerance'

207    response and thus, the agent gradually reinforces the 'tolerance' action (Figure 3c) (More

208    precisely, the rewarding value of each action is defined by the sum of discounted drive-

209    reductions during the 24hrs upon taking that action). This results in gradual fade of the ethanol-

210    induced deviation of temperature from setpoint (Figure 3e; see Figure 3 – source data 1 for

211    simulation details).

212    Clearly, if after this learning process cue-presentation is no longer followed by ethanol injection

213    (as in the first extinction trial, E1), the cue-triggered tolerance response increases the temperate

214    beyond the setpoint (Figure 3e).

215    In general, these results show that the tolerance response caused by predicted hypothermia is an

216    optimal behavior in terms of minimizing homeostatic deviation and thus, maximizing reward.

217     Thus, this optimal homeostatic maintenance policy is acquired by associative learning

218     mechanisms.

219     Our theory implies that animals are capable of learning not only Pavlovian (e.g. shivering, or

220     tolerance to ethanol), but also instrumental anticipatory responding (e.g., pressing a lever to

221     receive warmth, in response to a cold-predicting cue). This prediction is in contrast to the theory

222     of predictive homeostasis (also known as allostasis) where anticipatory behaviors are only

223     *reflexive* responses to the predicted homeostatic deprivation upon observing cues (Sterling, 2012;

224     Stephen C Woods & Ramsay, 2007).

225     **Behavioral plausibility of drive: accounting for key phenomena.** The definition of the drive

226     function (Equation 1) in our model has two degrees of freedom: $m$ and $n$ are free parameters

227     whose values determine the properties of the homeostatic space metric. Appropriate choice of $m$

228     and $n$ $(n > m > 2)$ permits our theory to account for the following four key behavioral

229     phenomena in a unified framework. First, it accounts for the fact that the reinforcing value of an

230     appetitive outcome increases as a function of its dose $(K_t)$ (Figure 4a):

$$\frac{\partial r(H_t, K_t)}{\partial k_{j,t}} > 0 \quad : \quad for \; K_t = \left(0, 0, \dots, k_{j,t}, \dots, 0\right) \; and \; k_{j,t} > 0 \tag{6}$$

231     This is supported by the fact that in progressive ratio schedules of reinforcement rats maintain

232     higher breakpoints when reinforced with bigger appetitive outcomes, reflecting higher

233     motivation toward them (Hodos, 1961; Skjoldager, Pierre, & Mittleman, 1993). Secondly, the

234     model accounts for the potentiating effect of the deprivation level on the reinforcing value (i.e.,

235     food will be more rewarding when the animal is hungrier) (Figure 4b, c):

$$\frac{\partial r(H_t, K_t)}{\partial |h_j^* - h_{j,t}|} > 0 \quad : \quad for \; K_t = \left(0, 0, \dots, k_{j,t}, \dots, 0\right) \; and \; k_{j,t} > 0 \tag{7}$$

236    This is consistent with experimental evidence showing that the level of food deprivation in rats

237    increases the breakpoint in a progressive ratio schedule (Hodos, 1961). Note that this point

238    effectively establishes a formal extension for the "incentive" concept as defined by incentive

239    salience theory (Berridge, 2012) (Discussed later).

240    Thirdly, the theory accounts for the inhibitory effect of irrelevant drives, which is consistent with

241    a large body of behavioral experiments showing competition between different motivational

242    systems (see (Dickinson & Balleine, 2002) for a review). In other words, as the deprivation level

243    for one need increases, it inhibits the rewarding value of other outcomes that satisfy irrelevant

244    motivational systems (Figure 4d):

$$\frac{\partial r(H_t, K_t)}{\partial |h_i^* - h_{i,t}|} > 0 \quad : \quad for\ all\ i \neq j, where\ K_t = \left(0,0, \dots, k_{j,t}, \dots, 0\right)\ and\ k_{j,t} > 0 \qquad (8$$

245    Intuitively, one does not play chess, or even search for sex, on an empty stomach. As some

246    examples, calcium deprivation reduces the appetite for phosphorus, and hunger inhibits sexual

247    behavior (Dickinson & Balleine, 2002).

248    Finally, the theory naturally captures the risk-aversive nature of behavior. The rewarding value

249    in our model is a concave function of the corresponding outcome magnitude:

$$\frac{\partial^2 r(H_t, K_t)}{\partial k_{j,t}^2} < 0 \quad : \quad for\ K_t = \left(0,0, \dots, k_{j,t}, \dots, 0\right)\ and\ k_{j,t} > 0 \qquad (9$$

250    It is well known that the concavity of the economic utility function is equivalent to risk aversion

251    (Mas-Colell, Whinston, & Green, 1995). Indeed, simulating the model shows that when faced

252    with two options with equal expected payoffs, the model learns to choose the more certain option

253    as opposed to the risky one (Figure 5; see Figure 5 - source data 1 for simulation details). This is

254    because frequent small deviations from the setpoint are preferable to rare drastic deviations. In

255 fact, our theory suggests the intuition that when the expected physiological instability caused by

256 two behavioral options are equal, organisms do not choose the risky option, because the severe,

257 though unlikely, physiological instabilities that it can cause might be life-threatening.

258 Our unified explanation for the above four behavioral patterns suggests that they may all arise

259 from the functional form of the mapping from the physiological to the motivational state. In this

260 sense, we propose that these behavioral phenomena are signatures of the coupling between the

261 homeostatic and the associative learning systems. We will discuss later that $m$, $n$, and $H^*$ can be

262 regarded as free parameters of an evolutionary process, which eventually determine the

263 equilibrium density of the species.

264 Note that the equations in this section hold only when the internal state remains below the

265 setpoint. However, the drive function is symmetric with respect to the setpoint and thus,

266 analogous conclusions can be derived for other three quarters of the homeostatic space.

267 **Stepping back from the brink:** Since learning requires experience, learning whether an action

268 in a certain internal state decreases or increases the drive (i.e. is rewarding or punishing,

269 respectively) would require our model to have experienced that internal state. Living organisms,

270 however, cannot just experience internal states with extreme and life threatening homeostatic

271 deviations in order to learn that the actions that cause them are bad. For example, once the body

272 temperature goes beyond 45°C, the organism can never return.

273 We now show how our model manages this problem; i.e., it avoids voluntarily experiencing

274 extreme homeostatic deviations and hence ensures that the animal does not voluntarily endanger

275 its physiological integrity (simulations in Figure 6). In the simplest case, let us assume that the

276 model is tabula rasa: it starts from absolute ignorance about the value of state-action pairs, and

277 can freely change its internal state in the homeostatic space. In a one-dimensional space, it means

13

278     that the agent can freely increase or decrease the internal state (Figure 6 - figure supplement 1).

279     As the value of 'increase' and 'decrease' actions at all internal states are initialized to zero, the

280     agent starts by performing a random walk in the homeostatic space. However, the probability of

281     choosing the same action for $z$ times in a row decreases exponentially as $z$ increases ($p(z) =$

282     $2^{-z}$): for example, the probability of choosing "increase" is $2^{-1} = 0.5$, the probability of

283     choosing two successive "increases" is $2^{-2} = 0.25$, the probability of choosing three successive

284     "increases" is $2^{-3} = 0.125$, and so on. Thus, it is highly likely for the agent to return at least one

285     step back, before getting too far from its starting point. When the agent returns to a state it had

286     previously experienced, going in the same deviation-increasing direction will be less likely than

287     the first time (i.e., than 50-50), since the agent has already experienced the punishment caused by

288     that state-action pair once. Repetition of this process results in the agent gradually getting more

289     and more attracted to the setpoint, without ever having experienced internal states that are

290     beyond a certain limit (i.e. the brink of death).

291     Simulating the model in a one-dimensional space shows that even after starting from a rather

292     deviated internal state (initial state= 30, setpoint= 0), the agent never visits states with a

293     deviation of more than 40 units after $10^6$ trials (every action is assumed to change the state by

294     one unit) (Figure 6a; See Figure 6 - figure supplements 1 and 2, and Figure 6 - source data 1 for

295     simulation details). Also, simulating $10^5$ agents over 1500 trials (starting from state 30) shows

296     that the mean value of the internal state across all agents converges to the setpoint (Figure 5c),

297     and its variance converges to a steady-state level (Figure 5d). This shows that all agents stay

298     within certain bounds around the setpoint (The maximum deviation from the setpoint among all

299     the $10^5$ agents over the 1500 trials was 61). Also, this property of the model is shown to be

300     insensitive to the parameters of the model, like the initial internal state (Figure 6 - figure

301      supplement 3), the rate of exploration (Figure 6 - figure supplement 4), $m$ and $n$ (Figure 6 -

302      figure supplement 5), or the discount factor (Figure 6 - figure supplements 6, 7). These

303      parameters only affect the rate of convergence or the distribution over visited states, but not the

304      general property of never-visiting-drastic-deviations (existence of a boundary). Moreover, this

305      property can be generalized to multi-dimensional homeostatic spaces. Therefore, our theory

306      suggests a potential normative explanation for how animals (who might be a priori naïve about

307      potential dangers of certain internal states) would learn to avoid extreme physiological

308      instability, without ever exploring how good or bad they are.

309      **Orosensory-based approximation of post-ingestive effects.** As mentioned, we hypothesize that

310      orosensory properties of food and water provide the animal with an estimate, $\widehat{K}_t$, of their true

311      post-ingestive effect, $K_t$, on the internal state. Such association between sensory and post-

312      ingestive properties could have been developed through prior learning (Beeler et al., 2012;

313      Swithers, Baker, & Davidson, 2009; Swithers, Martin, & Davidson, 2010) or evolutionary

314      mechanisms (Breslin, 2013). Based on this sensory approximation, the only information required

315      to compute the reward (and thus the reward prediction error) is the current physiological state

316      ($H_t$) and the sensory-based approximation of the nutritional content of the outcome ($\widehat{K}_t$):

$$r\big(H_t, \widehat{K}_t\big) = D(H_t) - D\big(H_t + \widehat{K}_t\big) \tag{10}$$

317      Clearly, the evolution of the internal state itself depends only on the actual ($K_t$) post-ingestive

318      effects of the outcome. That is $H_{t+1} = H_t + K_t$.

319      According to Equation 10, the reinforcing value of food and water outcomes can be

320      approximated as soon as they are sensed/consumed, without having to wait for the outcome to be

321      digested and the drive to be reduced. This proposition is compatible with the fact that dopamine

322  neurons exhibit instantaneous, rather than delayed, burst activity in response to unexpected food

323  reward (Schneider, 1989; Schultz, Dayan, & Montague, 1997). Moreover, it might provide a

324  formal explanations for the experimental fact that intravenous injection (and even intragastric

325  intubation, in some cases) of food is not rewarding even though its drive reduction effect is equal

326  to when it is ingested orally (Miller & Kessen, 1952) (*see also* (Ren et al., 2010)*)*. In fact, if the

327  post-ingestive effect of food is estimated by its sensory properties, the reinforcing value of

328  intravenously injected food that lacks sensory aspects will be effectively zero. In the same line of

329  reasoning, the theory suggests that animals' motivation toward palatable foods, such as

330  saccharine, that have no caloric content (and thus no need-reduction effect) is due to erroneous

331  over-estimation of their drive-reduction capacity, misguided by their taste or smell. Note that the

332  rationality of our theory, as shown in Equation 5, holds only as long as $\widehat{K}_t$ is an unbiased

333  estimation of $K_t$. Otherwise, pathological conditions could emerge.

334  Last but not least, the orosensory-based approximation provides a computational hypothesis for

335  the separation of reinforcement and satiation effects. A seminal series of experiments

336  (McFarland, 1969) demonstrated that the reinforcing and satiating (i.e., need reduction) effects of

337  drinking behavior, dissociable from one another, are governed by the orosensory and alimentary

338  components of the water, respectively. Two groups of water-deprived animals learned to press a

339  green key to self-administer water orally. After this pre-training session, pressing the green key

340  had no consequence anymore, whereas pressing a novel yellow key resulted in the oral delivery

341  of water in one group, and intragastric (through a fistula) delivery of water in the second group.

342  Results showed that the green key gradually extinguished in both groups (Figure 7a, b). During

343  this time, responding on the yellow key in the oral group initially increased but then gradually

344  extinguished (rise-fall pattern; Figure 7a). The second group, however, showed no motivation for

345    the yellow key (Figure 7b). This shows that only oral, but not intragastric, self-administration of

346    water is reinforcing for thirsty animals. Our model accounts for these behavioral dynamics.

347    Simulating the model shows that the agent's subjective probability of receiving water upon

348    pressing the green key gradually decreases to zero in both groups (Figure 8c, d). As this

349    predicted outcome (alimentary content) decreases, its approximated thirst-reduction effect (equal

350    to reward in our framework) decreases as well, resulting in the extinction of pressing the green

351    key (Figure 8a, b). As for the yellow key, the oral agent initially increases the rate of responding

352    (Figure 8a) as the subjective probability of receiving water upon pressing the yellow key

353    increases (Figure 8c). Gradually, however, the internal state of the animal reaches the

354    homeostatic setpoint (Figure 8e), resulting in diminishing motivation (thirst-reduction effect) of

355    seeking water (Figure 8a). Thus, our model shows that whereas the ascending limb of the

356    response curve represents a learning effect, the descending limb is due to mitigated homeostatic

357    imbalance (i.e., unlearning vs. satiation). Notably, classical RL models only explain the

358    ascending, and classical HR models only explain the descending pattern.

359    In contrast to the oral agent, the fistula agent never learns to press the yellow key (Figure 8b).

360    This is because the approximated alimentary content attributed to this response remains zero

361    (Figure 8d) and so does its drive-reduction effect.  Note that as above, the sensory-based

362    approximation ($\widehat{K}_t$) of the alimentary effect of water in the oral and fistula cases is assumed to be

363    equal to its actual effect ($K_t$) and zero, respectively (See Figure 8 - figure supplements 1 and 2,

364    and Figure 8 - source data 1 for simulation details).

365    Our theory also suggests that in contrast to reinforcement (above), satiation is independent of the

366    sensory aspects of water and only depends on its post-ingestive effects. In fact, experiments

367    show that when different proportions of water were delivered via the two routes in different

368    groups, satiation (i.e., suppression of responding) only depended on the total amount of water

369    ingested, regardless of the delivery route (McFarland, 1969).

370    Our model accounts for these data (Figure 9), since the evolution of the internal state only

371    depends on the actual water ingested. For example, whether water is administered completely

372    orally (Figure 9, left column) or half-orally-half-intragastrically (Figure 9, right column), the

373    agent stops seeking water when the setpoint is reached. As only oral delivery is sensed, the

374    subjective outcome magnitude converges to 1 (Figure 9c) and 0.5 (Figure 9d) units for the two

375    cases, respectively. When the setpoint is reached, consuming more water results in overshooting

376    the setpoint (increasing homeostatic deviation) and thus, is punishing. Therefore, both agents

377    self-administer the same total amount of water, equal to what is required for reaching the

378    setpoint.

379    However, as the sensed amount of water is bigger in the completely-oral case, water-seeking

380    behavior is approximated to have a higher thirst-reduction effect. As a result, the reinforcing

381    value of water-seeking is higher in the oral case (as compared to the half-oral-half- intragastric

382    case) and thus, the rate of responding is higher. This, in turn, results in faster convergence of the

383    internal state to the setpoint (compare Figure 9e and f). In this respect, we predict that the

384    oral/fistula proportion affects the speed of satiation: the higher the proportion is, the faster the

385    satiety state is reached and thus, the faster the descending limb of responding emerges.

## Discussion

387    Theories of conditioning are founded on the argument that animals seek reward, while reward

388    may be defined, at least in the behaviorist approach, as what animals seek. This apparently

389    circular argument relies on the hypothetical and out-of-reach axiom of reward-maximization as

390   the behavioral objective of animals. Physiological stability, however, is an observable fact. Here,

391   we develop a coherent mathematical theory where physiological stability is put as the basic

392   axiom, and reward is defined in physiological terms. We demonstrated that reinforcement

393   learning algorithms under such a definition of physiological reward lead to optimal policies that

394   both maximize reward collection and minimize homeostatic needs. This argues for behavioral

395   rationality of physiological integrity maintenance and further shows that temporal discounting of

396   rewards is paramount for homeostatic maintenance. Furthermore, we demonstrated that such

397   integration of the two systems can account for several behavioral phenomena, including

398   anticipatory responding, the rise-fall pattern of food-seeking response, risk-aversion, and

399   competition between motivational systems. Here we argue that our framework may also shed

400   light on the computational role of the interaction between the brain reward circuitry and the

401   homeostatic regulation system; namely, the modulation of midbrain dopaminergic activity by

402   hypothalamic signals.

403   **Neural substrates.** Homeostatic regulation critically depends on sensing the internal state. In the

404   case of energy regulation, for example, the arcuate nucleus of the hypothalamus integrates

405   peripheral hormones including leptin, insulin, and ghrelin, whose circulating levels reflect the

406   internal abundance of fat, abundance of carbohydrate, and hunger, respectively (Williams &

407   Elmquist, 2012). In our model, the deprivation level has an excitatory effect on the rewarding

408   value of outcomes (equation 7) and thus on the reward prediction error (RPE). Consistently,

409   recent evidence indicates neuronal pathways through which energy state-monitoring peptides

410   modulate the activity of midbrain dopamine neurons, which supposedly carry the RPE signal

411   (Palmiter, 2007).

412    Namely, orexin neurons, which project from the lateral hypothalamus area to several brain

413    regions including the ventral tegmental area (VTA) (Sakurai et al., 1998), have been shown to

414    have an excitatory effect on dopaminergic activity (Korotkova, Sergeeva, Eriksson, Haas, &

415    Brown, 2003; Narita et al., 2006), as well as feeding behavior (Rodgers et al., 2001). Orexin

416    neurons are responsive to peripheral metabolic signals as well as to the animal's deprivation

417    level (Burdakov, Gerasimenko, & Verkhratsky, 2005), as they are innervated by orexigenic and

418    anorexigenic neural populations in the arcuate nucleus where circulating peptides are sensed.

419    Accordingly, orexin neurons are suggested to act as an interface between internal states and the

420    reward learning circuit (Palmiter, 2007). In parallel with the orexinergic pathway, ghrelin, leptin

421    and insulin receptors are also expressed on the VTA dopamine neurons, providing a further

422    direct interface between the HR and RL systems. Consistently, whereas leptin and insulin inhibit

423    dopamine activity and feeding behavior, ghrelin has an excitatory effect on them (see (Palmiter,

424    2007) for a review).

425    The reinforcing value of food outcome (and thus RPE signal) in our theory is not only modulated

426    by the internal state, but also by the orosensory information that approximates the need-reduction

427    effects. In this respect, endogenous opioids and $\mu$-opioid receptors have long been implicated in

428    the hedonic aspects of food, signaled by its orosensory properties. Systemic administration of

429    opioid antagonists decreases subjective pleasantness rating and affective responses for palatable

430    foods in humans (Yeomans & Wright, 1991) and rats (Doyle, Berridge, & Gosnell, 1993),

431    respectively. Supposedly through modulating palatability, opioids also control food intake

432    (Sanger & McCarthy, 1980) as well as instrumental food-seeking behavior (Cleary, Weldon,

433    O'Hare, Billington, & Levine, 1996). For example, opioid antagonists decrease the breakpoint in

434    progressive ratio schedules of reinforcement with food (Barbano, Le Saux, & Cador, 2009),

435  whereas opioid agonists produce the opposite effect (Solinas & Goldberg, 2005). This reflects

436  the influence of orosensory information on the reinforcing effect of food. Consistent with our

437  model, these influences have mainly been attributed to the effect of opiates on increasing

438  extracellular dopamine levels in the Nucleus Accumbens (NAc) (Devine, Leone, & Wise, 1993)

439  through its action on $\mu$-opioid receptors in the VTA and NAc (Noel & Wise, 1993; M. Zhang &

440  Kelley, 1997).

441  Such orosensory-based approximation of nutritional content, as discussed before, could have

442  been obtained through evolutionary processes (Breslin, 2013), as well as through prior learning

443  (Beeler et al., 2012; Swithers et al., 2009, 2010). In the latter case, approximations based on

444  orosensory or contextual cues can be updated so as to match the true nutritional value, resulting

445  in a rational neural/behavioral response to food stimuli (De Araujo et al., 2008).

446  **Irrational behavior: the case of over-eating.** Above, we developed a normative theory for

447  reward-seeking behaviors that lead to homeostatic stability. However, animals do not always

448  follow rational behavioral patterns, notably as exemplified in eating disorders, drug addiction,

449  and many other psychiatric diseases. Here we discuss one prominent example of such irrational

450  behavior within the context of our theory.

451  Binge eating is a disorder characterized by compulsive eating even when the person is not

452  hungry. Among the many risk factors of developing binge eating, a prominent one is having easy

453  access to hyperpalatable foods, commonly defined as those loaded with fat, sugar, or salt (Rolls,

454  2007). As an attempt to explain this risk factor, we discuss one of the points of vulnerability of

455  our theory that can induce irrational choices and thus, pathological conditions.

456  Over-seeking of hyperpalatable foods is suggested to be caused by motivational systems

457  escaping homeostatic constraints, supposedly as a result of the inability of internal satiety signals

21

458    in blocking the opioid-based stimulation of DA neurons (M. Zhang & Kelley, 2000). Stimulation

459    of $\square$-opioid receptors in the NAc, for example, is demonstrated to preferentially increase the

460    intake of high-fat food (Glass, Grace, Cleary, Billington, & Levine, 1996; M. Zhang & Kelley,

461    2000), and hyperpalatable foods are shown to trigger potent release of DA into the NAc (Nestler,

462    2001). Moreover, stimulation of the brain reward circuitry (Will, Pratt, & Kelley, 2006), as well

463    as DA receptor agonists (Cornelius, Tippmann-Peikert, Slocumb, Frerichs, & Silber, 2010) are

464    shown to induce hedonic overeating long after energy requirements are met, suggesting the

465    hyper-palatability factor to be drive-independent.

466    Motivated by these neurobiological findings, one way to formulate the overriding of the

467    homeostatic satiety signals by hyperpalatable foods is to assume that the drive-reduction reward

468    for these outcomes is augmented by a drive-independent term, $T$ ($T > 0$ for palatable foods, and

469    $T = 0$ for 'normal' foods):

$$r(H_t, K_t) = D(H_t) - D(H_t + K_t) + T \qquad (11$$

470    In other words, even when the setpoint is reached and thus, the drive-reduction effect of food is

471    zero or even negative, the term $T$ overrides this signal and results in further motivation for

472    eating (see Materials and Methods for alternative formulations of equation 11).

473    Simulating this hypothesis shows that when a deprived agent (initial internal state$= -50$) is

474    given access to normal food, the internal state converges to the setpoint (Figure 10c). When

475    hyperpalatable food with equal caloric content ($K$ is the same for both types of food) is made

476    available instead, the steady level of the internal state goes beyond the setpoint (Figure 10c).

477    Moreover, the total consumption of food is higher in the latter case (Figure 8.d), reflecting

22

478    overeating. In fact, the inflated hedonic aspect of the hyperpalatable food causes it to be sought

479    and consumed to a certain extent, even after metabolic demands are fulfilled. One might

480    speculate that such persistent overshoot would result in excess energy storage, potentially

481    leading to obesity.

482    Simulating the model in another condition where the agent has 'concurrent' access to both types

483    of foods shows significant preference of the hyperpalatable food over the normal food (Figure

484    10e), and the internal state again converges to a higher-than-setpoint level (Figure 10f). This is

485    in agreement with the evidence showing that animals strongly prefer highly palatable to less

486    palatable foods (McCrory, Suen, & Roberts, 2002). (see Figure 10 - source data 1 for simulation

487    details)

488    **Relationship to classical drive-reduction theory.** Our model is inspired by the drive reduction

489    theory of motivation, initially proposed by Clark Hull (Hull, 1943), which became the dominant

490    theory of motivation in psychology during the 1940s and 1950s. However, major criticisms have

491    been leveled against this theory over the years (Berridge, 2004; McFarland, 1969; Savage, 2000;

492    Speakman et al., 2011). Here we propose that our formal theory alleviates some of major faults

493    of the classical drive-reduction. Firstly, the classical drive-reduction does not explain

494    anticipatory responding in which animals paradoxically voluntarily increase (rather than

495    decrease) their drive deviation, even in the absence of any physiological deficit. As we

496    demonstrated, such apparently maladaptive responses are optimal in terms of both reward-

497    seeking and ensuring physiological stability, and are thus acquired by animals.

498    Secondly, the drive reduction could not explain how secondary reinforcers (e.g., money, or a

499    light that predicts food) gain motivational value, since they do not reduce the drive *per se*.

500 Because our framework integrates an RL module with the HR reward computation, the drive

501 reduction-induced reward of primary reinforcers can be readily transferred through the learning

502 process to secondary reinforcers that predict them (i.e., Pavlovian conditioning) as well as to

503 behavioral policies that lead to them (i.e., instrumental conditioning).

504 Finally, the original Hull's theory is in contradiction with the fact that intravenous injection of

505 food is not rewarding, despite its drive-reduction effect. As we showed, this could be due to the

506 orosensory-based approximation mechanism required for computing the reward.

507 Despite its limitations (discussed later), we would suggest that our modern re-formulation of the

508 drive-reduction theory subject to specific assumptions (i.e., orosensory approximation,

509 connection to RL, drive form) can serve as a framework to understand the interaction between

510 internal states and motivated behaviors.

511 **Relationship to other theoretical models.** Several previous RL-based models have also tried to

512 incorporate the internal state into the computation of reward by proposing that reward increases

513 as a linear function of deprivation level. That is, $r = w\bar{r}$, where $\bar{r}$ is a constant and $w$ is

514 proportional to the deprivation level.

515 Interestingly, a linear approximation of our proposed drive-reduction reward is equivalent to

516 assuming that the rewarding value of outcomes is equal to the multiplication of the deprivation

517 level and the magnitude of the outcome. In fact, by rewriting equation 2 for the continuous case

518 we will have:

$$r(H_t, K_t) \equiv \frac{dD(H_t + K_t)}{dK_t} \tag{12}$$

519 Using Taylor expansion, this reward can be approximated by:

$$r(H_t, K_t) \cong -K_t . \nabla D_H(H_t) + O(\nabla^2 D_H(H_t)) \tag{13}$$

520 Where    is the gradient operator, and $\nabla^2$ is the Laplace operator. Thus, a linear approximation of

521 our proposed drive-reduction reward is equivalent to assuming that the rewarding value of

522 outcomes is linearly proportional to their need-reduction capacity ($K_t$), as well as a function (the

523 gradient of drive) of the deprivation level. In this respect, our framework generalizes and

524 provides a normative basis to multiplicative forms of deprivation-modulated reward (e.g.,

525 decision field theory (Busemeyer, Townsend, & Stout, 2002), intrinsically motivated RL theory

526 (Singh, Lewis, Barto, & Sorg, 2010), and MOTIVATOR theory (Dranias, Grossberg, & Bullock,

527 2008)), where reward increases as a linear function of deprivation level. Moreover, those

528 previous models cannot account for the non-linearities arising from our model; i.e., the inhibitory

529 effect of irrelevant drives and risk aversion.

530 Whether the brain implements a nonlinear drive-reduction reward (as in equation 2) or a linear

531 approximation of it (as in equation 13) can be examined experimentally. Assuming that an

532 animal is in a slightly deprived state (Figure 11a), a linear model predicts that as the magnitude

533 of the outcome increases, its rewarding value will increase linearly (Figure 11b). A non-linear

534 reward, however, predicts an inverted U-shaped economic utility function (Figure 11b). That is,

535 the rewarding value of a large outcome can be negative, if it results in overshooting the setpoint.

536 A more recent framework that also uses a multiplicative form of deprivation-modulated reward

537 is the incentive salience theory (Berridge, 2012; J. Zhang, Berridge, Tindell, Smith, & Aldridge,

538 2009). However, in contrast to the previous models and our framework, this model assumes that

539 the rewarding value of outcomes and conditioned stimuli is learned as if the animal is in a

540 reference internal state $(\psi = 1)$. Let's denote this reward by $r(s, \psi = 1)$ for state $s$. At the time

541 of encountering state $s$ in the future, the animal uses a factor, $\psi_t$, related to its current internal

542 state, to modulate the real-time motivation of the animal: $r(s, \psi_t) = \psi_t . r(s, \psi = 1)$. In the case

543      of conditioned tolerance to hypothermic agents, however, heat-producing response is motivated

544      at the time of cue presentation, when the hypothermic agent is not administered yet. At this time,

545      the animal's internal state is not deviated and thus, the motivational element, $\psi_t$, in the incentive

546      salience theory does not provoke the tolerance response. Therefore, in our reading and unlike our

547      framework, the incentive salience theory cannot give a computational account of anticipatory

548      responding.

549      Another approach to integrate responsiveness to both internal and external states appeals to

550      approximate inference techniques from statistical physics. The free energy theory of brain

551      (Friston, 2010) proposes that organisms optimize their actions in order to minimize 'surprise'.

552      Surprise is an information-theoretic notion measuring how inconceivable it is to the organism to

553      find itself in a certain state. Assume that evolutionary pressure has compelled a species to occupy

554      a restricted set of internal states, and $p(H_t)$ indicates the probability of occupying state $H_t$, after

555      the evolution of admissible states has converged to an equilibrium density. Surprise is defined as

556      the negative log-probability of $H_t$ occurring; $-\ln p(H_t)$.

557      We propose that our notion of drive is equivalent to surprise as utilized in the free energy

558      (Friston, 2010) and interoceptive inference (Seth, 2013) frameworks. In fact, we propose that an

559      organism has an equilibrium density, $p(.)$, with the following functional form:

$$p(H_t) \propto e^{-D(H_t)} = e^{-\sqrt[m]{\sum_{i=1}^{N}|h_i^*-h_{i,t}|^n}} \tag{14}$$

560      In order to stay faithful to this probability density (and ensure the survival of genes by remaining

561      within physiological bounds), the organism minimizes surprise, which is equal to $-\ln p(H_t) =$

562      $\sqrt[m]{\sum_{i=1}^{N}|h_i^* - h_{i,t}|^n}$. This specific form of surprise is equivalent to our definition of drive

563      (equation 1). The equivalency of reward maximization and physiological stability objectives in

26

564  our model (equation 5) shows that optimizing either homeostasis or sum of discounted rewards

565  corresponds to prescribing a principle of least action applied to the surprise function.

566  Although our homeostatic RL and the free-energy theory are similar in spirit, several major

567  differences can be mentioned. Most importantly, the two frameworks should be understood at

568  different levels of analysis (Marr, 1982): the free-energy theory is a computational framework,

569  whereas our theory fits in the algorithmic/representational level. In the same line, the two

570  theories use different mathematical tools as their optimization techniques. The free energy

571  approach uses variational Bayes inference. Thus, rationality in that model is bounded by the

572  simplifying assumptions for doing "approximate" inference (namely, factorization of the

573  variational distribution over some partition of the latent variables, Laplace approximation, etc.).

574  Our approach, however, depends on tools from optimal control theory and thus, rationality is

575  constrained by the capabilities and weaknesses of the variants of the RL algorithm being used

576  (e.g. model-based vs. model-free RL). In this sense, while the notion of reward is redundant in

577  the free energy formulation, and physiological stability is achieved through gradient descent

578  function, homeostasis in our model can only be achieved through computing reward. In fact, the

579  associative learning component in our model critically depends on receiving the approximated

580  reward from the upstream regulatory component. As a result, our model remains faithful to and

581  exploits the well-developed conditioning literature in behavioral psychology, with its strengths

582  and weaknesses.

583  A further approach toward adaptive homeostatic regulation is the predictive homeostasis

584  (otherwise known as allostasis) model (Sterling, 2012) where the classical negative-feedback

585  homeostatic models is coupled with an inference system capable of anticipating forthcoming

586  demands. In this framework, anticipated demands increase current homeostatic deviation (by

587    adjusting the setpoint level) and thus, prepare the organism to meet the predicted need. Again,

588    the concept of reward is redundant in this model and motivated behaviors are directly controlled

589    by homeostatic deviation, rather than by *a priori* computed and reinforced rewarding values.

590    As alternative to the homeostatic regulation theories phrased around maintenance of setpoints,

591    another theoretical approach toward modeling regulatory systems is the "settling point" theory

592    (Berridge, 2004; Müller, Bosy-Westphal, & Heymsfield, 2010; Speakman et al., 2011;

593    Wirtshafter & Davis, 1977). According to this theory, by viewing organisms as dynamical

594    systems, what looks like a homeostatic setpoint is just the stable state of the system caused by a

595    balance of different opposing effectors on the internal variables. However, one should notice that

596    mathematically, such dynamical systems can be re-formulated as a homeostatically regulated

597    system, by writing down a potential functional for the system (or an energy function). Such an

598    energy function is equivalent to our drive function whose setpoint corresponds to the settling

599    point of the dynamical system formulation. Thus, there is equivalence between the two methods,

600    and the setpoint approach summarizes the outcome of the underlying dynamical system on the

601    regulated variables. Note that nothing precludes our framework to treat the setpoint conceptually

602    as maintained internally by an underlying system of effectors and regulators. However, the

603    setpoint/drive-function formulation conveniently allows us to derive our normative theory.

604    **Predictions.** Here we list the testable predictions of our theory, some of which put our model to

605    test against alternative proposals. Firstly, as mentioned before (Figure 9), our theory predicts that

606    the oral vs. fistula proportion in the water self-administration task (McFarland, 1969) affects the

607    speed of satiation: the higher the oral portion is, the faster the setpoint will be reached.

608    Secondly, as discussed before, our model predicts an inverted U-shaped utility function (Figure

609    11a, b). This is in contrast to the multiplicative formulations of deprivation-modulated reward.

28

610    Thirdly, our model predicts that if animals are offered with two outcomes where one outcome

611    reduces the homeostatic deviation and the other increases the deviation, the animal chooses to

612    first take the deviation-reducing and then the deviation-increasing outcome (Figure 11c, green

613    sequence), but not the other way around (Figure 11c, red sequence). This is due to the fact that

614    future deviations (and rewards) are discounted. Thus, the animal tries to postpone further

615    deviations and expedite drive-reducing outcomes.

616    Fourthly, as explained earlier, we predict that animals are capable of learning not only Pavlovian,

617    but also instrumental anticipatory responding. This is in contrast to the prediction of the

618    predictive homeostasis theory (Sterling, 2012; Stephen C Woods & Ramsay, 2007).

619    Finally, our theory predicts that upon reducing the magnitude of the outcome, a transitory burst

620    of responding should be observed. We simulate both our model (Figure 12, left) and classical

621    homeostatic regulation models (Figure 12, right) in an artificial environment where pressing a

622    lever results in the agent receiving a big outcome (1g) during the first hour, and a significantly

623    smaller outcome (0.125g) during the second hour of the experiment. According to the classical

624    models, the corrective response (lever-press) is performed when the internal state drops below

625    the setpoint. Thus, during the first hour, the agent responds with a stable rate (Figure 12e, f) in

626    order maintain the internal state above the setpoint (Figure 12d). Upon decreasing the dose, the

627    agent waits until the internal state again drops below the setpoint. Thereafter, the agent presses

628    the lever with a new rate, corresponding to the new dose. Therefore, according to this class of

629    models, response rate switches from a stable low level to a stable high level, with no burst phase

630    in between (Figure 12f).

631    According to our model, however, when the unit dose decreases from 1g to 0.125g, the agent

632    requires at least some new experiences with the outcome in order to realize that this change has

633   happened (i.e., in order to update the expected outcome associated with every action). Thus, right

634   after the dose is decreased, the agent still expects to receive a big outcome upon pressing the

635   lever. Therefore, as the objective is to minimize deviation from the setpoint (rather that staying

636   above the setpoint), the agent waits for a period equal to the normal inter-infusion interval of the

637   1g unit-dose. During this period, the internal state reaches the same lower bound as in previous

638   trials (Figure 12a). Afterward, when the agent presses the lever for the first time, it receives an

639   unexpectedly small outcome, which is not sufficient for reaching the setpoint. Thus, several

640   further responses will be needed to reach the setpoint, resulting in a burst of responding after

641   decreasing the unit dose (Figure 12b, c). After the setpoint is achieved, the agent presses the

642   lever with a lower (-than-burst) rate, in order to keep the internal state close to the setpoint.  In

643   sum, in contrast to the classical HR models, our theory predicts a temporary burst of self-

644   administration after dose reduction (See Figure 11 - source data 1 for simulation details).

645   **Limitations and future directions.** From an evolutionary perspective, physiological stability

646   and thus survival may themselves be seen as means of guaranteeing reproduction. These

647   intermediate objectives can be even violated in specific conditions and be replaced with parental

648   sacrifice. Still, we believe that homeostatic maintenance can explain a significant proportion of

649   motivated behaviors in animals. It is also noteworthy that our theory only applies to rewards that

650   have a corresponding regulatory system. How to extend our theory to rewards without a

651   corresponding homeostatic regulation system (e.g., social rewards, novelty-induced reward, etc.)

652   remains a key challenge for the future.

653   In order to put forth our formal theory we had to put forward several key constraints and

654   assumptions. As further future directions, one could relax several constraining assumptions of

655   our formal setup of the theory. For example, redesigning the model in a *partially observable*

656 condition (as opposed to the fully-observable setup we used) where the internal state observation

657 is susceptible to noise could have important implications for understanding some psychiatric

658 diseases and self-perception distortion disorders, such as anorexia nervosa. Also, relaxing the

659 assumption that the setpoint is fixed and making it adaptive to the animal's experiences could

660 explain tolerance (as elevated perception of desired setpoint) and thus, drug addiction and

661 obesity. Furthermore, relaxing the restrictive functional form of the drive function and

662 introducing more general forms could explain behavioral patterns that our model does not yet

663 account for, like asymmetric risk-aversion toward gains vs. losses (Kahneman & Tversky, 1979).

664 **Conclusion.** In a nutshell, our theory incorporates a formal physiological definition of primary

665 rewards into a novel homeostatically regulated reinforcement learning theory, allowing us to

666 prove that economically rational behaviors ensure physiological integrity. Being inspired by the

667 classic drive-reduction theory of motivation, our mathematical treatment allows for quantitative

668 results to be obtained, predictions that make the theory testable, and logical coherence. The

669 theory, with its set of formal assumptions and proofs, does not purport to explain the full gamut

670 of animal behavior, yet we believe it to be a credible step toward developing a coherent

671 mathematical framework to understand behaviors that depend on motivations stemming from

672 internal states and needs of the individual. Furthermore, this work puts forth a meta-hypothesis

673 that a number of apparently irrational behaviors regain their rationality if the internal state of the

674 individual is taken into account. Among others, the relationship between our learning-based

675 theory and evolutionary processes that shape animal a priori preferences and influence

676 behavioral patterns remains a key challenge.

677 **Materials and Methods**

678 **Rationality of the theory.** Here we show analytically that maximizing rewards and minimizing

679 deviations from the setpoint are equivalent objective functions.

680 <u>Definition</u>: A "homeostatic trajectory", denoted by $p = \{K_0, K_1, K_2, ...\}$, is an ordered sequence

681 of transitions in the $v$-dimensional homeostatic space. Each $K_i$ is a $v$-dimensional vector,

682 determining the length and direction of one transition. We also define $\mathcal{P}(H_0)$ as the set of all

683 trajectories that if start from $H_0$, will end up at $H^*$. ∎

684 <u>Definition</u>: For each homeostatic trajectory $p$ that starts from the initial motivational state $H_0$ and

685 consists of $w$ elements, we define $SDD_p(H_0)$ as the "sum of discounted drives" through that

686 trajectory:

$$SDD_p(H_0) = \sum_{t=0}^{w-1} \gamma^t . D(H_{t+1}) \tag{S1}$$

687 Where $\gamma$ is the discount factor, and $D(.)$ is the drive function. Also, starting from $H_0$, the internal

688 state evolves by $H_{t+1} = H_t + K_t$. ∎

689 <u>Definition</u>: Similarly, for each homeostatic trajectory $p$ that starts from the initial motivational

690 state $H_0$ and consists of $m$ elements, we define $SDR_p(H_0)$ as the "sum of discounted rewards"

691 through that trajectory:

$$SDR_p(H_0) = \sum_{t=0}^{w-1} \gamma^t . r_t = \sum_{t=0}^{w-1} \gamma^t . \left( D(H_t) - D(H_{t+1}) \right) \tag{S2}$$

∎

692 <u>Proposition</u>: For any initial state $H_0$, if $\gamma < 1$, we will have:

$$\underset{p\in\mathcal{P}(H_0)}{\mathrm{argmin}}\ SDD_p(H_0) = \underset{p\in\mathcal{P}(H_0)}{\mathrm{argmax}}\ SDR_p(H_0) \tag{S3}$$

693    Roughly, this means that a policy that minimizes deviation from the setpoint, also maximizes

694    acquisition of reward, and vice versa.

695    Proof: Assume that $p_i \in \mathcal{P}(H_0)$ is a sample trajectory consisting of $w_i$ transitions. As a result of

696    these transitions, the internal state will take a sequence like: $\{H_{i,0} = H_0, H_{i,1}, H_{i,2}, \dots, H_{i,w} = H^*\}$.

697    Denoting $D(H_x)$ by $D_X$ for the sake of simplicity in notation, the drive value will take the

698    following sequence: $\{D_{i,0} = D_0, D_{i,1}, D_{i,2}, \dots, D_{i,w} = D^* = 0\}$. We have:

$$SDD_{p_i}(H_0) = D_{i,1} + \gamma.D_{i,2} + \gamma^2.D_{i,3} + \dots + \gamma^{w-1}.D^* \tag{S4}$$

699    We also have:

$$\begin{aligned}
SDR_{p_i}(H_0) &= r_{i,0} + \gamma.r_{i,1} + \gamma^2.r_{i,2} + \dots + \gamma^{w-1}.r_{i,w-1} \\[6pt]
&= (D_0 - D_{i,1}) + \gamma.(D_{i,1} - D_{i,2}) + \gamma^2.(D_{i,2} - D_{i,3}) + \dots \\[6pt]
&\qquad + \gamma^{w-1}.(D_{i,w-1} - D^*) \\[6pt]
&= D_0 + (\gamma - 1).(D_{i,1} + \gamma.D_{i,2} + \gamma^2.D_{i,3} + \dots + \gamma^{w-2}.D_{i,w-1}) \\[6pt]
&= D_0 + (\gamma - 1).SDD_{p_i}(H_0)
\end{aligned} \tag{S5}$$

700    Since $D_0$ has a fixed value and $\gamma - 1 < 0$, it can be concluded that if a certain trajectory from

701    $\mathcal{P}(H_0)$ maximizes $SDR(H_0)$, it will also minimize $SDD(H_0)$, and vice versa. Thus, the

702    trajectories that satisfy these two objectives are identical. ∎

703    **Hyper-palatability effect.** For the especial case that m/n = 1, equation 11 can be rewritten as

704    follows:

$$r(H_t, K_t) \quad = D(H_t) - D(H_t + K_t) + T$$

$$= (H_t - H^*)^2 - (H_t + K_t - H^*)^2 + T$$

(S6

$$= \left(H_t - \left(H^* + \frac{T}{2K_t}\right)\right)^2 - \left(H_t + K_t - \left(H^* + \frac{T}{2K_t}\right)\right)^2$$

705  This means that the effect of $T$ is equivalent to having a simple HRL system (without term $T$)

706  whose drive function is shifted such that the new setpoint is equal to $H^* + \frac{T}{2K_t}$, where $H^*$ is the

707  setpoint of the original system. This predicts that the bigger the hyper-palatability factor $T$ is, the

708  higher the new steady state is, and the higher the real nutritional content $K_t$ of the food outcome

709  is, the less divergence of the new setpoint from the original setpoint is.

710  Equation 5 can also be re-written as:

$$r(H_t, K_t) \quad = D(H_t) - D(H_t + K_t) + T$$

$$= (H_t - H^*)^2 - (H_t + K_t - H^*)^2 + T$$

(S7

$$= \left(\left(H_t - \frac{T}{2K_t}\right) - H^*\right)^2 - \left(\left(H_t - \frac{T}{2K_t} + K_t\right) - H^*\right)^2$$

711  This can be interpreted as the effect of $T$ being equivalent to a simple HRL system (without term

712  $T$) whose internal state $H_t$ is underestimated by $\frac{T}{2K_t}$ units. That is, hyper-palatability makes the

713  behavior look like as if the subject is hungrier than what they really are.

714

34

## References:

716 Barbano, M. F., Le Saux, M., & Cador, M. (2009). Involvement of dopamine and opioids in the
717      motivation to eat: influence of palatability, homeostatic state, and behavioral paradigms.
718      *Psychopharmacology*, *203*(3), 475–87.

719 Beeler, J. A., McCutcheon, J. E., Cao, Z. F. H., Murakami, M., Alexander, E., Roitman, M. F., & Zhuang,
720      X. (2012). Taste uncoupled from nutrition fails to sustain the reinforcing properties of food. *The*
721      *European journal of neuroscience*, *36*(4), 2533–46.

722 Bernard, C. (1957). Lectures on the physiological properties and the pathological alternations of the
723      liquids of the organism: Third lecture. In L. L. Langley (Ed.), *Homeostasis: Origins of the concept,*
724      *1973* (pp. 89–100). Stroudsberg, {PA}: Dowden, Hutchinson & Ross, Inc.

725 Berridge, K. C. (2004). Motivation concepts in behavioral neuroscience. *Physiology & Behavior*, *81*(2),
726      179–209.

727 Berridge, K. C. (2012). From prediction error to incentive salience: mesolimbic computation of reward
728      motivation. *The European journal of neuroscience*, *35*(7), 1124–43.

729 Breslin, P. A. S. (2013). An evolutionary perspective on food and human taste. *Current biology : CB*,
730      *23*(9), R409–18.

731 Burdakov, D., Gerasimenko, O., & Verkhratsky, A. (2005). Physiological changes in glucose
732      differentially modulate the excitability of hypothalamic melanin-concentrating hormone and orexin
733      neurons in situ. *The Journal of Neuroscience*, *25*(9), 2429–2433.

734 Busemeyer, J. R., Townsend, J. T., & Stout, J. C. (2002). Motivational underpinnings of utility in
735      decision making: decision field theory analysis of deprivation and satiation. In S. Moore & M.
736      Oaksford (Eds.), *Emotional cognition: from brain to behaviour* (pp. 197–218). Amsterdam: John
737      Benjamins.

738 Cabanac, M. (1971). Physiological Role of Pleasure. *Science*, *173*(4002), 1103–1107.

739 Cannon, W. B. (1929). Organization for physiological homeostasis. *Physiological Reviews*, *9*, 399–431.

740 Chung, S. H., & Herrnstein, R. J. (1967). Choice and delay of reinforcement. *Journal of the experimental*
741      *analysis of behavior*, *10*(1), 67–74.

742 Cleary, J., Weldon, D. T., O'Hare, E., Billington, C., & Levine, A. S. (1996). Naloxone effects on
743      sucrose-motivated behavior. *Psychopharmacology*, *126*(2), 110–4.

744 Conant, R. C., & Ashby, W. R. (1970). Every good regulator of a system must be a model of that system.
745      *Intl. J. Systems Science*, *1*(2), 89–97.

746 Cornelius, J. R., Tippmann-Peikert, M., Slocumb, N. L., Frerichs, C. F., & Silber, M. H. (2010). Impulse
747     control disorders with the use of dopaminergic agents in restless legs syndrome: a case-control
748     study. *Sleep*, *33*(1), 81–7.

749 De Araujo, I. E., Oliveira-Maia, A. J., Sotnikova, T. D., Gainetdinov, R. R., Caron, M. G., Nicolelis, M.
750     A. L., & Simon, S. A. (2008). Food reward in the absence of taste receptor signaling. *Neuron*, *57*(6),
751     930–41.

752 Devine, D. P., Leone, P., & Wise, R. A. (1993). Mesolimbic dopamine neurotransmission is increased by
753     administration of mu-opioid receptor antagonists. *European journal of pharmacology*, *243*(1), 55–
754     64.

755 Dickinson, A., & Balleine, B. W. (2002). The role of learning in motivation. In C. R. Gallistel (Ed.),
756     *Volume 3 of Steven's Handbook of Experimental Psychology: Learning, Motivation, and Emotion*
757     (3rd ed., pp. 497–533). New York: Wiley.

758 Doyle, T. G., Berridge, K. C., & Gosnell, B. A. (1993). Morphine enhances hedonic taste palatability in
759     rats. *Pharmacology, biochemistry, and behavior*, *46*(3), 745–9.

760 Dranias, M. R., Grossberg, S., & Bullock, D. (2008). Dopaminergic and non-dopaminergic value systems
761     in conditioning and outcome-specific revaluation. *Brain research*, *1238*, 239–87.

762 Friston, K. (2010). The free-energy principle: a unified brain theory? *Nature reviews. Neuroscience*,
763     *11*(2), 127–38.

764 Glass, M. J., Grace, M., Cleary, J. P., Billington, C. J., & Levine, A. S. (1996). Potency of naloxone's
765     anorectic effect in rats is dependent on diet preference. *The American journal of physiology*, *271*(1
766     Pt 2), R217–21.

767 Hjeresen, D. L., Reed, D. R., & Woods, S. C. (1986). Tolerance to hypothermia induced by ethanol
768     depends on specific drug effects. *Psychopharmacology*, *89*(1), 45–51.

769 Hodos, W. (1961). Progressive ratio as a measure of reward strength. *Science*, *134*, 943–944.

770 Hull, C. L. (1943). *Principles of behavior: an introduction to behavior theory*. New York: Appleton-
771     Century-Crofts.

772 Kahneman, D., & Tversky, A. (1979). Prospect Theory: An Analysis of Decision under Risk.
773     *Econometrica*, *47*(2), 263–291.

774 Korotkova, T. M., Sergeeva, O. A., Eriksson, K. S., Haas, H. L., & Brown, R. E. (2003). Excitation of
775     ventral tegmental area dopaminergic and nondopaminergic neurons by orexins/hypocretins. *The*
776     *Journal of Neuroscience*, *23*(1), 7–11.

777 Mansfield, J. G., Benedict, R. S., & Woods, S. C. (1983). Response specificity of behaviorally augmented
778     tolerance to ethanol supports a learning interpretation. *Psychopharmacology*, *79*(2-3), 94–98.

779 Mansfield, J. G., & Cunningham, C. L. (1980). Conditioning and extinction of tolerance to the
780     hypothermic effect of ethanol in rats. *Journal of Comparative and Physiological Psychology*, *94*(5),
781     962–969.

782 Marieb, E. N., & Hoehn, K. (2012). *Human Anatomy & Physiology* (9th ed., p. 1264). Benjamin
783     Cummings.

784 Marr, D. (1982). *Vision*. Cambridge, Massachusetts: MIT Press.

785 Mas-Colell, A., Whinston, M. D., & Green, J. R. (1995). *Microeconomic Theory*. Cambridge: Cambridge
786     Univ. Press.

787 McCrory, M. A., Suen, V. M. M., & Roberts, S. B. (2002). Biobehavioral influences on energy intake and
788     adult weight gain. *The Journal of nutrition*, *132*(12), 3830S–3834S.

789 McFarland, D. (1969). Separation of satiating and rewarding consequences of drinking. *Physiology &*
790     *Behavior*, *4*(6), 987–989.

791 Miller, N. E., & Kessen, M. L. (1952). Reward effects of food via stomach fistula compared with those of
792     food via mouth. *Journal of Comparative and Physiological Psychology*, *45*(6), 555–564.

793 Mowrer, O. H. (1960). *Learning theory and behavior*. New York: Wiley.

794 Müller, M. J., Bosy-Westphal, A., & Heymsfield, S. B. (2010). Is there evidence for a set point that
795     regulates human body weight? *F1000 medicine reports*, *2*, 59.

796 Narita, M., Nagumo, Y., Hashimoto, S., Narita, M., Khotib, J., Miyatake, M., Sakurai, T., et al. (2006).
797     Direct involvement of orexinergic systems in the activation of the mesolimbic dopamine pathway
798     and related behaviors induced by morphine. *The Journal of neuroscience*, *26*(2), 398–405.

799 Nestler, E. J. (2001). Molecular basis of long-term plasticity underlying addiction. *Nature reviews.*
800     *Neuroscience*, *2*(2), 119–28.

801 Noel, M. B., & Wise, R. A. (1993). Ventral tegmental injections of morphine but not U-50,488H enhance
802     feeding in food-deprived rats. *Brain research*, *632*(1-2), 68–73.

803 Palmiter, R. D. (2007). Is dopamine a physiologically relevant mediator of feeding behavior? *Trends in*
804     *neurosciences*, *30*(8), 375–81.

805 Rangel, A. (2013). Regulation of dietary choice by the decision-making circuitry. *Nature neuroscience*,
806     *16*(12), 1717–24.

807 Rangel, A., Camerer, C., & Montague, P. R. (2008). A framework for studying the neurobiology of value-
808     based decision making. *Nature reviews. Neuroscience*, *9*(7), 545–56.

809 Ren, X., Ferreira, J. G., Zhou, L., Shammah-Lagnado, S. J., Yeckel, C. W., & De Araujo, I. E. (2010).
810     Nutrient selection in the absence of taste receptor signaling. *The Journal of Neuroscience*, *30*(23),
811     8012–23.

812 Rodgers, R. J., Halford, J. C., Nunes de Souza, R. L., Canto de Souza, A. L., Piper, D. C., Arch, J. R.,
813     Upton, N., et al. (2001). SB-334867, a selective orexin-1 receptor antagonist, enhances behavioural
814     satiety and blocks the hyperphagic effect of orexin-A in rats. *The European journal of neuroscience*,
815     *13*(7), 1444–52.

816 Rolls, E. T. (2007). Understanding the mechanisms of food intake and obesity. *Obesity reviews*, *8*, 67–72.

817 Sakurai, T., Amemiya, A., Ishii, M., Matsuzaki, I., Chemelli, R. M., Tanaka, H., Williams, S. C., et al.
818     (1998). Orexins and orexin receptors: a family of hypothalamic neuropeptides and G protein-
819     coupled receptors that regulate feeding behavior. *Cell*, *92*(5), 573–585.

820 Sanger, D. J., & McCarthy, P. S. (1980). Differential effects of morphine on food and water intake in
821     food deprived and freely-feeding rats. *Psychopharmacology*, *72*(1), 103–6.

822 Savage, T. (2000). Artificial motives: A review of motivation in artificial creatures. *Connection Science*,
823     *12*(3-4), 211–277.

824 Schneider, L. H. (1989). Orosensory self-stimulation by sucrose involves brain dopaminergic
825     mechanisms. *Annals of the New York Academy of Sciences*, *575*, 307–19.

826 Schultz, W., Dayan, P., & Montague, P. R. (1997). A neural substrate of prediction and reward. *Science*,
827     *275*(5306), 1593–1599.

828 Seth, A. K. (2013). Interoceptive inference, emotion, and the embodied self. *Trends in cognitive sciences*,
829     *17*(11), 565–73.

830 Sibly, R. M., & McFarland, D. J. (1974). *State Space Approach to Motivation, Motivational Control*
831     *System Analysis*. Academic Press.

832 Singh, S., Lewis, R. L., Barto, A. G., & Sorg, J. (2010). Intrinsically Motivated Reinforcement Learning:
833     An Evolutionary Perspective. *IEEE Transactions on Autonomous Mental Development*, *2*(2), 70–82.

834 Skjoldager, P., Pierre, P. J., & Mittleman, G. (1993). Reinforcer Magnitude and Progressive Ratio
835     Responding in the Rat: Effects of Increased Effort, Prefeeding, and Extinction. *Learning and*
836     *Motivation*, *24*(3), 303–343.

837 Solinas, M., & Goldberg, S. R. (2005). Motivational effects of cannabinoids and opioids on food
838     reinforcement depend on simultaneous activation of cannabinoid and opioid systems.
839     *Neuropsychopharmacology*, *30*(11), 2035–45.

840 Speakman, J. R., Levitsky, D. A., Allison, D. B., Bray, M. S., De Castro, J. M., Clegg, D. J., Clapham, J.
841     C., et al. (2011). Set points, settling points and some alternative models: theoretical options to
842     understand how genes and environments combine to regulate body adiposity. *Disease models &*
843     *mechanisms*, *4*(6), 733–45.

844 Spence, K. W. (1956). *Behavior theory and conditioning*. Westport: Greenwood Press.

845 Sterling, P. (2012). Allostasis: A model of predictive regulation. *Physiology & behavior*, *106*(1), 5–15.

846     Sutton, R. S., & Barto, A. G. (1998). *Reinforcement Learning: An Introduction*. Cambridge: MIT Press.

847     Swithers, S. E., Baker, C. R., & Davidson, T. L. (2009). General and persistent effects of high-intensity
848         sweeteners on body weight gain and caloric compensation in rats. *Behavioral neuroscience*, *123*(4),
849         772–80.

850     Swithers, S. E., Martin, A. A., & Davidson, T. L. (2010). High-intensity sweeteners and energy balance.
851         *Physiology & behavior*, *100*(1), 55–62.

852     Will, M. J., Pratt, W. E., & Kelley, A. E. (2006). Pharmacological characterization of high-fat feeding
853         induced by opioid stimulation of the ventral striatum. *Physiology & behavior*, *89*(2), 226–34.

854     Williams, K. W., & Elmquist, J. K. (2012). From neuroanatomy to behavior: central integration of
855         peripheral signals regulating feeding behavior. *Nature neuroscience*, *15*(10), 1350–5.

856     Wirtshafter, D., & Davis, J. D. (1977). Set points, settling points, and the control of body weight.
857         *Physiology & behavior*, *19*(1), 75–8.

858     Woods, S C. (1991). The eating paradox: how we tolerate food. *Psychological Review*, *98*(4), 488–505.

859     Woods, S C, & Seeley, R. J. (2002). Hunger and energy homeostasis. In C. R. Gallistel (Ed.), *Volume 3 of*
860         *Steven's Handbook of Experimental Psychology: Learning, Motivation, and Emotion* (3rd ed., pp.
861         633–68). New York: Wiley.

862     Woods, Stephen C, & Ramsay, D. S. (2007). Homeostasis: beyond Curt Richter. *Appetite*, *49*(2), 388–
863         398.

864     Yeo, G. S. H., & Heisler, L. K. (2012). Unraveling the brain regulation of appetite: lessons from genetics.
865         *Nature neuroscience*, *15*(10), 1343–9.

866     Yeomans, M. R., & Wright, P. (1991). Lower pleasantness of palatable foods in nalmefene-treated human
867         volunteers. *Appetite*, *16*(3), 249–59.

868     Zhang, J., Berridge, K. C., Tindell, A. J., Smith, K. S., & Aldridge, J. W. (2009). A Neural Computational
869         Model of Incentive Salience. *PLoS computational biology*, *5*(7).

870     Zhang, M., & Kelley, A. E. (1997). Opiate agonists microinjected into the nucleus accumbens enhance
871         sucrose drinking in rats. *Psychopharmacology*, *132*(4), 350–60.

872     Zhang, M., & Kelley, A. E. (2000). Enhanced intake of high-fat food following striatal mu-opioid
873         stimulation: microinjection mapping and fos expression. *Neuroscience*, *99*(2), 267–77.

874

881

**Figures:**

883    **Figure 1.** Schematics of the model in an exemplary two-dimensional homeostatic space. Upon

884    performing an action, the animal receives an outcome $K_t$ from the environment. The rewarding

885    value of this outcome depends on its ability to make the internal state, $H_t$, closer to the

886    homeostatic setpoint, $H^*$, and thus reduce the drive level (the vertical axis). This experienced

887    reward, denoted by $r(H_t, K_t)$, is then learned by an RL algorithm. Here a model-free RL

888    algorithm is shown in which a reward prediction error signal is computed by comparing the

889    realized reward and the expected rewarding value of the performed response. This signal is then

890    used to update the subjective value attributed to the corresponding response. Subjective values of

891    alternative choices bias the action selection process.

892

893    **Figure 2.** Experimental results (adapted from (Mansfield & Cunningham, 1980)) on the

894    acquisition and extinction of conditioned tolerance response to ethanol. (a) In each block (day) of

895    the experiment, the animal received ethanol injection after the presentation of the stimulus. (b)

896    The change in the body temperature was measured 30, 60, 90, and 120 minutes after ethanol

897    administration. Initially, the hypothermic effect of ethanol decreased the body temperature of

898    animals. After several training days, however, animals learned to activate a tolerance response

899    upon observing the stimulus, resulting in smaller deviations from the temperature setpoint. If the

900    stimulus was not followed by ethanol injection, as in the first day of extinction (E1), the

901    activation of the conditioned tolerance response resulted in an increase in body temperature. The

902    tolerance response was weakened after several (four) extinction sessions, resulting in increased

903    deviation from the setpoint in the first day of re-acquisition (R1), where presentation of the cue

904    was again followed by ethanol injection.

905

**Figure 3.** Simulation result on anticipatory responding. (a) In every trial, the simulated agent can choose between initiating a tolerance response and doing nothing, upon observing the stimulus. Regardless of the agent's choice, ethanol is administered after one hour, followed by four temperature measurements every 30 minutes. (b) Dynamics of temperature upon ethanol injection. (c) Learning curve for choosing the 'tolerance' response. (d) Dynamics of temperature upon initiating the tolerance response. (e) Temperature profile during several simulated trails. (f) Dynamics of temperature upon initiating the tolerance response, followed by ethanol administration. Plots c and e are averaged over 500 simulated agents.

914

**Figure 3 - source data 1.** Free parameters for the anticipatory responding simulation.

916

**Figure 4.** Schematic illustration of the behavioral properties of the drive function. (1) excitatory effect of the dose of outcome on its rewarding value. (b,c) excitatory effect of deprivation level on the rewarding value of outcomes: Increased deprivation increases the rewarding value of reducing drive (b), and increases the punishing value of increasing drive (c). (d) inhibitory effect of irrelevant drives on the rewarding value of outcomes.

922

**Figure 5.** Risk aversion simulation. In a conditioned place preference paradigm, the agent's presence in the left and the right compartments has equal expected payoffs, but different levels of risk (a). Panel b shows the Markov decision process of the same task. In fact, in every trial, the agent chooses whether to stay it the current compartment, or transit to the other one. The average input of energy per trial, regardless of the animal's choice, is set such that it is equal to the

928     animal's normal energy expenditure. Thus, the internal state stays close to its initial level, which

929     is equal to the setpoint here (d). The model learns to prefer the non-risky over the risky

930     compartments (c) in order to avoid severe deviations from the setpoint.

931

932     **Figure 5 - source data 1.** Free parameters for the risk-aversion simulations.

933

934     **Figure 6.** Simulations showing that the model avoids extreme deviations. Starting from 30, the

935     agent can either decrease or increase its internal state by one unit in each trial. (a) The number of

936     visits at each internal state after $10^6$ trials. (b) The drive function in the one-dimensional

937     homeostatic space. (setpoint= 0). The mean (c) and standard deviation (d) of the internal state of

938     $10^5$ agents, along 1500 trials.

939

940     **Figure 6 - figure supplement 1.** The Markov Decision Process used for simulation results

941     presented in Figure 6 and Figure 6 - figure supplements 2-7.

942

943     **Figure 6 - figure supplement 2.** Value function and choice preferences for state-action pairs

944     after simulating one agent for $10^6$ trials (as in Figure 6). The parameters of the model where as

945     follows: $\alpha = 0.4, \beta = 0.05, \gamma = 0.9, m = 3, n = 4$.

946

947     **Figure 6 - figure supplement 3.** Simulation results replicating Figure 6, with the difference that

948     the initial internal state was zero.

949

950    **Figure 6 - figure supplement 4.** Simulation results replicating Figure 6, with the difference that

951    the initial internal state was zero, and the rate of exploration, $\beta$, was 0.03.

952

953    **Figure 6 - figure supplement 5.** Simulation results replicating Figure 6, with the difference that

954    the initial internal state was zero, and also $m = n = 1$.

955

956    **Figure 6 - figure supplement 6.** Simulation results replicating Figure 6, with the difference that

957    the initial internal state was zero, and the discount factor, $\gamma$, was zero.

958

959    **Figure 6 - figure supplement 7.** Simulation results replicating Figure 6, with the difference that

960    the initial internal state was zero, and the discount factor, $\gamma$, was one (no discounting).

961

962    **Figure 6 - source data 1.** Free parameters for the simulations showing that the model avoids

963    extreme homeostatic deviations.

964

965    **Figure 7.** Experimental results (adapted from (McFarland, 1969)) on learning the reinforcing

966    effect of oral vs. intragastric delivery of water. Thirsty animals were initially trained to peck at a

967    green key to receive water orally. In the next phase, pecking at the green key had no

968    consequence, while pecking at a novel yellow key resulted in oral delivery of water in one group

969    (a), and intragastric injection of the same amount of water through a fistula in a second group

970    (b). In the first group, responding was rapidly transferred from the green to the yellow key, and

971    then suppressed. In the fistula group, the yellow key was not reinforced.

972

973   **Figure 8.** Simulation results replicating the data from (McFarland, 1969) on learning the

974   reinforcing effect of oral vs. intragastric delivery of water. As in the experiment, two groups of

975   simulated agents were pre-trained to respond on the green key to receive oral delivery of water.

976   During the test phase, the green key had no consequence, whereas a novel yellow key resulted in

977   oral delivery in one group (a) and intragastric injection in the second group (b). All agents started

978   this phase in a thirsty state (initial internal state $= 0$; setpoint $= 50$). In the oral group,

979   responding transferred rapidly from the green to the yellow key and was then suppressed (a) as

980   the internal state approached the setpoint (e). This transfer is due to gradually updating the

981   subjective probability of receiving water outcome upon responding on either key (c). In the

982   fistula group, as the water was not sensed, the outcome expectation converged to zero for both

983   keys (d) and thus, responding was extinguished (b). As a result, the internal state changed only

984   slightly (f).

985

986   **Figure 8 - figure supplement 1.** A model-based homeostatic RL system. Upon performing an

987   action in a certain state, the agent receives an outcome, $K_t$, which results in the internal state to

988   shift from $H_t$ to $H_t + K_t$. At the same time, sensory properties of the outcome are sensed by the

989   agent. Based on this information, the agent updates the state-action-outcome associations. In fact,

990   the agent learns to predict the sensory properties, $\widehat{K}_t$, of the outcome that is expected to be

991   received upon performing a certain action. Having learned these associations, the agent can

992   estimate the rewarding value of different options. That is, when the agent is in a certain state, it

993   predicts the outcome $\widehat{K}_t$, expected to result from each behavioral policy. Based on $\widehat{K}_t$ and the

994   internal state $H_t$, the agent can approximate the drive-reduction reward.

995

996   **Figure 8 - figure supplement 2.** The Markov Decision Process used for simulating the

997   reinforcing vs. satiation effects of water. At each time point, the agent can choose between doing

998   nothing (*nul*) or pecking at either the green or the yellow key.

999

1000   **Figure 8 - source data 1.** Free parameters for the reinforcing vs. satiation simulations.

1001

1002   **Figure 9.** Simulation results of the satiation test. Left column shows results for the case where

1003   water was received only orally. Rate of responding drops rapidly (a) as the internal state

1004   approaches the setpoint (e). Also, the agent learns rapidly that upon every key pecking, it

1005   receives 1.0 unit of water (c). On the right column, upon every key-peck, 0.5 unit of water is

1006   received orally, and 0.5 unit is received via the fistula. As only oral delivery is sensed by the

1007   agent, the subjective outcome-magnitude converges to 0.5 (d). As a result, the reinforcing value

1008   of key-pecking is less than that of the oral case and thus, the rate of responding is lower (b). This

1009   in turn results in slower convergence of the internal state to the setpoint (f). The MDP and the

1010   free parameters used for simulation are the same as in Figure 8.

1011

1012   **Figure 10.** Simulating over-eating of hyperpalatable vs. normal food. (a) The simulated agent

1013   can consume normal ($T = 0$) or hyperpalatable ($T > 0$) food. The nutritional content, $K$, of both

1014   foods are equal. In the single-option task (c, d), one group of animals can only choose between

1015   normal food and nothing (*nul*), whereas the other group can choose between hyperpalatable food

1016   and nothing. Starting the task in a deprived state (initial internal state=-50), the internal state of

1017   the second, but not the first, group converges to a level above the setpoint (c) and the total

1018   consumption of food is higher in this group (d). In the multiple-choice task, the agents can

46

1019    choose between normal food, hyperpalatable food, and nothing (b). Results show that the

1020    hyperpalatable food is preferred over the normal food (e) and the internal state is defended at a

1021    level beyond the setpoint (f). See **Figure 10 - figure supplement 1** for simulation details.

1022

1023    **Figure 10 - source data 1.** Free parameters for the over-eating simulations.

1024

1025    **Figure 11.** Behavioral predictions of the model. (a) Differential predictions of the multiplicative

1026    (linear) and drive-reduction (non-linear) forms of reward. In our model, assuming that the

1027    internal state is at $h_t$ (a), outcomes larger than $h^* - h_t$ result in overshooting the setpoint and

1028    thus a declining trend of the rewarding value (b). Previous models, however, predict the

1029    rewarding value to increase linearly as the outcome increases in magnitude. (c) Our model

1030    predicts that when given a choice between two options with equal net effects on the internal

1031    state, animals choose the option that first results in reducing the homeostatic deviation and then

1032    is followed by an increase in deviation (green), as compared to a reversed-order option (red).

1033

1034    **Figure 12.** Simulation results, predicting a transitory burst of responding upon reducing the dose

1035    of outcome. Our model (left column) and negative-feedback models (right column) are simulated

1036    is a process where responding yields big and small outcomes, during the first and second hours

1037    of the experiment, respectively. Our model predicts a short-term burst of responding after the

1038    dose reduction, followed by regular and escalated response rate (b, c). Classical HR models,

1039    however, predict an immediate transition from a steady low to a steady high response rate (e, f).

1040    See **Figure 12 - figure supplements 1 and 2** for simulation details.

1041

1042     **Figure 12 - figure supplement 1.** The Markov Decision Process used for the within-session

1043     dose-change simulation.

1044

1045     **Figure 12 - source data 1.** Free parameters for the within-session dose-change simulation.

Organism

drive

$d(H_t)$

$d(H_{t+1})$

$h_2$ : temperature

$H_t$   $K_t$   $H_{t+1}$

$k_{1,t}$

$H^* = (h_1^*, h_2^*)$

$h_1$ : glucose

$r(H_t, K_t)$

Reward   **+**   **−**

$\otimes$

Prediction error

Value(State)

Estimated values

Values

Action Selection

Outcome        State   Action

**Environment**

**a**



**b**



Experimental Data

**a** 60' 30' 30' 30' 30' 21 hrs

tolerance

$s_1$

$s_0$

nul

$s_2$

30' 60' 90' 120'

$s_3$

**b** Ethanol injection

Ethanol injection

Hour

Change of temperature

**c**

Response probability

Blocks

**d** Tolerance response

Tolerance response

Hour

Change of temperature

**e**

30'
60'
90'
120'

Change of temperature

Acquisition Blocks

**f** Tolerance response + Ethanol injection

30'
60'
90'
120'

Ethanol injection
Tolerance response

Hour

Change of temperature

**a**

$h_2$ : temperature

$k_{1,t}$

$k_{2,t}$

$h_1$: glucose

$(h_1^*, h_2^*)$

**b**

$h_2$ : temperature

$k_{1,t}$

$k_{1,t}$

$(h_1, h_2)$

$(h_1', h_2)$

$h_1$: glucose

$(h_1^*, h_2^*)$

**c**

$h_2$ : temperature

$k_{1,t}$

$H$

$k_{1,t}$

$H$

$h_1$: glucose

$(h_1^*, h_2^*)$

**d**

$h_2$ : temperature

$(h_1', h_2)$

$k_{1,t}$

$(h_1, h_2')$

$k_{1,t}$

$h_1$: glucose

$(h_1^*, h_2^*)$

**a**

2 units of energy with *p=1*

8 units of energy with *p=0.25*

2 units of energy

**b**

$a'$

$a$ $s$ $s$ $a$

$a'$

**c**

Non-risky state
Risky state

Place preference

0 100 200 300
Time

**d**

Internal state

0 100 200 300
Time

**a**



**b**



**c**



**d**

**a** Oral reward

**b** Fistula reward

**a** Oral reward

**b** Fistula reward

**a** Oral reward

**b** Oral/Fistula reward

**a**

T = 0        T > 0

**b**

*nul*

*s*

*eat tasty food*        *eat normal food*

**c**

Tasty food ——
Normal food ——

Internal state

50
25
0
-25
-50

0    100   200   300   400
Time

**d**

Total consumption

370
360
350
340

Normal food    Hyperpalatable food

**e**

Tasty food ——
Normal food ——
nul ——

Choice probabilities

1.0
0.8
0.6
0.4
0.2
0.0

0    100   200   300   400
Time

**f**

Internal state

50
25
0
-25
-50

0    100   200   300   400
Time

**a**



$D(h_t)$

Current internal state

$h_t$    $h^*$    $h$

**b**

$r(h_t, k_t)$

$k_t$

— Multiplicative reward
— Drive reduction reward

**c**

$h_2$

$H^*$

$H'$

$A$

$B$

$H$

$h_1$

**a** Our model

**d** Previous models

Legend:
- 1.000 g/response (green)
- 0.125 g/response (blue)
- Burst phase (red)
- Blindness phase (olive/dark yellow)