

# 1 Where does value come from?

2

3 Keno Juechems\* and Christopher Summerfield

4

5 Department of Experimental Psychology

6 University of Oxford

7 Radcliffe Observatory Quarter

8 Woodstock Road, Oxford, UK

9

10

## 11 **Abstract**

12

13 The computational framework of reinforcement learning (RL) has allowed us to both understand  
14 biological brains and build successful artificial agents. However, in this article we highlight open  
15 challenges for RL as a model of animal behaviour in natural environments. We ask how the external  
16 reward function is designed for biological systems, and how we can account for the context sensitivity  
17 of valuation. We summarise both old and new theories proposing that animals track current and  
18 desired internal states and seek to minimise the distance to goal across multiple value dimensions. We  
19 suggest that this framework can readily account for canonical phenomena observed in the fields of  
20 psychology, behavioural ecology, and economics, and recent findings from brain imaging studies of  
21 value-guided decision-making.

22

23

24 \*Correspondence: [keno.juechems@psy.ox.ac.uk](mailto:keno.juechems@psy.ox.ac.uk)

25

## 26 **Keywords**

27 Reinforcement learning; value; homeostasis; reward; goal-directed decision-making; medial  
28 prefrontal cortex

## 29 The reward hypothesis and the reward paradox

30

31 What is it that animals seek to achieve by their behaviour? For evolutionary biologists, there is a simple  
32 answer to this question – that animal behaviour has evolved to maximise reproductive fitness, i.e. the  
33 propensity to produce offspring that will carry forth our genes. However, for those studying the  
34 behaviour of humans and animals in the fields of ethology, economics, psychology and neuroscience,  
35 this answer is only partly satisfying. Critically, it stops short of specifying how animals can learn which  
36 behaviours will increase or decrease their fitness. Instead, behavioural scientists propose that the  
37 environment furnishes reinforcement signals that indicate the likely costs or benefits of an action and  
38 assume that these signals can be directly sensed by the agent, allowing them to behave in ways that  
39 maximise utility over the short or long term. In machine learning research, this is sometimes called the  
40 “reward hypothesis” [1].

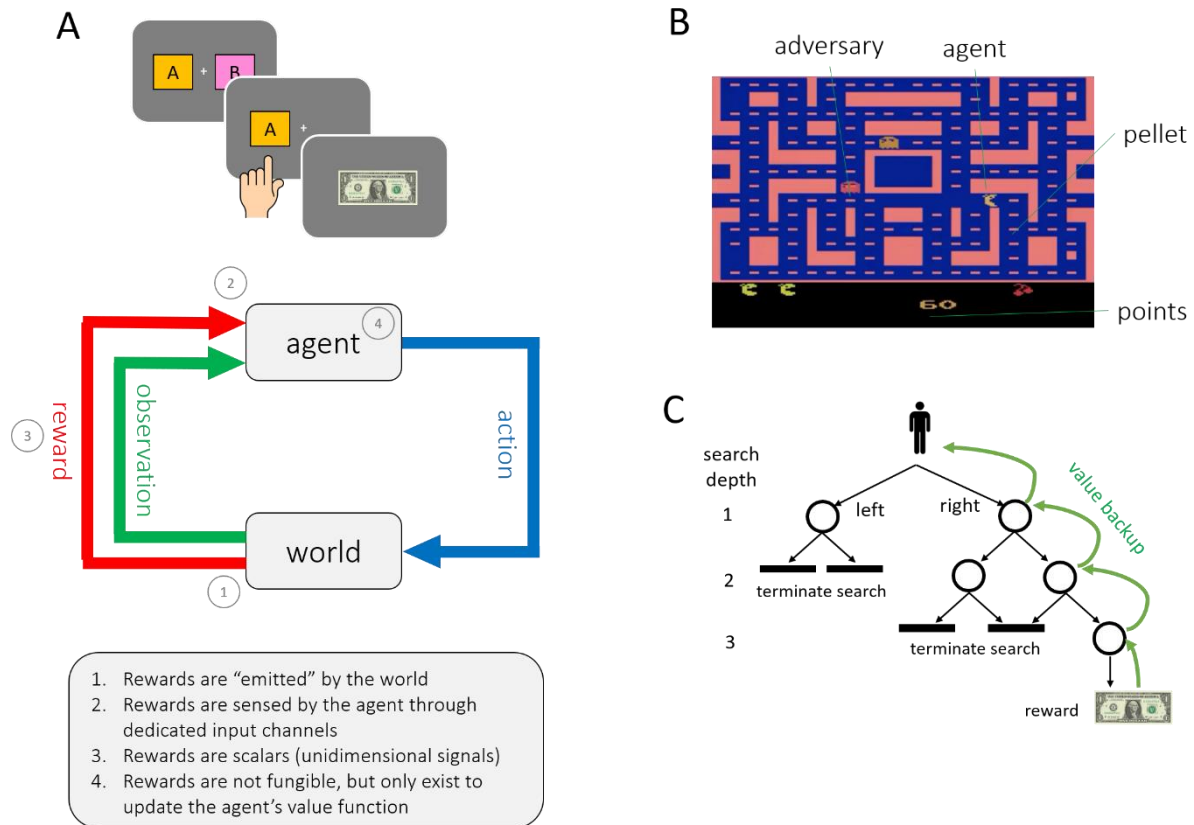
41

42 Over the past half century, the fields psychology and artificial intelligence (AI) have jointly developed a  
43 normative theory, known as reinforcement learning (RL), that is founded on the reward hypothesis [2].  
44 RL conceives of each encounter with the world as involving three distinct stages: (i) an agent receives  
45 sensory observations from the environment; (ii) the agent takes actions that influence its future states;  
46 and (iii) the agent receives a scalar reward signal that is emitted by the environment, and that is  
47 processed by the agent via a dedicated input channel (**Fig. 1a**). This conceptualisation lends itself  
48 naturally to modelling the canonical laboratory paradigm for operant behaviour, whereby an animal  
49 receives a sensory stimulus (e.g. a visual signal or tone), takes an action (e.g. presses a button or licks a  
50 spout) and receives a positive or negative reinforcer (e.g. liquid, food, or money). The framework  
51 successfully accounts for a wide array of habit-based and goal-directed behaviours in humans and other  
52 animals [3], and prominently explains the neural signals that accompany the expectation and delivery  
53 of reinforcers [4]. RL has furnished the canonical computational framework for understanding value-  
54 guided decision-making in humans and other primates, and has been widely used to explain how  
55 different brain regions participate in valuation and choice [2]. Moreover, RL models are currently  
56 offering some of the most exciting advances in AI research, producing artificial agents that can behave  
57 intelligently in complex, dynamic environments such as board and video games, thereby providing a  
58 proof of concept that intelligent behaviour can emerge *de novo* under the assumptions of the reward  
59 hypothesis [5, 6].

60

61 The RL framework offers an important set of computational tools for understanding animal learning  
62 and behaviour. However, as a theory of biological intelligence, RL focusses on how rewards influence  
63 behaviour but does not specify the provenance of the external reward function against which behaviour  
64 is optimised. In machine learning, the reward function for an RL model is hand-designed by the  
65 researcher. For example, when training an autonomous vehicle, she might specify that it needs to reach  
66 its destination as quickly as possible whilst obeying the traffic laws, by assigning benefits and costs for  
67 desired and undesired behaviours. However, in natural environments, no external entity exists that can  
68 directly quantify the consequences of each action, like the points that are awarded in a video game for  
69 completing levels or shooting monsters (**Fig. 1b**). Nor is it obvious that biological systems have a  
70 dedicated channel for receipt of external rewards that is distinct from the classical senses. Rather,  
71 rewards and punishments *are* sensory observations – the taste of an apple, the warmth of an embrace  
72 – and so stimulus value must be inferred by the agent, not conferred by the world. In other words,  
73 rewards must be intrinsic, not extrinsic (**Box 1**). We call this “the reward paradox”.

74



75  
76

77 **Figure 1. The RL framework** **A.** Top: a typical “bandit” task in which participants choose between options (A and  
78 B) with uncertain reward probability, incurring financial outcomes. Below: in the RL framework, an agent receives  
79 observations (green), takes an action (blue) and receives feedback in the form of reward (red). The numbers  
80 highlight some assumptions of the basic RL framework that are less well suited to explaining motivated behaviour  
81 in biological organisms (grey box below). **B.** The RL framework applied to video games (here, Ms. Pac-Man). The  
82 agent acts in order to maximise game score, which is a scalar quantity. Unlike humans playing the game, the agent  
83 receives the score differential directly (e.g. points for eating pellets) as reward signal via a dedicated input  
84 channel. **C.** Schematic of model-based RL, where the agent searches through a tree of possible future states and  
85 actions, sometimes stopping if no reward is expected from this branch of the search tree, and ultimately choosing  
86 the trajectory through the tree that leads to the highest return.

87

88 In this *Opinion*, we consider the challenges of applying the reward hypothesis to the study of biological  
89 agents in the natural world, with a focus on the neural and computational mechanisms by which  
90 humans and other animals make value-guided choices. In doing so, we highlight some ways in which  
91 the RL framework might be modified, adapted or nuanced to account for the reward paradox. We  
92 appeal to a diversity of theories concerning motivated behaviour and ask how they can explain the  
93 behaviour and brain activity observed during human decisions in cognitive, social and economic  
94 settings.

**Box 1. Intrinsic Motivation.** Machine learning and AI researchers know that designing a reward function by hand is costly and does not scale to real world settings. To avoid this problem, agents have been built that generate their own intrinsic motivational signals [37, 57, 58]. For these agents, the cost function is populated by additional terms that encourage the agent to pursue adaptive behaviours even in the absence of reward. Many of these additional terms draw inspiration from psychological theory, which has long noted that humans and other animals often pursue activities for their own sake (e.g. a kitten furiously chasing a piece of string, or a commuter completing the crossword on her journey to work). For example, one key intrinsic motivation may be to explore the world in a structured fashion, satisfying curiosity and avoiding boredom [59]. This allows agents to seek out experiences that may accelerate learning without directly enhancing reward, a phenomenon known as active learning. Other intrinsic motivation signals may encourage us to make actions with predictable consequences, or that exert maximal control over the environment or the behaviour of conspecifics, or those that minimise discrepant or dissonant inputs (see [57] for a review). Augmenting the reward function with these intrinsic motivational signals often leads to faster learning and more stable patterns of behaviour in artificial agents [60]. However, in most reports, intrinsic motivation signals supplement (rather than replace) rewards from the external environment, leaving the thorny question of where rewards come from only partially answered. Other approaches that avoid any external feedback by using purely unsupervised or self-supervised methods – where the agent is optimised to maximise information gain or reduce sensory surprise – hold undoubted promise but have not yet been observed to yield complex behaviours in dynamic environments that resemble the natural world [61].

96

97 **Context-dependent valuation and the representation of internal state**

98

99 To circumvent the reward paradox, it is sometimes assumed that the reward function has been  
 100 designed by evolution. In other words, reward computation may rely exclusively on fast, mandatory,  
 101 phylogenetically prespecified mechanisms that automatically convert sensory observations into  
 102 hedonic signals that act as a proxy for “external” rewards. Indeed, there is no doubt that dedicated  
 103 mechanisms exist for those physiological drives that lie at the bottom of our “hierarchy of needs” [7],  
 104 such as hunger or fatigue. These needs are signalled by projections from the hypothalamus to neural  
 105 systems that respond to rewards, including dopaminergic neurons in the ventral tegmental area [8].  
 106 On this basis, proponents of RL models have dubbed dopamine neurons the “reward retina” as if they  
 107 were responsible for directly sensing reinforcement from the environment [9].

108

109 How, then, is the reward value of a stimulus is computed based on the agent’s internal state as well as  
 110 the context provided by the external environment? Intuitively, the value a potentially rewarding  
 111 stimulus (e.g. a hamburger) will depend on the agent’s ongoing bodily needs (am I hungry?). Under the  
 112 RL framework, this sensitivity to internal context might be accounted for by assuming that bodily states  
 113 part of the environment, i.e. that the state of “I am hungry” is external to the agent in just the same  
 114 way that the presence or absence of an apple is a property of the world, not just the mind. However,  
 115 this argument becomes harder to defend when one considers the range of affective and cognitive  
 116 factors that can potentially modulate value. For example, a hamburger might be less valuable to  
 117 someone who is nervous before an exam, or to a committed vegetarian. It seems harder to argue that  
 118 private mental phenomena such as moods and or moral attitudes are properties of the external  
 119 environment.

120

121 More generally, standard RL models assume that reward is a purely hedonic signal that is used solely  
 122 to update the agent’s value function. Reward is thus not *fungible* – that is, once acquired, it is not  
 123 exchangeable for other assets, and cannot be used in a way that would alter future observations or  
 124 facilitate future actions. For example, when RL agents are learning to play video games, the rewards  
 125 they receive (i.e. game score differential) evaporate instantly without providing future opportunities  
 126 for disbursement within the game. This is unlike natural settings, where rewards are typically

127 accumulated because they lead to persistent changes in state. For example, bodies act as repositories  
128 for energetic resources that are harvested by the agent during trophic behaviours, and bank accounts  
129 are used to store for income acquired in exchange for work. The agent can thus choose to exchange  
130 one stored asset for another that is needed, the defining feature of economic behaviour in the  
131 marketplace. In the laboratory, rewards may be implicitly accumulated as experimental animals or  
132 human participants become more sated or wealthy as a result of the liquid rewards or financial  
133 incentives that they receive for performing the task, but these changes in state are typically considered  
134 to be nuisance variables by the experimenter. For example, when a monkey working for liquid reward  
135 is no longer thirsty enough to work, the researcher will typically terminate the experiment. To fully  
136 account for biological behaviour, our models of value learning need to account for the internal state of  
137 the organism.

138  
139 Psychologists and neuroscientists who study appetitive responses to primary reinforcers such as food  
140 and liquid have long appreciated that theories motivated behaviour need to include a representation  
141 of internal state. One prominent model, known as incentive salience theory [10], proposes that whilst  
142 an animal might “like” a stimulus on the basis of cached state values (e.g. sugar tastes good) as  
143 proposed in standard RL, the extent to which it “wants” a stimulus also depends on its physiological  
144 state (e.g. hunger), which adjusts valuation via a gain control mechanism. This ensures that stimuli that  
145 satisfy natural appetites will be consumed more readily (e.g. food when hungry and water when thirsty).  
146 Incentive salience can account for a wide range of empirical phenomena. For example, the rate of  
147 approach towards food depends on hunger levels, and adapts immediately according to internal levels  
148 of satiety irrespective of prior training [11]. Similarly, food items that have been administered during  
149 states of food deprivation are preferred over control items during a later test when the animal is sated,  
150 suggesting they have acquired greater value [12]. However, incentive salience does not describe exactly  
151 how internal state is represented or specify how the animal chooses which appetite it wishes to satisfy.  
152 More widely, a rich literature has studied the consequences of satiety on food choices, often invoking  
153 model-based RL to describe why animals take actions that elicit a reinforcer that has not been devalued  
154 over one that has. However, this literature has largely been silent about the mechanisms by which  
155 devaluation itself occurs, focussing instead on how animals learn the causal structure of the task [13].

156  
157 In the remainder of this article, we explore the relevance of these considerations for the study of human  
158 decision-making. Humans routinely make decisions about complex stimuli or abstract plans of action in  
159 cognitive, social and economic settings. Our argument is that a representation of internal state may be  
160 computationally relevant even for decisions between consumer products, social opportunities, or  
161 lifestyle alternatives. We adapt models that have been previously proposed to account for choices  
162 among basic needs to these more intricate decision settings and describe how they might account for  
163 some intriguing findings concerning the neurobiology of human choice.

## 164 165 **Goals as cognitive setpoints**

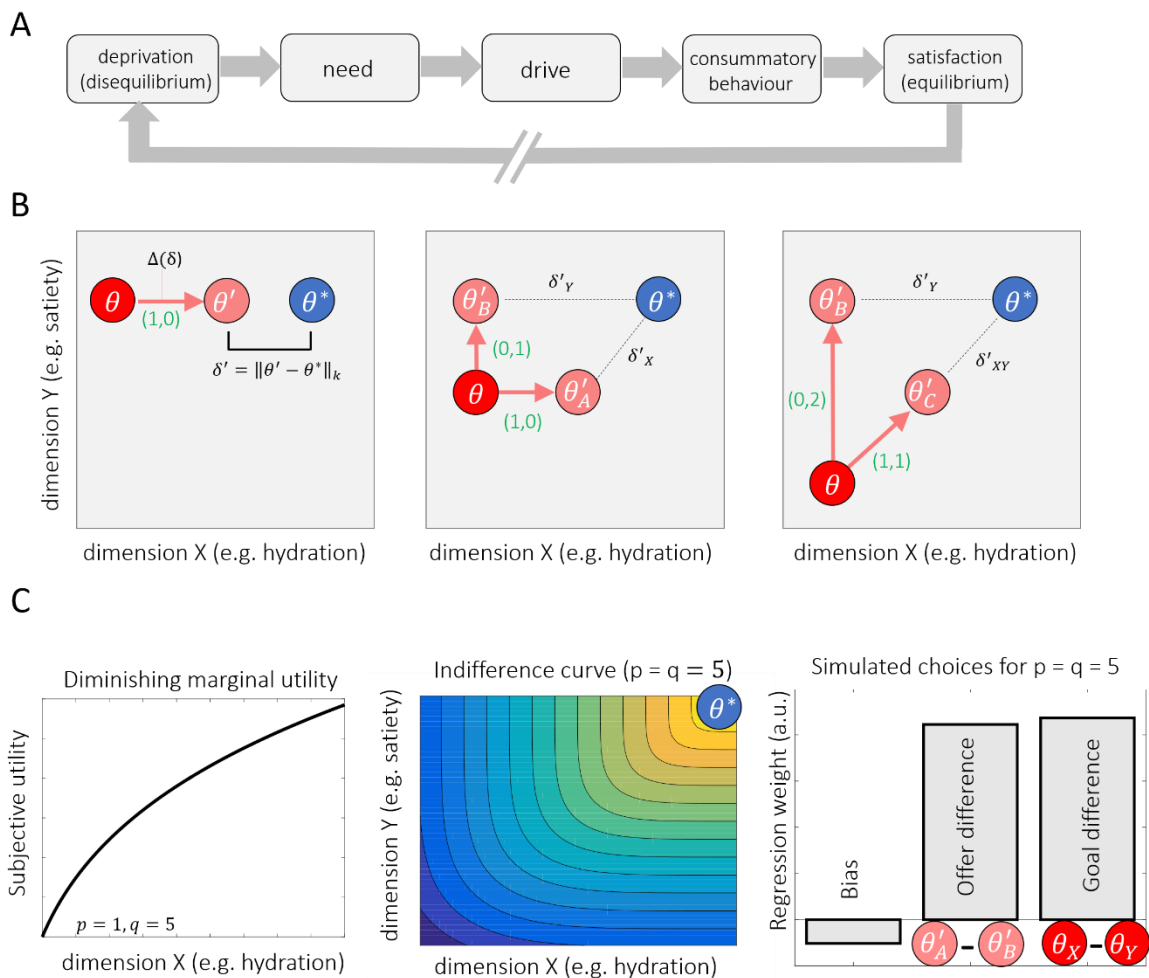
166  
167 Our argument begins with the observation that much human motivated behaviour is intrinsically goal-  
168 directed. Humans will often work tirelessly to complete tasks whose only apparent reward is the  
169 satisfaction of task completion itself, for example when seeking to win a race, solve a puzzle, or climb  
170 a mountain [14]. Indeed, this purposive, self-consistent property of human behaviour is so powerful  
171 that humans will seek to complete tasks even when the outcome is known to be a net loss, a  
172 phenomenon known as the sunk cost fallacy [15]. Although the RL framework embraces goal-directed  
173 behaviours, it does so under the assumption that animals explicitly learn the transition probabilities  
174 between states and use mental tree search to identify the most rewarding trajectories, which then  
175 guide “model-based” decisions (see **Fig. 1c**). Thus, for RL, goals are not the objects of computation  
176 themselves but emerge as a by-product of policies which optimise expected future reward [16]. Other  
177 elaborations of the RL framework explicitly propose additions to the cost function, such as a desire for

178 novelty, agency or consistency in behaviour, which draw upon the assumption from cognitive  
 179 psychology that humans do more than simply maximise reward (see **Box 1**).

180

181 In considering these phenomena, we turn to one of the earliest theories of motivated behaviour – one  
 182 which assumes that animals seek to keep internal drives in equilibrium [17, 18]. This theory builds on  
 183 early models of homeostatic function [19] and found its first prominent expression in the 1940s with  
 184 Hull’s Drive Reduction Theory [20]. These homeostatic models argue that motivated behaviours seek  
 185 to restore imbalance among different physiological needs, such as warmth, satiety, and hydration (**Fig.**  
 186 **2a**). A central idea is that neural circuits encoded desired states or “setpoints” and drives result from  
 187 the disparity between current and desired state. Thus, rewards are multidimensional quantities  
 188 signalling the drive reduction that occurs when an imbalance among the organism’s multiple internal  
 189 needs is redressed.

190



191

192 **Figure 2. Theories of Motivation.** **A.** Schematic of Drive Reduction Theory. Behaviours are initiated to reduce  
 193 drives, which in turn are provoked by needs that arise following homeostatic imbalance. **B.** Illustration of  
 194 homeostatically regulated reinforcement learning (HRRL) [18], here in two example dimensions (labelled satiety  
 195 and hydration). The agent’s current state is denoted  $\theta$ , candidate next states  $\theta'$  and the desired state (setpoint)  
 196  $\theta^*$ . Offers in the (x,y) dimension are shown in green. In the leftmost panel, the agent is sated but thirsty, as in a  
 197 laboratory experiment where a monkey is working for liquid. An action reduces the distance  $\delta$  between the  
 198 current and desired state  $\|\theta - \theta^*\|_k$  (where  $k$  indicates the Minkowski distance; we use this here for simplicity,  
 199 but note that a more general, two-parameter version depicted in panel **C** is needed to account for non-linearities  
 200 in valuation, see **Box 2** for details). Notice that the rate of reduction in distance is proportional to what would (in  
 201 the RL framework) be construed as an external reward signal. In the middle panel, an agent who is more thirsty

202 than hungry is offered an equivalent quantity of food (0,1) or water (1,0). Water reduces distance to goal  
203 (Euclidean,  $k = 2$ ) more than food and is thus preferred. In the rightmost panel, the agent is both very hungry  
204 and very thirsty. An offer comprising a bundle of both food and water (1,1) is preferred over a large quantity of  
205 food (0,2). **C.** Predictions of HRRL with two parameters,  $p$  and  $q$ , governing non-linearities. The model can  
206 reproduce the ubiquitous finding of diminishing marginal utility (left panel). As  $p$  and  $q$  increase, distance  
207 estimates depend more strongly on just one dimension, indicating that maximum weight is given to the currently  
208 most needed asset (middle panel). Intermediate values  $p$  and  $q$  can reproduce one of the central findings of  
209 Juechems et al., 2019: that choices are both driven by offers and needs when assets need to be kept in equilibrium  
210 (see also Fig. 4).

211  
212 Here, we make a connection between homeostatic models of reward, which have emphasised balance  
213 among basic drives and interoceptive processes, and recent work on value-guided decision-making in  
214 humans and other primates [21-23]. We argue that motivated behaviour in cognitive, economic or  
215 social settings can similarly be understood as seeking to restore a balance among competing “setpoints”  
216 that correspond to higher-order goals. Unlike currently popular theories of value-guided choice, which  
217 are largely inspired by RL models [2] or expected utility theory from economics [24], our view  
218 emphasises the multidimensional nature of value signals, the dependence of motivated behaviour on  
219 the internal state of the agent, and the desire to explicitly realise goal states (constraint satisfaction)  
220 rather than maximise reward. In what follows, we show how extant theories of motivation can explain  
221 some otherwise puzzling behavioural and neural findings in the literature on value-guided decision-  
222 making and economic choice.

223

## 224 Homeostatically regulated reinforcement learning for cognitive behaviour

225

226 We suggest that during learning, humans form new setpoints pertaining to cognitive goals. For  
227 example, we might represent current and desired states on axes pertaining to financial stability, moral  
228 worth, or physical health as well as hunger, thirst or temperature. To sketch the computational  
229 underpinnings of this process, we appeal to an existing framework for understanding motivated  
230 behaviour for basic physiological processes [18, 25, 26], known as homeostatically regulated  
231 reinforcement learning (HRRL; **Box 2**). This theory proposes that current states and goals are encoded  
232 in a multidimensional “value map”. Motivated behaviour can then be seen as an attempt to minimise  
233 the maximum distance to setpoints in this value space. Repurposing this framework for cognitive  
234 settings, agents commit to policies which focus on purposively driving the current state towards  
235 setpoints on a particular goal dimension, such as caching resources, building a shelter, obtaining a mate,  
236 or enhancing professional status. In doing so, their ultimate goal is to maintain equilibrium among all  
237 goal states, achieving what might be popularly characterised as a state of “wellbeing”.

238

239 In HRRL, reward emerges naturally as a consequence of the computations required for learning, rather  
240 than being furnished by the external environment. This allows the “reward paradox” highlighted above  
241 to be sidestepped. To illustrate, consider an animal whose internal state is currently described by  
242 parameters  $\theta$  and who seeks to achieve a desired state  $\theta^*$ . Let us assume that the agent has a  
243 biologically plausible functional form, such as a neural network. In order to calculate the gradients that  
244 will allow network weights to be optimised for future behaviour, it is necessary to compute the loss  
245 term  $\Delta(\|\theta^* - \theta\|_k)$ , which indicates the change in distance between a current and desired state that  
246 is afforded by any action. Now consider a typical laboratory study in which a single reward dimension  
247 is relevant, for example because the experimenter reduces prior liquid intake and offers only water in  
248 exchange for behaviour. For a thirsty animal, this loss term is now mathematically identical to what one  
249 would construe as “reward” in the RL framework – i.e. the volume of liquid administered on a trial [18].  
250 This equivalence will break down as  $\theta$  tends towards  $\theta^*$ , but by this point the researcher will typically  
251 have ended the experiment because the animal is fully hydrated. Thus, although in this experiment  
252 value is artificially constrained to be unidimensional, we note that from the agent’s perspective,

253 multiple value dimensions may continue to be relevant. For example, a human participant might be  
 254 tempted to neglect the task in favour of checking their social media account, in order to enhance social  
 255 capital. However, this individual would most likely be excluded from the analysis.  
 256 Homeostatically regulated reinforcement learning readily incorporates other accounts of reference-  
 257 dependent evaluation by adjusting how the agent computes  $\theta^*$ . For instance, it may do so by using a  
 258 running average over past options and anchor its behaviour to this reference. When  $k > 1$ , simply  
 259 assuming that other drives exist that cannot currently be satisfied immediately gives rise to the  
 260 ubiquitous phenomenon of diminishing marginal utility. In other words, as the agent acquires the

### Box 2: Homeostatically Regulated Reinforcement Learning

Let us assume that sensory observations,  $\mathbf{X}_t$  are mapped onto a persistent internal state  $\theta_t$ :

$$\theta_t = \mathbf{X}_t \times \mathbf{U} + \theta_{t-1} \times \mathbf{A}_\theta \quad \text{Eq.1}$$

Where  $\mathbf{U}$  is a matrix mapping the chosen option  $\mathbf{X}_t$  onto  $\theta_t$ . The internal state,  $\theta_t$ , exhibits dynamics whereby it will change (decay) in the absence of any input, e.g. credit card debt will increase if it is not repaid, just as thirst will increase if not slaked. These dynamics are captured by the matrix  $\mathbf{A}_\theta$ , which governs the intrinsic changes over time.

Note that this implies that the correspondence between internal and external states must be learned, i.e. our preferences are not simply given to us but acquired with experience, and that we might be uncertain about whether we are hungry, tired or happy. This relates to the framework of “active inference”, under which reward values are inferred rather than being directly conferred [17]. We conceive of this internal state space as a multi-dimensional map whose dimensions variously pertain to both physiological states (e.g. satiety) and cognitive goals (e.g. social status). Internal states are subject to momentum via recurrent connections so that states (or moods) may persist in the face of a change in external circumstance.

In HRRL [18], the agent represents a set of desired (or aversive) points  $\theta^*$  on this space, allowing it to estimate the distance (and contributing features) of its current position to the goal state. In doing so, the agent needs to combine across the relevant goal dimensions and weight these according to their relative contribution to this distance by computing the Minkowski distance:

$$\delta_{t+1} = \|\theta_{t+1} - \theta^*\|_k = \sqrt[p]{\sum |\theta_{t+1} - \theta^*|^q}, \text{ for } p = q \quad \text{Eq.2}$$

The parameter  $k$  allows the agent to prefer bundles that only consist of one single asset ( $k < 1$ ), be indifferent between all bundles ( $k = 1$ ) or prefer mixed bundles ( $k > 1$ ). The model is thus well placed to account for commonly observed indifference curves between two or more economic goods. If we further assume that there are two non-linearities,  $p$  and  $q$ , rather than just  $k$ , our model can further account for diminishing marginal utility, **Fig. 2c**. The implications of this computational step are described more fully in ref [18].

Finally, where, then, is value (or: reward) in this framework? Value arises naturally when the agent compares its current distance to its desired state to its previous distance, which summarizes that the agent has progressed on at least one of its goals since the last time-point:

$$R_{t+1} = \delta_t - \delta_{t+1} \quad \text{Eq.3}$$

Value is thus a summary of whether the agent is approaching or retreating from its goals. This model encapsulates situations in which multi-attribute objects need to be integrated into a single reward signal, such that  $R = \mathbf{X} \times \mathbf{U}$  if  $\mathbf{A}_\theta = \mathbf{0}$  and  $k = 1$ .

261 currently relevant asset (such as liquid reward) its attention will progressively tend towards other  
262 relevant assets, diminishing the contribution of liquid reward to its well-being (Fig. 2c). However,  
263 allowing  $k$  to vary freely, a range of different policies can be accommodated, including those in which  
264 one goal (e.g. attainment of a desired drug in addiction) dominates all others [27].

265

### 266 **Implications of HRRL for cognitive decisions**

267

268 More generally, we argue that some of the most complex and abstract decisions that humans make  
269 might be better described by a process that optimises over states, rather than rewards. For example,  
270 consider a high school student choosing a career path. Under the (model-based) RL framework, she  
271 must consider an impossibly large number of potential futures and select whichever is going to be most  
272 rewarding. This seems to imply the devotion of disproportionate levels of computational resources to  
273 the search problem. The approach advocated here implies that she first selects a goal state (e.g.  
274 become a lawyer) than then takes actions which minimise distance to that goal. For example, she seeks  
275 to go to law school; to maximise her chances of acceptance, she first studies hard for her exams; this  
276 in turn influences decisions about whether to socialise with friends. This explanation seems to accord  
277 better with our common sense intuition of how the complex choices faced by humans are made.  
278 However, the computations involved may build upon more phylogenetically ancient mechanisms. For  
279 example, one of the most prominent theories of insect navigation proposes that in order to reach their  
280 home base, central place foragers such as honey bees (and desert ants) initially encode an egocentric  
281 snapshot of their base and subsequently, on the return journey, use a similarity-matching process to  
282 gradually reach their goal [28]. This implies they are similarly performing gradient descent over states,  
283 akin to the process proposed here.

284

285 Our view relates to an earlier proposal that animals are motivated to satisfy a “current concern” [29]  
286 and understands disorders of mental health as disrupted goal selection, pursuit and appraisal [30]. A  
287 fleshed out computational theory of this process will require the specification of how agents learn the  
288 relationship between individual states and the degree of goal realisation. One model proposes a  
289 distinction between “primary” and “learned” value, where the former resembles the traditional reward  
290 signal in the RL framework and the latter quantifies the eligibility of states for goal realisation [14]. A  
291 very appealing aspect of this framework is that it provides a natural way to understand the affective  
292 states that pervade our everyday mental landscape, including satisfaction (goal completion), frustration  
293 (goal obstruction), and disappointment (goal abandonment), which have largely eluded computational  
294 description thus far [14]. A more directly related machine learning approach directly computes the  
295 feasibility of attainment of goal states, showing how it can be used for compositional inference under  
296 the nonstationary reward functions characteristic of real world environments [31].

297

298 Many of the more complex behaviours exhibited by humans involve decision about money. Money has  
299 the peculiar property of being fluidly exchangeable for other assets (except, as the Beatles remind us,  
300 for love). Our here might seem to break down in financial settings, because it seems incongruous to  
301 propose a “setpoint” for wealth – anecdotally at least, the accumulation of wealth does not always  
302 blunt the desire to acquire new wealth. However, there is evidence that even wealth accrual may be  
303 subject to some of the constraints identified here. For example, taxi drivers on shifts of variable length  
304 will tend to work until a fixed return has been achieved, returning home even when this “wealth  
305 setpoint” falls during a period of high yield – for example, even if customers are plentiful and fare rates  
306 are high [32]. It is also observed that humans use mental accounting to allocate wealth to different  
307 potential value dimensions in advance [33]. Thus, although money can be deployed fluidly to facilitate  
308 movement in any direction in value space, humans mentally pre-allocate funds to the different  
309 directions of advancement, treating an unexpected excess in one direction with profligacy rather than  
310 reallocating it facilitate progress along other dimensions. This latter phenomenon is known as the  
311 “house money” effect.

312

313 Pure drive reduction theories fell from favour with the suggestion that animals will work for rewarding  
314 stimuli when their drives are satisfied, i.e. purely for their hedonic value [34, 35]. For example, animals  
315 will continue to eat after food has been administered intra-gastrically [34], and will consume stimuli  
316 with no nutritional value, such as saccharine water, when they are hungry [36]. Although alternative  
317 explanations for these phenomena exist [19], we caveat our argument with the acknowledgement that  
318 some behaviours are driven by “liking” alone, irrespective of any relation to the current setpoint.  
319 Similarly, it is unlikely that all human cognitive decisions are based on minimising distance to a desired  
320 goal – the mechanisms proposed here may exist alongside other processes that allow behaviours that  
321 have been pleasurable or successful in the past to be repeated in the future. Thus, we might buy a pair  
322 of shoes on impulse, sneak an extra bite of dessert even when we are full, or linger at a party even after  
323 our social duties have been fulfilled. However, our argument is that the human tendency to optimise  
324 towards desired states has been overlooked in the face of enthusiasm about rewards as the primary  
325 drivers of behaviour.

326

### 327 A map of internal state in the OFC

328

329 In the final sections, we consider the implications of the perspective advanced here for the  
330 neurobiology of decision-making. For valuation to depend on the internal state of the agent, these  
331 states must be explicitly represented in brain signals. However, most neural studies of value-guided  
332 choice have used a paradigm in which the animal or human participant makes a succession of unrelated,  
333 independent choices between food items or monetary gambles. This paradigm is not well suited to  
334 measuring how variation in internal state (how satiated am I right now? How much wealth have I  
335 accumulated?) might be coded in neural circuits. However, recently researchers have begun to employ  
336 more complex paradigms that may shed some light on this issue.

337

338 At the neural level, the integrity of the orbitofrontal cortex (OFC) is critical for goal-driven reward  
339 behaviour [13]. Interestingly, a recent theory proposes that the OFC constitutes a “map” of both  
340 observable and latent states within a task [37, 38], similar to the framework proposed here (c.f. Eq.1 in  
341 Box 2). Moreover, OFC may code endogenous (internal) as well as exogenous (external) value, as firing  
342 rates [39] and BOLD signals [40] can be choice-predictive even in the baseline period prior to stimulus  
343 onset. An emerging theory suggests that OFC encodes affective states that pertain to ongoing positive  
344 or negative value, i.e. moods or “momentum” in internal value states [41, 42]. Together, these findings  
345 imply that OFC is a likely candidate for representing ongoing estimates of position on the value map.

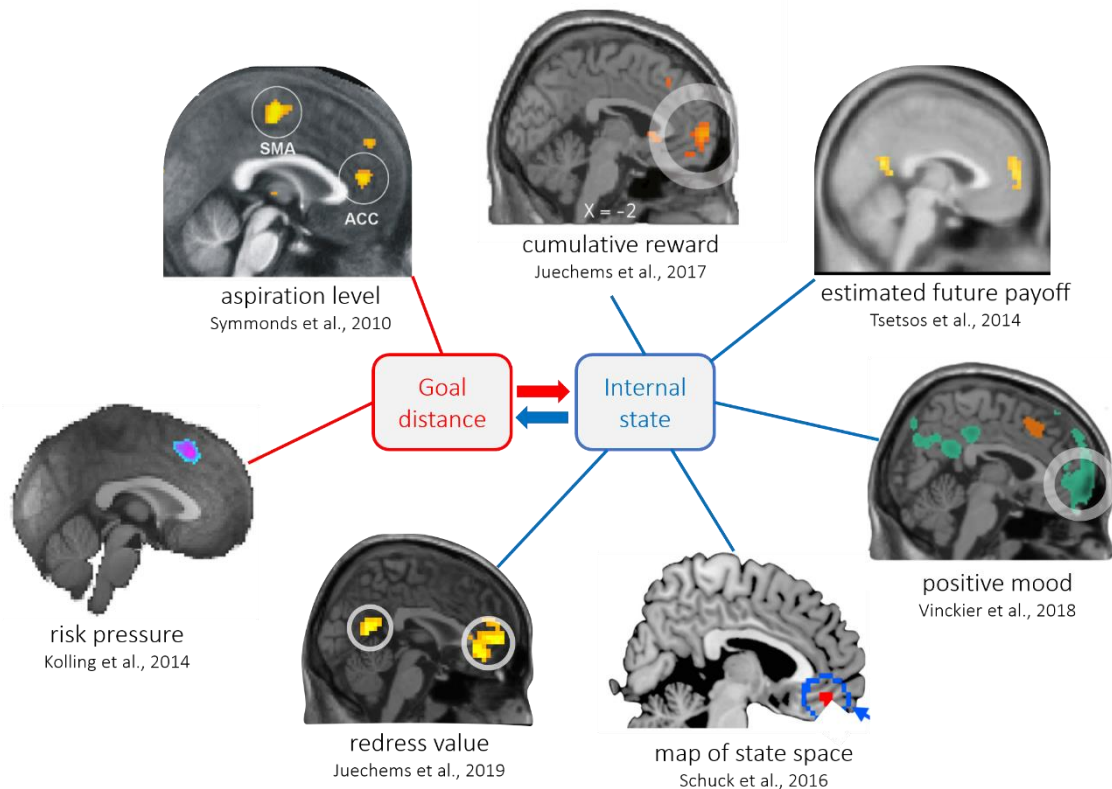
346

347 One recent study directly tested whether BOLD signals in OFC tracked the internal accumulation of  
348 financial resources across an episode in a way consistent with the ongoing internal state representation  
349 proposed here [43]. Participants performed an economic task that involved gambling for financial  
350 incentives over successive, independent trials. For a rational agent performing this task, there is no  
351 imperative to track the accumulated level of payoff incurred by choices – to maximise return, one  
352 simply needs to select, on each trial, the gamble with the highest expected value. However, when  
353 humans performed this gambling task over multiple discrete episodes, each signalled by a prominent  
354 contextual cue that was structurally irrelevant to the task, brain activity in the medial OFC tracked the  
355 level of accumulated resources (or “wealth”) over the context (**Fig. 3**). This “value accumulation” signal  
356 occurred even though the accumulated rewards were never overtly signalled to participants, nor were  
357 they instructed or incentivised to calculate their ongoing wealth.

358

359 Another brain imaging study revealed that a nearby region of medial prefrontal cortex tracks internal  
360 states during wealth accumulation [44]. Participants were asked to choose among four alternatives  
361 (bandits) that either paid or charged an uncertain sum of money. Critically, each selected bandit was  
362 fungible – it continued to incur financial costs or benefits across a short block of trials, as if the agent

363 had purchased that asset and was continually profiting (or otherwise) from its possession (“rule in”  
 364 condition). In a different condition, the bandits began in the agent’s possession and decisions had to  
 365 be made about which (if any) bandits to “rule out”. This allowed us to assess how neural encoding of  
 366 value depended on whether an asset was being bought or sold. Critically, BOLD signals in the  
 367 ventromedial prefrontal cortex (vmPFC) encoded the cumulative expected value of assets across the  
 368 block (Fig. 3), and only coded for momentary rewards when decisions were made with future  
 369 consequences for reward (e.g. when a bandit was selected in the “rule in” condition, or not selected in  
 370 the “rule out” condition). Together, these studies suggest that the reward system tracks level of  
 371 accumulated resources over time, for example so that future decision policies can adapt according to  
 372 the current internal state, in a way that is compatible with the theory advanced here.  
 373



374 **Figure 3. Neural coding of internal state and goal distance.** Cortical areas on the medial surface correlating with  
 375 diverse aspects of internal state and goal distance. Reprinted with permission from refs [37, 42-44, 46, 47, 55].  
 376

377  
 378 **Budget rules: aspiration and avoidance points**

379  
 380 The natural world is structured in such a way that some states are critical for survival or have substantial  
 381 impact on long-run future outcomes. For example, the student introduced above might work hard to  
 382 pass her exams in the knowledge that it will open up interesting career opportunities. These states are  
 383 often attained when accumulated resources reach, or fall below, a critical threshold. Behavioural  
 384 ecologists have argued that animals’ risky foraging behaviour adapts to satisfy a “budget rule” that  
 385 seeks to maintain energetic resources at aspirational levels that safely offset future scarcity. For  
 386 example, birds make risky foraging choices at dusk in order to accrue sufficient energy to survive a cold  
 387 night [45]. This view is neatly accommodated within the framework proposed here, in that the  
 388 aspiration level reflects the setpoint against which current resource levels are compared, and the driver  
 389 of behaviour is the disparity between current state and goal.  
 390

391 Recent work with brain imaging [46, 47] suggests that this goal distance signal might be computed in a  
392 different medial prefrontal cortex region, the dorsal anterior cingulate cortex (dACC). Symmonds [47]  
393 asked participants to make a series of gambles which only paid out if a critical threshold or aspiration  
394 point was reached. BOLD signals in the vmPFC scaled with the expected future return over a series of  
395 gambles, similar to [44], whereas dACC coded for the aspiration point itself (**Fig. 3**). Using a very similar  
396 design, Kolling [46] asked human participants to accept or reject monetary gambles over a short block,  
397 with cumulative earnings increased by a multiplier if they reached a fixed aspiration level. BOLD signals  
398 in the dACC coded for number of steps to goal, as well as the discrepancy between accumulated return  
399 and the aspiration level, scaled by the distance to the end of the block. This signal, which the authors  
400 call “risk pressure”, seems to be signalling the agent’s current position with respect to a goal, rather  
401 than merely the arrival of a punctate reward. In another study involving planning in a virtual subway  
402 environment, dACC BOLD signals coded for the distance to goal in number of context switches [48].  
403 Indeed, the dACC is one brain region where single-cell responses gradually build up in association with  
404 a series of actions that precede reward delivery [49], and where lesions provoke failures of task  
405 persistence [50], as if this region participated more generally computing the disparity signal  $\theta^* - \theta$  in  
406 the service of goal-directed behaviour.

407  
408

### 409 **The multidimensional nature of value**

410

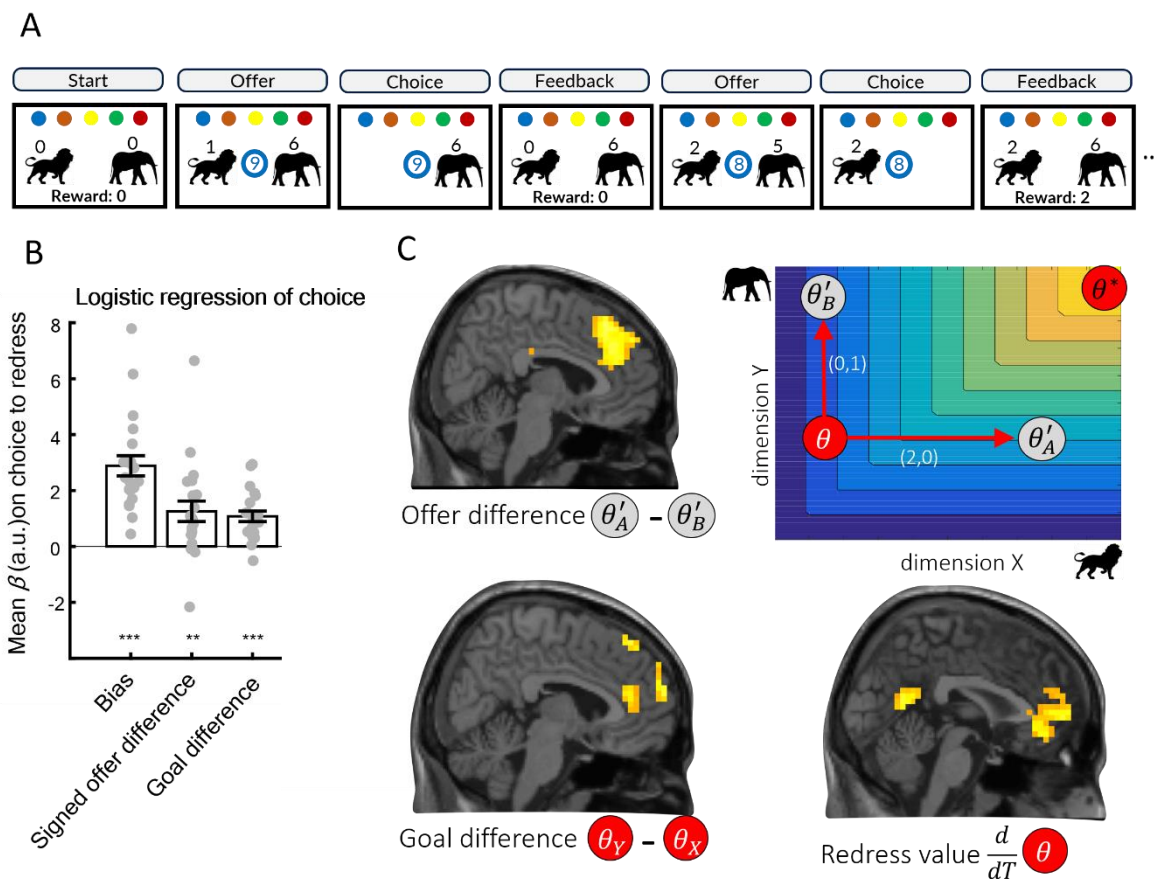
411 In the standard RL framework, each state-action pair is associated with a scalar value, often called a Q-  
412 value [1]. In other words, most RL models assume that value is a unidimensional quantity (although  
413 “multi-objective” RL models relax this assumption [51]). This seems incongruous when considering the  
414 challenges that animals face in natural environments, which involve jointly satisfying many different  
415 constraints (for example, maintaining levels of satiety, hydration, warmth and social capital). Rewards  
416 take many different forms – from food items to social and sexual signals to secondary reinforcers such  
417 as money. A popular view in the hybrid field of neuroeconomics is that neural signals allow potentially  
418 incommensurable rewards to be encoded in a single neural “common currency” [23]. For example,  
419 single cells in OFC code for the “offer value” of different liquid rewards irrespective of their taste [52].  
420 The existence of a common neural code for different assets allows economists to understand rational  
421 decision-making in consumer settings, which might involve choosing between a meal in a restaurant or  
422 tickets to the theatre. However, HRRL suggests a different perspective, whereby assets are explicitly  
423 encoded along different dimensions (needs) in the value map, and decisions are made according to  
424 whichever need is currently greatest.

425

426 Without further assumptions, HRRL predicts that animals will treat different assets as if they are  
427 imperfect substitutions for one another. Multidimensional prospects (or “bundles”) should be  
428 preferred if they contain a mixture of different assets, because this allows multiple drives to be satisfied  
429 simultaneously (**Fig. 2b**). For example, it is better to be neither too hungry nor too thirsty than to be  
430 starving or dehydrated. Several longstanding empirical results support this view. First of all, in early  
431 conditioning experiments it was observed that a cue that is associated either simultaneously or  
432 successively with multiple assets – such as food and water - has greater value than an equivalent cue  
433 associated with only a single asset (or drive) [53]. Finally, a preference for compromise decisions is also  
434 apparent from studies of so-called “decoy” effects in multi-attribute, multi-alternative choices. Faced  
435 with a choice between two iso-preferred items that have complementary costs and benefits (e.g. an  
436 expensive but high quality consumer product A, or a cheaper but lower quality product B, the  
437 introduction of a “decoy” that is higher in price/quality than A will bias preferences towards A, whereas  
438 an item that is lower in price/quality than B will bias choices towards B, as if participants had an intrinsic  
439 preference for the compromise solution [54]. This preference is predicted by a need to keep goals in  
440 equilibrium, by assuming that agents aspire to high quality and low price simultaneously, rendering the  
441 agent’s indifference curve over options convex to the origin (**Fig. 2b**).

442  
 443  
 444  
 445  
 446  
 447  
 448  
 449  
 450  
 451  
 452  
 453  
 454  
 455  
 456  
 457  
 458  
 459  
 460  
 461

In natural environments we are faced with choices between offers of different quantity and quality, in variable states of need. In such scenarios, our proposal predicts that both goal difference (the relative level of two internal resources A and B) and offer difference (the relative quality of offer of A and B) will predict choices, **Fig. 2c**. This prediction was tested in a recent study in which participants made decisions that directly trade off multiple assets. Juechems [55] asked participants to perform a task that involved managing a virtual zoo containing lions and elephants. On each trial they populated the zoo by choosing among offers of variable numbers of each animal, and the financial success of their zoo was proportional to the minimum number of either lions or elephants present on each trial (which was not directly observed), **Fig. 4a**. Each choice was classified according to whether it “redressed” the imbalance in cumulative lions and elephants and decisions to “redress” were predicted by both goal difference and offer difference, **Fig. 4b**. Critically, BOLD signals in the vmPFC signalled the level of redress incurred by a choice – the extent to which a choice brought the numbers of animals into equilibrium, but independent of the total animal numbers and thus of any overall financial return. Interestingly, the dACC coded for the level of imbalance among assets, i.e. the pressure to redress this imbalance, reminiscent of the “risk pressure” signal described above [46], **Fig. 4c**. Whilst the conclusions that one can draw from this study may be limited by the fact that (by experimental necessity) lions and elephants were ultimately converted to a single currency (money), we think it provides preliminary evidence for the framework advanced here.



462  
 463  
 464  
 465  
 466  
 467  
 468  
 469

**Figure 4. Neural signals for computing goal equilibrium.** **A.** Experimental task used in Juechems et al., 2019. Participants managed virtual zoos in which they needed to accumulate as many lions and elephants as possible. Starting out from an empty zoo, participants were offered varying amounts of lions and elephants to acquire. Importantly, reward was calculated as the lower tally of animals currently owned within the zoo (e.g. when participants owned 2 lions and 6 elephants, they earned 2 units of reward). This enforced a rule by which the two tallies should be held approximately in equilibrium across time. Participants completed five zoos per block with their respective length indicated as a countdown in the centre of the screen. **B.** Logistic regression predicting a

470 “redress” choice – i.e. whether participants chose to increase the lower tally of animals currently in the zoo (as  
471 given by the feedback screen). Participants placed approximately equal weight on the offer and goal (tally)  
472 differences, as in the simulation in Fig. 2c. C. Indifference map over the two assets (lions and elephants) and three  
473 key neural correlates. Regions in the dACC and vmPFC encoded the offer difference, goal difference, and the  
474 redress (the decrease in distance to goal between trials). Reprinted with permission from ref [55].

475  
476

## 477 **Where do goals come from?**

478

479 Our focus here has been on “where value comes from”. Drawing upon homeostatic theories of  
480 motivation, we argue that the value of actions taken in cognitive, social and economic settings depends  
481 at least in part on the extent to which they drive an agent towards – or away from – a setpoint.  
482 However, this begs another, potentially more thorny question: how are the setpoints set in the first  
483 place? In those use cases to which HRRL has been applied thus far, pertaining to physiological needs,  
484 this is relatively uncontroversial, as it is safe to assume that setpoints for satiation and hydration are  
485 innate. What, then, of our desire to climb a mountain or become a lawyer?

486

487 In the RL framework, goals are chosen because of their potential reward value. But if reward is a  
488 reflection of the distance to a goal, then in our hands this argument becomes circular. One possibility  
489 is that goal-setting involves a hierarchical means-end reasoning process, by which proximal needs beget  
490 simple goals, the achievement of which requires more complex goals, and so on in a hierarchy of  
491 abstraction. In fact, the earliest theories of executive function and the prefrontal cortex proposed just  
492 such a mechanism for goal-directed behaviour [56]. Another factor that should not be overlooked is  
493 the role of social influence in goal setting. We often seek to achieve goals because they were suggested  
494 by our friends and associates, or by commercial advertising. However, we have to admit that we do not  
495 know how goals are computed. We think it is likely that the process by which we set ourselves goals is  
496 at least partly arbitrary, and otherwise intricately linked with the most intricate and ephemeral aspects  
497 of human cognition that pertain to our identity and sense of self.

498

## 499 **Concluding Remarks**

500

501 Our contribution here is twofold. Firstly, we highlight some challenges for RL as a computational theory  
502 of motivated behaviour in biological agents. The RL framework as developed within machine learning  
503 assumes that rewards originate in the environment, that they are sensed by dedicated input channels,  
504 that they are not fungible, and (typically) unidimensional. Secondly, we suggest that an alternative  
505 framework [14, 18, 20, 29] that treats reward as the by-product of computing distance to largely self-  
506 defined (*intrinsic*) goals, may have rich explanatory power for understanding human decisions. We  
507 show how this framework readily accounts for commonly observed effects in human and animal  
508 behaviour, such as reward-maximization, reference-dependence, diminishing marginal utility, and  
509 convex indifference between goods. Empirical tests of the theory’s wider predictions will require tasks  
510 (or field observations) that involve accumulation of multiple assets. Such tasks are rare in the literature  
511 concerning value-guided choices, and where multiple value dimensions inadvertently come into play,  
512 they are often even discarded as nuisance variables. We hope that the ideas advanced here will suggest  
513 ways in which to probe goal-directed behaviour using richer tasks which more closely emulate the fact  
514 that human (as well as animal) behaviour is driven by multiple, competing goals.

515

516

## 517 **Acknowledgements**

518

519 This work was funded by the European Research Council (grant REP-725937 to C.S.) and also received  
520 funding from the European Union’s Horizon 2020 Framework Programme for Research and Innovation

521 under Specific Grant Agreement 785907 (Human Brain Project SGA2, to C.S.). We are grateful to a large  
522 number of colleagues for interesting discussions on this topic and constructive feedback about a  
523 preprint version of the MS (and especially Tom Ringstrom); to Mehdi Keramati and two anonymous  
524 reviewers for helpful and constructive critique; and to Boris Gutkin, who thoughtfully provided an  
525 informal review of the revised submission.

526

527

528

529

530

531

532

533    **References**

- 534    1. Sutton, R. and Barto, A. (1998) Reinforcement Learning, MIT press.  
535    2. Dayan, P. and Daw, N.D. (2008) Decision theory, reinforcement learning, and the brain. *Cogn*  
536    *Affect Behav Neurosci* 8 (4), 429-53.  
537    3. Dolan, R.J. and Dayan, P. (2013) Goals and habits in the brain. *Neuron* 80 (2), 312-25.  
538    4. Schultz, W. et al. (1997) A neural substrate of prediction and reward. *Science* 275 (5306), 1593-9.  
539    5. Mnih, V. et al. (2015) Human-level control through deep reinforcement learning. *Nature* 518  
540    (7540), 529-33.  
541    6. Silver, D. et al. (2017) Mastering the game of Go without human knowledge. *Nature* 550 (7676),  
542    354-359.  
543    7. Maslow, A.H. (1943) A theory of human motivation. *Psychological Review* 50, 370–96.  
544    8. Fadel, J. and Deutch, A.Y. (2002) Anatomical substrates of orexin-dopamine interactions: lateral  
545    hypothalamic projections to the ventral tegmental area. *Neuroscience* 111 (2), 379-87.  
546    9. Schultz, W. (2015) Neuronal Reward and Decision Signals: From Theories to Data. *Physiol Rev* 95  
547    (3), 853-951.  
548    10. Berridge, K.C. (2012) From prediction error to incentive salience: mesolimbic computation of  
549    reward motivation. *Eur J Neurosci* 35 (7), 1124-43.  
550    11. Mollenauer, S.O. (1971) Shifts in deprivation level: different effects depending on the amount of  
551    preshift training. *Learning and Motivation* 2, 58-66.  
552    12. Pompilio, L. and Kacelnik, A. (2005) State-dependent learning and suboptimal choice: When  
553    starlings prefer long over short delays to food. *Animal Behaviour* 70, 571–578.  
554    13. Burke, K.A. et al. (2008) The role of the orbitofrontal cortex in the pursuit of happiness and more  
555    specific rewards. *Nature* 454 (7202), 340-4.  
556    14. O'Reilly, R.C. et al. (2014) Goal-Driven Cognition in the Brain: A Computational Framework. arXiv  
557    arXiv:1404.7591.  
558    15. Knox, R.E. and Inkster, J.A. (1968) Postdecision dissonance at post time. *Journal of Personality*  
559    and *Social Psychology* 8, 319–323.  
560    16. Daw, N.D. et al. (2005) Uncertainty-based competition between prefrontal and dorsolateral  
561    striatal systems for behavioral control. *Nat Neurosci* 8 (12), 1704-11.  
562    17. Pezzulo, G. et al. (2015) Active Inference, homeostatic regulation and adaptive behavioural  
563    control. *Prog Neurobiol* 134, 17-35.  
564    18. Keramati, M. and Gutkin, B. (2014) Homeostatic reinforcement learning for integrating reward  
565    collection and physiological stability. *Elife* 3.  
566    19. Berridge, K.C. (2004) Motivation concepts in behavioral neuroscience. *Physiol Behav* 81 (2), 179-  
567    209.  
568    20. Hull, C.L. (1943) *Principles of Behavior: An Introduction to Behavior Theory*, Appleton-Century-  
569    Croft.  
570    21. Rushworth, M.F. et al. (2011) Frontal cortex and reward-guided learning and decision-making.  
571    *Neuron* 70 (6), 1054-69.  
572    22. Rangel, A. et al. (2008) A framework for studying the neurobiology of value-based decision  
573    making. *Nat Rev Neurosci* 9 (7), 545-56.  
574    23. Kable, J.W. and Glimcher, P.W. (2009) The neurobiology of decision: consensus and controversy.  
575    *Neuron* 63 (6), 733-45.  
576    24. Glimcher, P.W. (2004) *Decision, Uncertainty and the Brain: : The Science of Neuroeconomics*,  
577    MIT Press.  
578    25. Keramati, M. and Gutkin, B., *A Reinforcement Learning Theory for Homeostatic Regulation*,  
579    *Advances in Neural Information Processing Systems*, 2011.  
580    26. Hulme, O. et al. (2019) Neurocomputational Theories of Homeostatic Control.  
581    <https://psyarxiv.com/s2q46/>.  
582    27. Keramati, M. et al. (2017) Cocaine addiction as a homeostatic reinforcement learning disorder.  
583    *Psychol Rev* 124 (2), 130-153.

584 28. Cartwright, B.A. and Collett, T.S. (1983) Landmark learning in bees. *Journal of comparative*  
585 *physiology* 151, 521-543.

586 29. Klinger, E. (1975) Consequences of Commitment to and Disengagement from Incentives. *Psychol*  
587 *Rev* 82, 1-25.

588 30. Keramati, M. et al. (2017) Misdeed of the need: towards computational accounts of transition to  
589 addiction. *Curr Opin Neurobiol* 46, 142-153.

590 31. Ringstrom, T. and Schrater, P. (2019) Constraint Satisfaction Propagation: Non-stationary Policy  
591 Synthesis for Temporal Logic Planning. arXiv:1901.10405.

592 32. Camerer, C. et al. (1997) Labor supply of New York City cabdrivers: one day at a time. *Quarterly*  
593 *Journal of Economics*, 407-441.

594 33. Thaler, R. (2015) *Misbehaving: The Making of Behavioural Economics*, W. W. Norton & Company.

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

597 35. Olds, J. and Milner, P. (1954) Positive reinforcement produced by electrical stimulation of septal  
598 area and other regions of rat brain. *Journal of Comparative and Physiological Psychology* 47, 419-  
599 427.

600 36. Young, P.T. (1966) Hedonic organization and regulation of behavior. *Psychol Rev* 73 (1), 59-86.

601 37. Schuck, N.W. et al. (2016) Human Orbitofrontal Cortex Represents a Cognitive Map of State  
602 Space. *Neuron* 91 (6), 1402-1412.

603 38. Wilson, R.C. et al. (2014) Orbitofrontal cortex as a cognitive map of task space. *Neuron* 81 (2),  
604 267-79.

605 39. Padoa-Schioppa, C. (2013) Neuronal origins of choice variability in economic decisions. *Neuron*  
606 80 (5), 1322-36.

607 40. Abitbol, R. et al. (2015) Neural mechanisms underlying contextual dependency of subjective  
608 values: converging evidence from monkeys and humans. *J Neurosci* 35 (5), 2308-20.

609 41. Eldar, E. et al. (2016) Mood as Representation of Momentum. *Trends Cogn Sci* 20 (1), 15-24.

610 42. Vinckier, F. et al. (2018) Neuro-computational account of how mood fluctuations arise and affect  
611 decision making. *Nat Commun* 9 (1), 1708.

612 43. Juechems, K. et al. (2017) Ventromedial Prefrontal Cortex Encodes a Latent Estimate of  
613 Cumulative Reward. *Neuron* 93 (3), 705-714 e4.

614 44. Tsetsos, K. et al. (2014) Neural mechanisms of economic commitment in the human medial  
615 prefrontal cortex. *Elife* 3.

616 45. Caraco, T. (1981) Energy budgets, risk and foraging preferences in dark-eyed juncos (*Junco*  
617 *hyemalis*). *Behav Ecol Sociobiol.* 8, 213–217.

618 46. Kolling, N. et al. (2014) Multiple neural mechanisms of decision making and their competition  
619 under changing risk pressure. *Neuron* 81 (5), 1190-1202.

620 47. Symmonds, M. et al. (2010) A behavioral and neural evaluation of prospective decision-making  
621 under risk. *J Neurosci* 30 (43), 14380-9.

622 48. Balaguer, J. et al. (2016) Neural Mechanisms of Hierarchical Planning in a Virtual Subway  
623 Network. *Neuron* 90 (4), 893-903.

624 49. Shidara, M. and Richmond, B.J. (2002) Anterior cingulate: single neuronal signals related to  
625 degree of reward expectancy. *Science* 296 (5573), 1709-11.

626 50. Kennerley, S.W. et al. (2006) Optimal decision making and the anterior cingulate cortex. *Nat*  
627 *Neurosci* 9 (7), 940-7.

628 51. Roijers, D.M. et al. (2013) A Survey of Multi-Objective Sequential Decision-Making. *Journal of*  
629 *Artificial Intelligence Research* 48, 67-113.

630 52. Padoa-Schioppa, C. and Assad, J.A. (2006) Neurons in the orbitofrontal cortex encode economic  
631 value. *Nature* 441 (7090), 223-6.

632 53. Wike, E.L. and Barrientos, G. (1958) Secondary reinforcement and multiple drive reduction.  
633 *Journal of Comparative and Physiological Psychology* 51, 640-643.

634 54. Simonson, I. (1989) Choice Based on Reasons: The Case of Attraction and Compromise Effects.  
635 Journal of Consumer Research 16, 158-174.

636 55. Juechems, K. et al. (2019) A Network for Computing Value Equilibrium in the Human Medial  
637 Prefrontal Cortex. Neuron 101 (5), 977-987 e3.

638 56. Miller, G.A. et al. (1960) Plans and the structure of behavior, Holt, Rhinehart, & Winston.

639 57. Oudeyer, P. and Kaplan, F. (2007) What is Intrinsic Motivation? A Typology of Computational  
640 Approaches. Front Neurobotics 1:6.

641 58. Ryan, R.M. and Deci, E.L. (2000) Intrinsic and Extrinsic Motivations: Classic Definitions and New  
642 Directions. Contemporary Educational Psychology 25, 54–67.

643 59. Gottlieb, J. and Oudeyer, P.Y. (2018) Towards a neuroscience of active sampling and curiosity.  
644 Nat Rev Neurosci 19 (12), 758-770.

645 60. Guo, X. et al. (2016) Deep Learning for Reward Design to Improve Monte Carlo Tree Search in  
646 ATARI Games. arXiv.

647 61. Friston, K. (2010) The free-energy principle: a unified brain theory? Nat Rev Neurosci 11 (2), 127-  
648 38.

649