

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33

Nutrient-sensitive reinforcement learning in monkeys

Fei-Yang Huang¹ and Fabian Grabenhorst^{1,2,*}

¹Department of Physiology, Development and Neuroscience, University of Cambridge, Cambridge CB2 3DY, UK

²Wellcome Trust-MRC Institute of Metabolic Science, Addenbrooke's Hospital, Cambridge CB2 0QQ, UK

*For correspondence: fg292@cam.ac.uk

ABSTRACT

Animals make adaptive food choices to acquire nutrients that are essential for survival. In reinforcement learning (RL), animals choose by assigning values to options and update these values with new experiences. This framework has been instrumental for identifying fundamental learning and decision variables, and their neural substrates. However, canonical RL models do not explain how learning depends on biologically critical intrinsic reward components, such as nutrients, and related homeostatic regulation. Here, we investigated this question in monkeys making choices for nutrient-defined food rewards under varying reward probabilities. We found that the nutrient composition of rewards strongly influenced monkeys' choices and learning. The animals preferred rewards high in nutrient content and showed individual preferences for specific nutrients (sugar, fat). These nutrient preferences affected how the animals adapted to changing reward probabilities: the monkeys learned faster from preferred nutrient rewards and chose them frequently even when they were associated with lower reward probability. Although more recently experienced rewards generally had a stronger influence on monkeys' choices, the impact of reward history depended on the rewards' specific nutrient composition. A nutrient-sensitive RL model captured these processes. It updated the value of individual sugar and fat components of expected rewards from experience and integrated them into scalar values that explained the monkeys' choices. Our findings indicate that nutrients constitute important reward components that influence subjective valuation, learning and choice. Incorporating nutrient-value functions into RL models may enhance their biological validity and help reveal unrecognized nutrient-specific learning and decision computations.

34 INTRODUCTION

35 According to the influential Reinforcement Learning (RL) framework, animals learn by updating reward
36 values based on experience and chose by comparing these values between options¹. The RL framework has
37 been critical for identifying fundamental learning and decision variables that guide animals' behaviour,
38 including object values and action values, which provide essential decision inputs, and the reward prediction
39 error, which updates values from experience. Direct physical implementations of these theoretical constructs
40 have been discovered in the activity of neurons in primate dopamine neurons²⁻⁵, striatum^{6,7}, amygdala^{8,9}, and
41 prefrontal cortex¹⁰⁻¹³. Despite its broad explanatory power, the RL framework does not explain how learning
42 and choice depend on specific reward properties. For example, nutrients are biologically critical, intrinsic
43 components of food rewards, and an animal's survival depends on its ability to make adaptive food choices
44 that acquire specific nutrients. Investigating how nutrient rewards influence learning and choice could not only
45 enhance the biological validity of RL models. It may also guide the discovery of so-far unrecognized nutrient-
46 specific learning and decision computations, and their neuronal implementations.

47 Because nutrients are mainly acquired from food intake, an animal's ability to adapt its food choice to
48 changing nutrient availabilities critically determines its nutrient balance and long-term health. To optimize
49 nutrient intake, foraging animals adapt their feeding patterns in response to regional and seasonal variations of
50 food resources¹⁴⁻¹⁶. For instance, monkeys spend more time in food patches associated with a high probability
51 of nutritious foods (e.g., nuts) while ignoring more frequent low-nutrient foods (e.g., leaves). Primates,
52 including humans, also exhibit individual subjective preferences for specific nutrients and sensory food
53 qualities to regulate nutrient intake¹⁷⁻²⁴. Thus, ecological data suggest that animals consider both the nutritional
54 value of food and the food's availability. However, the specific learning and decision computations underlying
55 such nutrient-sensitive food choices remain unclear. Here, we examined the food choices of rhesus monkeys
56 (*Macaca mulatta*) in a dynamic foraging task that involved choices between rewards with different nutrient
57 (fat, sugar) components under varying reward probabilities.

58 Previous studies examined how monkeys adapt to changing reward probabilities^{9-13,25-27}. In probabilistic
59 learning tasks, monkeys track the high-probability option based on past choices and reward outcomes and
60 distribute their choices according to the reward probability of both options. This learning strategy has been
61 modelled by linking subjectively weighted recent rewards to current choices ('reward history') using logistic
62 regression^{25,26} and by dynamic updating of option values based on reward outcomes via RL mechanisms¹. We
63 followed these approaches and examined whether monkeys assigned higher value to more nutritious foods
64 during learning and learned faster from high-nutrient rewards.

65 First, we characterized monkeys' nutrient preferences and learning during probabilistic reward-based
66 choices. If the monkeys preferred specific nutrients, they should choose high-nutrient rewards more frequently
67 and track their changing probability more closely to maximize intake of the specific nutrient. We recently
68 showed in a nutrient-choice task without learning requirement that macaques' choices reflect underlying,
69 stable nutrient-value functions²². Accordingly, we hypothesized that nutrient-value functions also govern
70 choices during probabilistic reward learning.

71 Next, we examined whether monkeys demonstrated nutrient-specific learning. We followed established
72 approaches for characterizing the integration of past reward experiences into subjective values using logistic-
73 regression and RL frameworks^{10,11,25,26} to examine whether nutrient preferences modulated reward learning.
74 To account for nutrient-specific learning, influences of recent reward and choice histories on current choice
75 should be higher for high-nutrient reward. Accordingly, the value function in a formal RL model should
76 incorporate higher preferences for high-nutrient rewards ('nutrient-value function'). In addition, the animals
77 may assign higher weights to reward outcomes with particular nutrient content, as reflected by influences on
78 learning rate ('nutrient-specific learning rates').

79 Finally, based on behavioral evidence for nutrient-sensitive reinforcement learning, we propose candidate
80 neuronal mechanisms necessary to implement nutrient-specific learning and decision computations, as a
81 framework to guide future neurophysiological recordings.

82

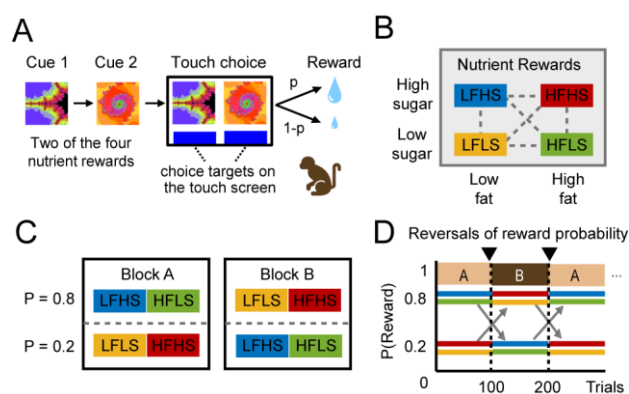


Fig. 1. Nutrient foraging task. A) Task structure. In each trial, two visual cues appeared sequentially on a touch screen before reappearing in a left-right arrangement as choice targets. Following the touch choices, the monkeys received the liquid reward associated with the chosen cue. The amount of the delivered reward depended on a prespecified reward probability (p). B) Four types of liquids with 2×2 factorial fat and sugar levels were offered to the monkeys: the low-fat low-sugar (LFLS) liquid, the high-fat low-sugar (HFLS) liquid, the low-fat high-sugar (LFHS) liquid, and the high-fat high-sugar (HFHS) liquid. C) Reward probabilities associated with the different reward types reversed between blocks of trials in a testing session. In block A, LFHS and HFLS were associated with a high probability ($p = 0.8$) of receiving the large reward, LFLS and HFHS were associated with a low probability ($p = 0.2$) of large reward; these probabilities reversed in block B. D) Each session started with either block A or block B and the reward probabilities reversed every 100 trials between the two block types, with typically 3-5 reversals per session.

RESULTS

Two monkeys performed in a dynamic foraging task to obtain different nutrient-defined liquid rewards (Fig. 1A). In each choice trial, the monkeys were presented with two visual cues from a set of four, chose between the two cues, and received either a large amount ('rewarded') or a small amount ('non-rewarded') of the cue-associated liquid reward, depending on a prespecified reward probability (p). We used new, untrained visual cues in each session to avoid influences of prior experience. Session-specific visual cues were each associated with one of four different rewards; cue-reward associations were fixed within each session. To examine whether fat and sugar biased learning from reward outcomes, we used liquid rewards from a 2×2 factorial design with fat and sugar levels as factors (Fig. 1B; LFLS: low-fat low-sugar; HFLS: high-fat low-sugar; LFHS: low-fat high-sugar; HFHS: high-fat high-sugar). At the start of each session, two rewards (LFLS/HFHS or LFHS/HFLS) were associated with a high probability of obtaining a large reward ($p = 0.8$), and the other two rewards were associated with a low reward probability ($p = 0.2$) (Fig. 1C, block A or block B). We reversed the reward probabilities every 100 trials throughout the session ($p = 0.2 \rightarrow 0.8$; $p = 0.8 \rightarrow 0.2$) to encourage continual learning from reward outcomes (Fig. 1D). Notably, this design offered the monkeys equal availability of fat and sugar in all choice trials irrespective of block type because there were always two high-probability and two low-probability options for both high-fat and high-sugar rewards. All liquids were matched in flavour (blackcurrant or peach) and other ingredients (protein, salt, etc); therefore, differential learning and choice patterns could be attributed to the nutrient content of the rewards.

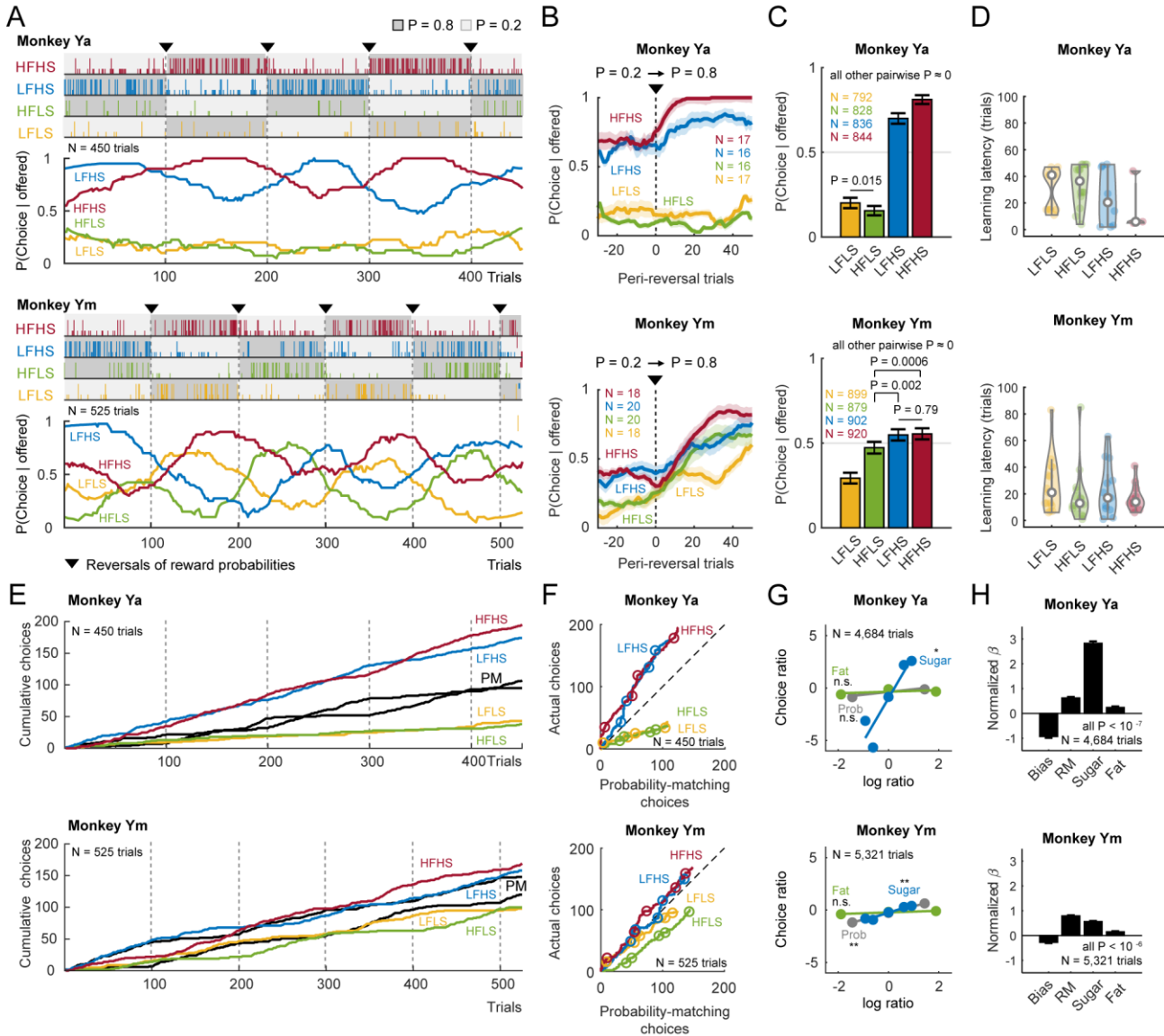
Nutrients bias reward learning and food choices

The behaviour in two example sessions (Fig. 2A) showed that both monkeys exhibited preferences for specific nutrients while tracking changing reward probabilities. Monkey Ya's choices (Fig. 2A, top) were dominated by a general preference for high-sugar rewards, with a smaller impact of reward probability on choice. Specifically, monkey Ya chose the high-sugar rewards frequently even when they were associated with a lower probability of obtaining a large reward amount; in addition, choice frequencies tracked changing reward probabilities, particularly for the high-sugar rewards. By contrast, monkey Ym's choices (Fig. 2A, bottom) reflected both a preference for high-nutrient content and a strong dependence on reward probability. Specifically, within a given trial block, monkey Ym preferred high-nutrient rewards over low-nutrient rewards with matched reward probabilities (compare red and yellow curves) but would reduce his choices for the high-nutrient reward when it was associated with a relatively lower reward probability.

The patterns observed in single sessions were also observed in averaged data across sessions. Overall, the monkeys' choice probabilities increased when reward probabilities switched from low ($p = 0.2$) to high ($p = 0.8$), as evident by averaged choice probabilities around probability-reversal points (Fig. 2B). Importantly, the monkeys responded differently to probability changes for rewards that differed in fat and sugar content, with more pronounced probability increases for high-nutrient rewards and specifically high-sugar rewards (Fig. 2B).

4

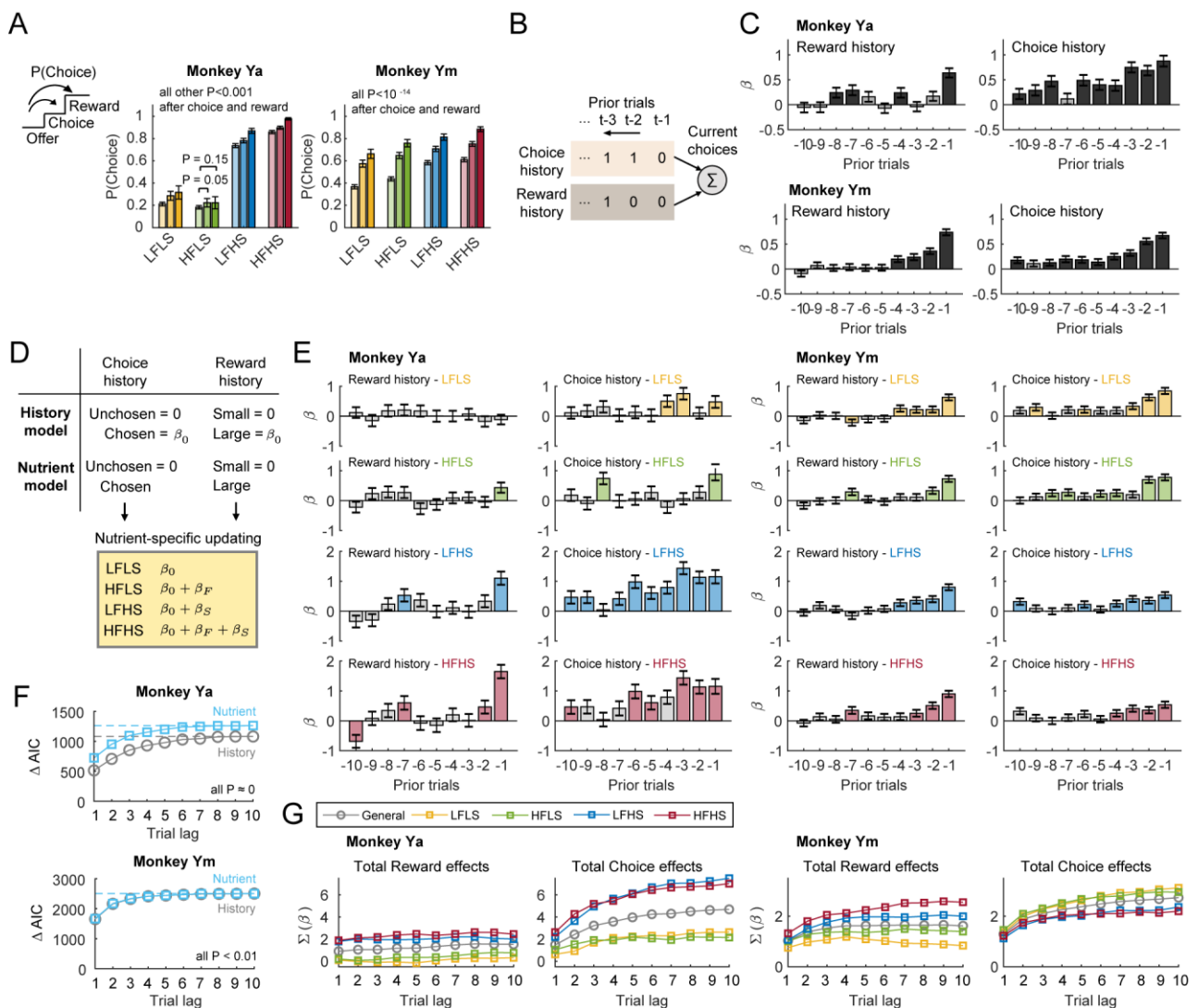
131 When reward probabilities were stable (between reversal points), monkey Ya showed a strong preference for
 132 the high-sugar rewards irrespective of fat level, whereas monkey Ym showed graded preferences for both high-
 133 fat and high-sugar rewards over the low-nutrient option (Fig. 2C). Immediately following the probability
 134 reversals, the monkeys had shorter learning latencies for high-nutrient rewards: they adjust their choices more
 135 quickly to the changed reward probabilities when high-sugar and high-fat rewards were offered, which
 136 indicated that learning was sensitive to the nutrient content of reward outcomes (Fig. 2D). Thus, the monkeys
 137 preferred high-nutrient rewards, tracked changing reward probabilities in a nutrient-dependent manner, and
 138 learned faster from high-nutrient reward outcomes.
 139



140
 141
 142
 143
 144
 145
 146
 147
 148
 149
 150
 151
 152
 153
 154
 155
 156

Fig. 2. Nutrient-sensitive learning and choice in monkeys. A) Choices and reward outcomes in a single session for monkey Ya (top) and monkey Ym (bottom). Each tick mark represents a choice of a specific reward type; long marks indicate large reward outcome, short marks indicate small reward outcomes. Reward types in dark-gray blocks were associated with high reward probability ($p = 0.8$) and those in light-gray blocks were associated with low reward probability ($p = 0.2$). Choice curves showed running-average choice patterns of each reward. B) Learning curves. Mean running-averaged choice frequencies aligned to probability reversals ($p = 0.2 \rightarrow 0.8$) indicate how choices depend on both reward-probability changes and nutrient content. (N: number of tested sessions). C) Reward preferences. Average choice frequencies indicate preferences among the four reward types. The choice frequencies were computed after sessions were truncated, including only probability-balanced trials for all reward types. (mean \pm s.e.m.) (N: number of trials). D) Learning latency. The number of trials from probability reversal to the first significant change point in the cumulative choice record (see *Methods*) indicates latency to adapt choices after probability changes. E) Monkeys' single-session cumulative choice records deviate from the pure probability-matching strategy. PM: probability-matching choice strategy, calculated by matching choices to the past ratio of large/small rewards, irrespective of reward type. F) Direct comparisons of monkeys' choices with probability-matching choices. Circles indicate probability reversals. G) Nutrient-sensitive matching behavior. Correlations of choice ratios with fat, sugar, and probability ratios, respectively. H) Normalized regression coefficients of probability ratios, fat ratios, and sugar ratios on choice ratios. (mean \pm s.e.m.).

5



157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184

Fig. 3. Nutrient-specific reward and choice histories influence monkeys' choices. A) Choice probabilities for different rewards depended on recent experience. Choice probabilities when the same reward was chosen on the previous trials ('choice'), when a large reward was received on the previous trial ('reward'), and irrespective of last-trial choice and reward outcomes ('offer'). B) History model for explaining the choice. The history model included regressors for choice history and reward history, with choice history = 1 (chosen) or 0 (not chosen) and reward history = 1 (rewarded) or 0 (non-rewarded) for the past 10 trials. C) Results, history model. Logistic-regression coefficients show influences of reward (left) and choice history (right) on current choices. (mean \pm s.e.m.; dark-gray bars: $P < 0.05$; light-gray bars: non-significant) D) Nutrient model for explaining choices. In contrast to the assumption of uniform history effects across reward types in the history model, the nutrient model examined nutrient-specific history effects by including additional nutrient-specific history regressors. E) Results, nutrient model. Nutrient-specific logistic regression coefficients for current choices. ($P < 0.05$, yellow: LFLS; green: HFSL; blue: LFHS; red: HFHS; light-gray bars: non-significant) F) Model performances and history lengths. Model performance improved with history length based on $\Delta AIC = AIC(\text{trial lag} = 0) - AIC(\text{trial lag} = i, i = 1, 2, \dots, 10)$. AIC = Akaike Information Criteria. History length-matched nutrient models and history models were compared using the loglikelihood test. Higher ΔAIC values indicated that the nutrient model outperformed the history model in all history length-matched comparisons. G) Aggregated effects of reward and choice history increased with history lengths and reflected nutrient composition, indicated by the cumulative reward or choice history regression coefficients over recent trials.

The preferences for fat and sugar biased the monkeys' choices away from a pure probability-matching (PM) strategy, which predicted distributed choices according to the relative frequency of receiving large rewards from each option. In the two example sessions, choices for the high-sugar rewards accumulated more rapidly than predicted by the PM strategy, whereas choices for low-sugar rewards accumulated more slowly (Fig. 2E). Specifically, compared to the PM strategy, monkey Ya significantly over-matched the high-sugar rewards and under-matched the low-sugar rewards, irrespective of reward-fat level. These patterns were much less pronounced in monkey Ym (Fig. 2F). Specifically, the choice ratios of monkey Ya were dominated by the sugar ratios but those in monkey Ym were jointly determined by the probability ratios and sugar ratios (Fig. 2G). Multiple regression confirmed that, in addition to the probability ratios, both the fat and sugar ratios significantly influenced the choice ratios (Fig. 2H). Notably, both monkeys' choices were explained by similar

185 effect sizes of the probability ratios and the fat ratios. However, the effects of sugar ratios were particularly
186 strong in monkey Ya but slightly weaker than the influences of probability ratios in monkey Ym.

187 Taken together, these results suggested that the specific nutrient composition of food rewards and the
188 animals' individual preferences for sugar and fat biased learning and choice.

189

190 **Nutrient-specific reward history and choice history influence monkeys' choices**

191 One strategy to respond to unsignaled changes in reward probabilities is to choose based on recent choices
192 and reward outcomes. Because the choice outcomes reflect the underlying reward probability, this strategy
193 adapts choices to the changing reward probabilities and can help to optimize reward rate and nutrient-intake
194 levels. Consistent with these notions, we found that monkey Ym tended to repeat his choices, particularly after
195 receiving a large reward on the previous choice; this effect was evident across all reward types (**Fig. 3A**, right).
196 By contrast, the tendency to repeat choices was less pronounced for the low-sugar rewards in monkey Ya (**Fig.**
197 **3A**, left). This result suggested that both recent choices and the reward outcomes increased choice repetition,
198 but the influences depended on individual nutrient preferences.

199 To formally characterize the learning from recent choices and reward outcomes, we modelled the trial-
200 by-trial choices in a logistic regression model (history model, see *Method*) that accounted for whether the
201 option was chosen in previous offers (choice history) and whether the previous choices were rewarded (reward
202 history) (**Fig. 3B**). The regression coefficients showed that both the choice and reward history reinforced
203 current choices and that these effects decayed for more remote past trials (**Fig. 3C**). Given the monkeys'
204 preferences for fat and sugar, we next examined whether these reward- and choice-history effects also
205 depended on the nutrient composition of reward outcomes and choice offers. We tested this possibility by
206 including nutrient-history interaction regressors in the history model (nutrient model, see *Method*). These
207 interaction terms would capture any additional reinforcing effects from specific nutrients by decomposing the
208 aggregated reward and history effects into the effects of baseline low-nutrient liquid (β_0), high-fat content (β_F),
209 and high-sugar content (β_S), depending on the fat and sugar levels of the offered reward types (**Fig. 3D-E**).
210 Larger history regression coefficients for sugar compared to fat suggested that recently obtained high-sugar
211 reward outcomes had a stronger impact on current-trial choice than recently obtained high-fat rewards in both
212 monkeys. However, the two monkeys differed in their tendency to repeat choices for high-fat and high-sugar
213 liquids, as indicated by the nutrient-specific choice-history coefficients. Monkey Ya repeated the high-sugar
214 choices more frequently than choices for low-nutrient rewards and high-fat rewards. By contrast, monkey Ym
215 repeated choices slightly less frequently for the high-sugar rewards. Importantly, although the explanatory
216 power of both models increased with history length, the nutrient model outperformed the history model in all
217 history length-matched comparisons (**Fig. 3F**). These history effects showed distinct temporal dynamics in the
218 two monkeys although they both decayed either in the history model or the nutrient model (**Fig. 3G**).

219 These results indicated that both monkeys' choices depended on the recent histories of obtaining and
220 choosing rewards with specific nutrient content.

221

222 **Reinforcement learning based on nutrient-specific values**

223 The temporal dynamics of the nutrient-specific reward- and choice-history effects suggested that the
224 monkeys constantly updated their choices based on recent choices and reward outcomes. RL models that
225 update trial-by-trial reward values for each option based on the reward outcomes are well-suited to model such
226 adaptive choices. However, canonical RL models typically do not account for the nutrient composition of food
227 rewards, and accordingly cannot explain the presently observed nutrient preferences and nutrient-specific
228 learning effects. Therefore, we developed a nutrient-sensitive RL model that incorporated subjective nutrient
229 values to model how specific nutrients (fat, sugar) differentially influenced the trial-by-trial updating of
230 expected reward values and their influence on the choice (**Fig 4A**). Instead of updating the value of the chosen
231 reward with a binary reward outcome, our model updated reward values based on the nutrient composition of
232 each reward type as given below,
233

$$234 \quad Q_i(t+1) = Q_i(t) + \alpha \cdot [V_i(t) - Q_i(t)], \quad V_i(t) = \begin{cases} 1/(V_F \cdot V_S \cdot V_{FS}) & , i(t) = LFLS \\ V_F/(V_F \cdot V_S \cdot V_{FS}) & , i(t) = HFLS \\ V_S/(V_F \cdot V_S \cdot V_{FS}) & , i(t) = LFHS \\ 1 & , i(t) = HFHS \end{cases}$$

235 , where the value for reward i , Q_i , was updated depending on the chosen reward type on trial t , $i(t)$ and its
236 nutrient-specific reward value, $V_i(t)$. V_F , V_S , and V_{FS} denoted the subjective value of high-fat content, high-

7

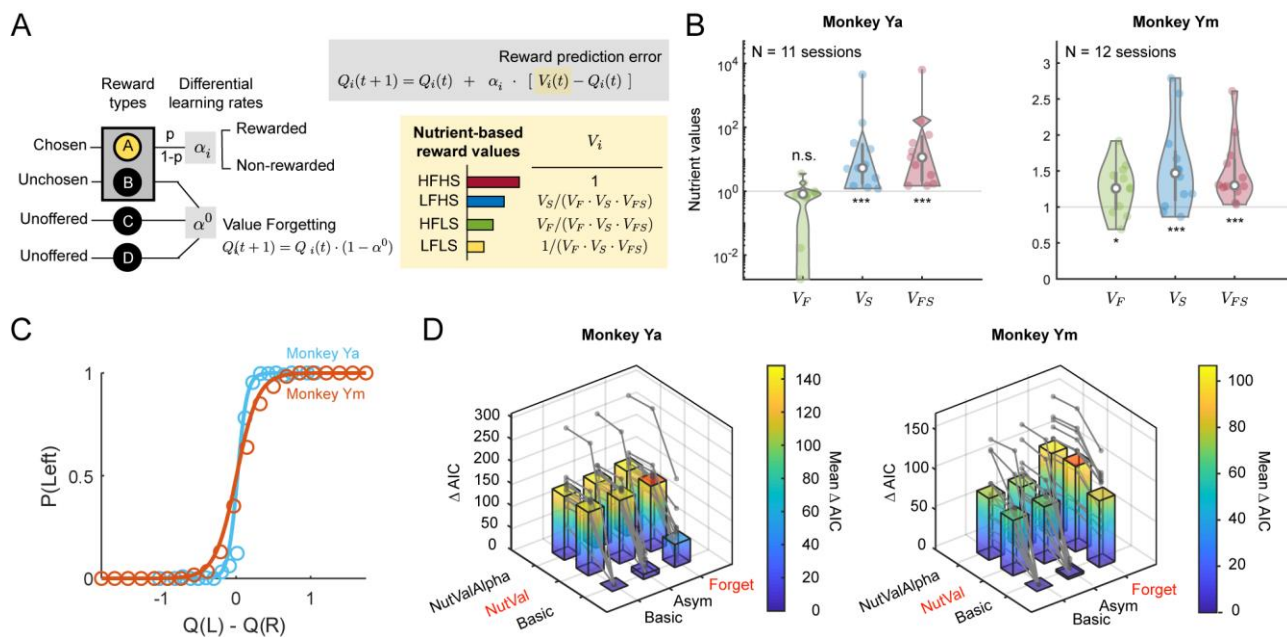
237 sugar content, and their interaction, respectively, on the common scale of the low-nutrient reward value.
 238 Therefore, any nutrient value larger than 1 suggested a preference for the specific nutrient; values for V_{FS}
 239 larger than 1 indicated supra-additive values of fat and sugar. Without loss of generality, we normalized all
 240 reward values to the highest nutrient value, ($V_F \cdot V_S \cdot V_{FS}$), to constrain all reward values between 0 and 1. For
 241 the unchosen and unoffered rewards, we allowed the values to decay as follows,

$$242 \quad 243 \quad Q_j(t+1) = Q_j(t) \cdot (1 - \alpha^0), \quad \forall j \neq i(t)$$

244 , where the values of the unchosen and unoffered rewards, $Q_j(t)$, were discounted according to a forgetting
 245 rate (α^0), which would be 0 for perfect (but biologically implausible) value memory.

246 The results of fitting this nutrient-sensitive RL model to each monkeys' choices and reward outcomes in
 247 each session confirmed that both monkeys assigned higher values to the high-sugar choice options and that
 248 monkey Ym assigned higher value to fat but monkey Ya did not (**Fig. 4B**). The high-fat high-sugar reward
 249 was also valued higher than the low-nutrient reference, but the fat values and the sugar values did not show
 250 supra-additive effects in monkey Ya but negative interactions in monkey Ym when determining the reward
 251 values (**Fig. S1**). The model-derived subjective values for fat and sugar accurately predicted the monkeys'
 252 choices (**Fig. 4C**). The nutrient-sensitive RL model outperformed alternative RL models involving
 253 combinatorial differential learning rates and nutrient-specific parameters (**Fig. 4D**; see *Methods*). Notably,
 254 there was no evidence for nutrient-specific learning rates but only a significant but small forgetting rate for
 255 monkey Ym (**Fig. S2**).

256 Thus, the monkeys' stochastic choices for rewards with specific nutrient compositions were well
 257 explained by a nutrient-sensitive RL model that assigned nutrient-specific values to reward outcomes.
 258
 259



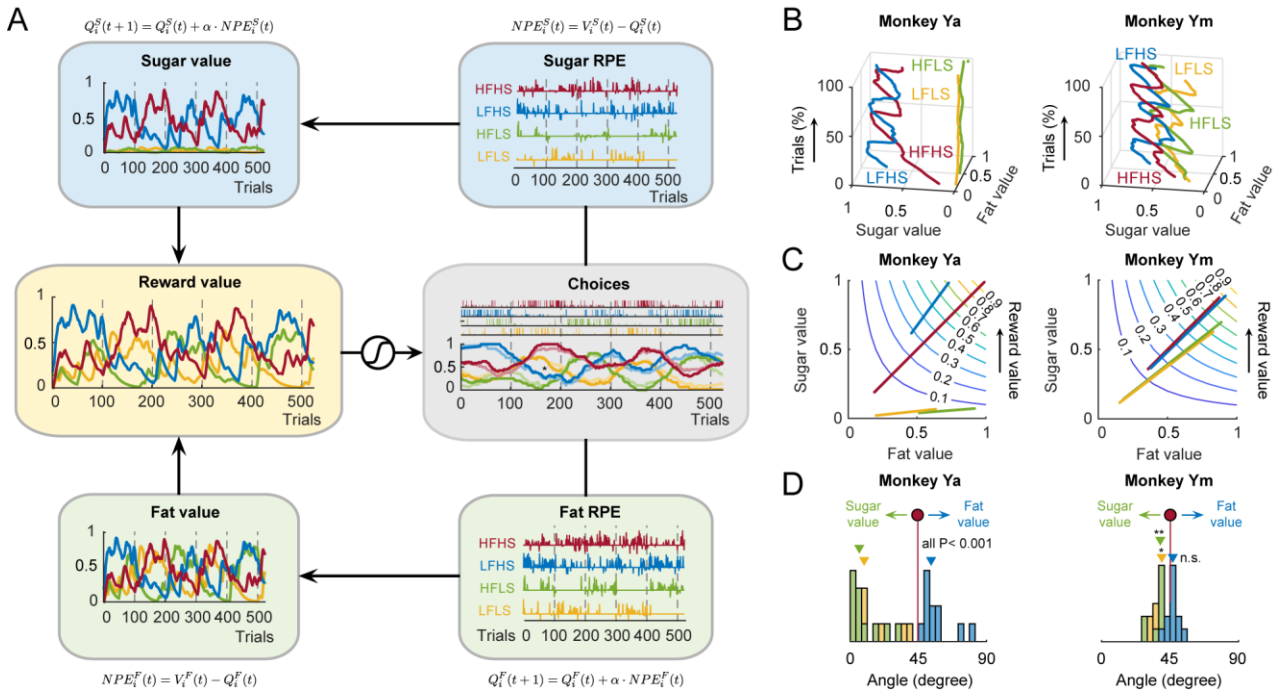
260
 261
 262 **Fig. 4. Nutrient-sensitive reinforcement learning models.** A) The nutrient-value RL model (NV-RL model). Reward values were
 263 updated based on the nutrient-specific values, $V_i(t)$, with V_F indicating values for the high-fat content, V_S for the high-sugar content,
 264 V_{FS} for the high fat-sugar combination, and 1 for the low-nutrient value reference. We normalized all reward values to the highest
 265 reward value ($V_F \cdot V_S \cdot V_{FS}$) to constrain all reward values between 0 and 1. B) Nutrient-specific reward values. The distributions of
 266 fitted nutrient-specific reward values across trials (monkey Ya: log scale; monkey Ym: linear scale). All reward values were tested
 267 against equal values for all reward types (nutrient values =1), Wilcoxon signed-rank test. C) Nutrient-value functions. Psychometric
 268 curves based on integrated values, calculated with the nutrient-value RL model, indicate that both monkeys' choices depended on
 269 nutrient-dependent value differences between choice options. D) Model comparisons. The main nutrient-value RL model (NutVal-
 270 Forget model) was systematically compared with alternative RL models involving combinations of differential learning rates (*NutVal*
 271 = nutrient-specific values; *NutValAlpha* = nutrient-specific values + learning rates, Figure S2A) and nutrient-specific parameters (*Asym*
 272 = independent learning rate for the non-rewarded chosen option; *Forget* = value-forgetting for unchosen and unoffered options). Models
 273 were compared using Akaike Information Criterion (AIC). All model AICs were subtracted from the AIC of the basic RL model
 274 ($\Delta AIC = AIC_{basic} - AIC$) for comparison. The higher mean ΔAIC indicated better model performance (red: the best fitting model).
 275
 276

277 **Value updating based on distinct sugar and fat value components**

278 The nutrient-sensitive RL model implied that the animals can independently track values for specific fat and
 279 sugar nutrients, and integrate them into a scalar value that guided choices. To better understand the dynamics
 280 of this nutrient-specific value tracking and updating, we modelled the dynamic learning of individual nutrient
 281 values in a nutrient prediction error-based RL model (NPE-RL) in which the reward value on trial t , $Q_i(t)$, was
 282 jointly determined by individual fat value and sugar value components (**Fig. 5A**, see *Methods*).

283 The NPE-RL model characterized how fat and sugar values could (i) separately adapt to changes in reward
 284 probabilities as indicated by experienced outcomes, and (ii) flexibly determine the integrated reward values
 285 for specific choice options, based on their nutrient composition. Specifically, the fat and sugar RPEs for each
 286 reward updated the fat and sugar values, respectively, which were then combined into integrated reward values
 287 to guide choices (**Fig. 5A**). Decomposing the reward values into two independent nutrient components revealed
 288 each animal's idiosyncratic sensitivity of reward values to individual nutrient constituents. To illustrate the
 289 dynamic, nutrient-specific value updating, we plotted the evolving value trajectories within a session in a space
 290 defined by the separate fat and sugar value components (**Fig. 5B**). These trajectories indicated that the updating
 291 of reward values in monkey Ya was primarily based on the sugar value component over the fat value
 292 component, whereas both fat and sugar value components contributed to value learning in monkey Ym (**Fig.**
 293 **5B**).

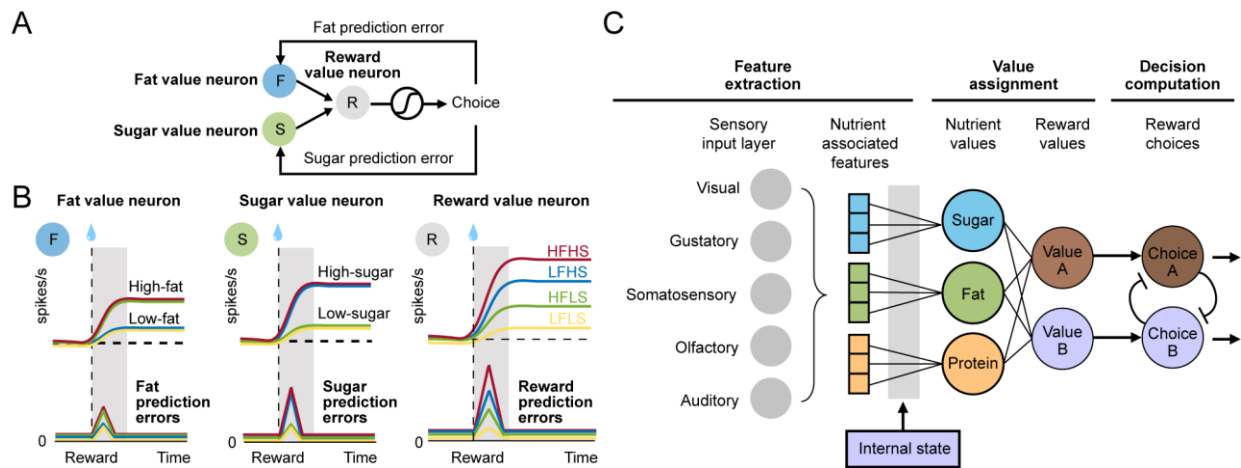
294



295

296

297 **Fig. 5. Dynamics of sugar and fat value components in nutrient-sensitive reinforcement learning.** A) Nutrient-specific value
 298 updating. Nutrient-specific values for sugar (top) and fat (bottom) were updated based on discrepancies between previous choice
 299 outcomes and predicted nutrient rewards (nutrient prediction errors, NPE); sugar and fat values were integrated into composite reward
 300 values that guided choices. B) Trajectories of nutrient-specific values within sessions. The value trajectories tracked the evolving
 301 reward values and their nutrient components with choice trials. C) Projected reward-value trajectories and iso-value contour curves.
 302 Each segment showed the ranges and orientations of the fluctuating reward values in the nutrient value space. The diagonal line
 303 represented equal contributions of the nutrient components to the reward values. D) Nutrient sensitivities of reward values.
 304 Distributions of the rotating angles quantified the relative changes of nutrient values during reward value updating ($\Delta V_S / \Delta V_F$) across
 305 sessions. Wilcoxon signed-rank test.



306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

328

329

330

331

332

333

334

335

336

337

338

339

340

341

342

343

344

345

346

Fig. 6. Neuronal mechanisms for nutrient-sensitive reinforcement learning and choice. A) Nutrient-sensitive reinforcement learning architecture. Fat-value neurons (F) and sugar-value neurons (S) each update the fat and sugar components of the value predictions and provide input to the reward-value neurons (R) that code integrated values for decision computations. B) Predicted neuronal responses of fat-value, sugar-value, and reward-value neurons. Fat-value neurons (F) are updated based on fat-specific reinforcement learning (fat prediction errors) from delivered rewards, with higher responses to high-fat compared to low-fat rewards. An equivalent process operates for sugar-value neurons (S). Together, these nutrient-value neurons converge onto reward-value neurons to code scalar value signals in a common currency for downstream decision computations. C) Nutrient-sensitive decision-making neural network. Sensory properties of foods are detected via multiple sensory channels and integrated into nutrient-associated feature representations that determine nutrient values depending on the internal physiological state. Nutrient values then flexibly inform reward values for decision computations, based on the nutrient composition of food rewards.

The distinct sensitivities of reward values to specific nutrient components were illustrated by projections of dynamic reward value trajectories onto the nutrient value space, where ‘iso-value contours’ visualized levels of equal reward values (Fig. 5C). If reward values were equally sensitive to both the sugar and fat nutrient components, the value trajectories should fall onto the 45-degree diagonal line in nutrient value space. Because we normalized the nutrient values to the HFHS reward, the value trajectory for HFHS would be at the diagonal for both monkeys (Fig. 5C, red). However, higher sensitivity to the sugar value components compressed the low-sugar value trajectories along the sugar value axis and rotated the trajectories towards the fat value axis (clockwise); similarly, higher sensitivity to the fat value components rotated the low-fat trajectories towards the sugar value axis (counterclockwise). For example, monkey Ya showed a slight counterclockwise-rotated LFHS trajectory and marked clockwise-rotated low-sugar trajectories, indicating his weak preference for fat and the strong preference for sugar, respectively. In contrast, monkey Ym showed only mild clockwise-rotated low-sugar trajectories and negligible rotation for the LFHS value trajectory, reflecting his mild sugar preference and non-significant fat preference.

The rotating angles of the value trajectories in the nutrient value space quantified the relative changes of sugar and fat values on the value trajectories ($\Delta V_S / \Delta V_F$), therefore highlighting the contributions of each nutrient value to the overall reward values. Compared to the steepest 45-degree value gradient, value trajectories with rotating angles larger or smaller than 45 degrees updated reward values with distinct contributions of each nutrient value. Specifically, reward values were mostly updated from the fat values when the angles were smaller than 45 degrees, but more from the sugar values if the angles were larger than 45 degrees. Across sessions, the nutrient-specific contributions of reward values, indicated by the orientations of the value trajectories, recapitulated the subjective nutrient values estimated by the nutrient-sensitive RL model (Fig. 5D).

Thus, subjective nutrient-value functions guided the dynamic updating and the integration of reward values based on individual nutrient-specific components.

347 **DISCUSSION**

348 We investigated monkeys' choices for different nutrient-defined rewards under varying reward
349 probabilities. We found that the nutrient composition of rewards strongly influenced choices and learning. The
350 animals generally preferred rewards that were high in nutrient content but also showed individual preferences
351 for sugar and fat, consistent with the assignment of subjective values to choice options. The animals' nutrient
352 preferences affected how they adapted their choices to changing reward probabilities. Specifically, the
353 monkeys learned faster from preferred nutrient-rewards and chose them frequently even under low reward
354 probabilities (i.e., low probability of obtaining large reward amounts). Influences of past rewards on current
355 choice were well described by a reward-history analysis. As in previous studies^{11,25}, more recent rewards had
356 a stronger influence on the monkeys' choices. Critically, we also found that the impact of reward history
357 depended on the nutrient composition of past rewards: the effect of past rewards high in preferred sugar content
358 was stronger compared to that of less preferred low-nutrient or fat rewards. The history of past choices,
359 irrespective of reward outcomes, also had a significant and nutrient-dependent effect on choice, with stronger
360 effects of past choices for preferred nutrient rewards. We proposed a nutrient-sensitive RL model that captured
361 the influences of preferred nutrients on learning and choice. The model updated the value of individual sugar
362 and fat components of expected rewards trial by trial, based on recently experienced rewards, and integrated
363 these components into scalar values that explained the monkeys' choices. These results suggest that nutrients
364 constitute important reward components that influence subjective valuation, learning and choice, and that
365 canonical RL models can be usefully extended to capture such nutrient-specific values.

366 Previous studies of reinforcement learning in macaques revealed important influences on learning and
367 choice, including effects of reward and choice history^{5,7,11-13,25-27}, the variance of recent rewards¹⁰, novelty and
368 reward rarity^{9,28}, and social observations⁸. Importantly, these studies did not vary the composition of reward
369 outcomes and thus could not test whether specific reward components differentially affected learning and
370 choice. We reasoned that nutrients are biologically critical reward components that are essential for survival
371 and that monkeys should prefer high-nutrient rewards and adapt their choices to optimise nutrient intake. By
372 manipulating the sugar and fat content of our liquid rewards, we confirmed that the monkeys' learned
373 differently from these different rewards.

374 Previous studies demonstrated that macaques have sophisticated preferences for different reward types
375 that comply with principles of economic choice theory^{2,29-31} but did not examine how different rewards affect
376 learning. Here we showed that subjective preferences for specific nutrients influenced how monkeys tracked
377 the changing reward probabilities of choice options. Specifically, both animals learned faster from preferred
378 nutrient rewards. Moreover, they based their choices on both subjective valuations of offered reward types and
379 estimates of current reward probabilities. This latter finding confirms the result from a previous study that
380 macaques integrate reward type and probability information to express subjective preferences²⁹; different from
381 that study, our monkeys were required to derive probability information from past reward experiences rather
382 than from explicit visual cues.

383 Crucially, by varying the nutrient composition of rewards, we investigated reinforcement learning and
384 choice for biologically important, universal reward components. Nutrients are basic building blocks of foods
385 that are sensed by dedicated taste and oral-texture mechanisms^{22,32-34} and engage physiological and homeostatic
386 processes^{35,36}. Moreover, evidence from ecology and human metabolic sciences points to specific behavioral
387 mechanisms that regulate nutrient intake. For example, ecological studies identify a 'nutrient-balancing
388 mechanism' in wild macaques that promotes reproductive and survival success^{14,37-39}. In humans, reduced
389 protein in ultra-processed foods increases energy intake by 'protein leveraging', a mechanism that regulates
390 food choice to counter protein deficits^{35,40-42}. A 'fat-appetite mechanism' emerges in human monogenic obesity
391 affecting melanocortin-signalling²¹. We recently showed that in macaques, nutrients and sensory food qualities
392 (taste, viscosity, oral friction) shape human-like economic preferences²². Our approach makes a first step
393 towards integrating the influential RL framework with these nutrient-dependent behavioral processes and thus
394 enhance its biological validity.

395 The concept of nutrient homeostasis in metabolic sciences suggests that internal states modulate nutrient
396 values to guide state-dependent food choices. Recent homeostatic RL models explain the value of rewards as
397 discrepancies between the current state and physiological setpoints⁴³. This approach views reward values as
398 physiological signals that serve to maintain homeostasis. However, fat and sugar can be preferred even without
399 corresponding nutrient deficits^{22,23,44}; therefore, the hedonic values of foods cannot be explained solely by
400 homeostatic regulations of nutrient deficits. Future experiments could challenge the nutrient states of animals
401 during food choices to estimate empirical nutrient-value functions from state-dependent choice patterns to
402 refine these models.

403 We described a nutrient-specific learning mechanism that updates value estimates for separate fat and
404 sugar reward-components and integrates this information to guide adaptive food choices. This mechanism
405 implies parallel nutrient valuation systems that detect and evaluate the nutrient components depending on
406 internal states. The neuronal implementation of this mechanism would require neurons that encode individual
407 nutrient values (nutrient-value neurons) and dynamically update these nutrient values via nutrient prediction
408 error signals (**Fig. 6A**). At a neural-network level (**Fig. 6B**), these nutrient-value neurons would extract
409 nutrient-specific features from a food's sensory properties to guide food choices. Importantly, physiological-
410 state signals could modulate the neural representations of nutrient values to allow for state-dependent valuation
411 of food rewards. Therefore, we propose nutrient-value neurons and nutrient prediction error signals as potential
412 substrates for nutrient-sensitive learning and choice.

413 Our findings within a nutrient-based RL paradigm and our proposed computational framework have
414 implications for value-based learning and decision theories and underlying neural mechanisms. Because
415 nutrients provide energy and serve physiological functions for survival, animal reward systems should be
416 shaped by nutrient availability in the environment and evolved dedicated mechanisms for adaptive nutrient-
417 sensitive decision-making. By decomposing the trial-by-trial reward values that guide reinforcement learning
418 into nutrient-value components, we identified candidate signals that could be encoded by neurons in the reward
419 and decision systems of the primate brain. The midbrain dopamine neurons, orbitofrontal cortex and amygdala
420 participate in decision-making, reinforcement learning, and food evaluation^{2,3,8,31,45-48} and thus constitute
421 suitable targets for testing these hypotheses experimentally.

422

423

424

425 **METHODS**

426

427 **Animals.** Two adult male rhesus macaques (*Macaca mulatta*) were trained in the study: monkey Ya (weight
428 during the experiments: 17-19 kg, age: 6 years) and monkey Ym (12-13 kg, age 6 years). The animals were
429 trained and tested approximately one to two hours per day and five days per week for 6 months. Both monkeys
430 participated in another nutrient choice study using the same dairy-based nutrient rewards as in this study. The
431 animals were on a standard diet for laboratory macaques, composed of high-protein dry pellets (% calories
432 provided by protein: 30.36%, fat: 13.29%, carbohydrates: 56.34%), dried fruits, seeds, nuts, and fresh fruits
433 and vegetables. We monitored the monkeys' health condition and body weights to ensure their welfare after
434 introducing high-calorie rewards. No effects of these rewards on the animals' health were observed. Each
435 testing day, the animals had free access to the standard diet before and after the experiments and received their
436 main liquid intake in the laboratory. The animals' body weights increased as expected for growing animals.

437 All animal procedures conformed to US National Institutes of Health Guidelines. The experiments have
438 been regulated, ethically reviewed and supervised by the following UK and University of Cambridge (UCam)
439 institutions and individuals: UK Home Office, implementing the Animals (Scientific Procedures) Act 1986,
440 Amendment Regulations 2012, and represented by the local UK Home Office Inspector; UK Animals in
441 Science Committee; UCam Animal Welfare and Ethical Review Body (AWERB); UK National Centre for
442 Replacement, Refinement and Reduction of Animal Experiments (NC3Rs); UCam Biomedical Service (UBS)
443 Certificate Holder; UCam Welfare Officer; UCam Governance and Strategy Committee; UCam Named
444 Veterinary Surgeon (NVS); UCam Named Animal Care and Welfare Officer (NACWO).

445

446 **Experimental Design**

447

448 **Nutrient rewards.** We prepared nutrient-controlled liquids with 2×2 fat and sugar levels to examine whether
449 fat and sugar biased learning from reward outcomes (**Fig. 1B**; LFLS: low-fat low-sugar; HFSL: high-fat low-
450 sugar; LFHS: low-fat high-sugar; HFHS: high-fat high-sugar). The liquids were matched in flavor (peach or
451 blackcurrant), temperature, protein, salt and other ingredients (see²² for detailed liquid compositions). We used
452 commercial skimmed milk and whole milk (British skimmed milk and British whole milk, Sainsbury's
453 Supermarkets Ltd., UK) as baseline low-fat and high-fat liquids and flavored the liquids with fruit juice to
454 increase palatability

455

456 **Nutrient foraging task.** The four nutrient reward types were associated with four untrained visual cues,
457 respectively, in each session. When a choice trial started, the monkeys were first presented with two of the
458 four visual cues, made a touch-monitor choice between the two cues, and then received either a large amount
459 ('rewarded') or a small amount ('non-rewarded') of the cue-associated liquids depending on its prespecified
460 reward probability (p) (**Fig. 1A**). When the session started, two of the rewards (LFLS/HFHS or LFHS/HFSL)
461 were offered in high reward probabilities ($p=0.8$), and the other two rewards in low reward probabilities ($p=0.2$)
462 (**Fig. 1C**, block A or block B). The reward probabilities were reversed every 100 trials ($p=0.2 \rightarrow 0.8$;
463 $p=0.8 \rightarrow 0.2$) (**Fig. 1D**).

464

465 **Data Analysis**

466 All data were analyzed using Matlab 2017 (Mathworks).

467

468 **Learning curve.** The learning curves were plotted by aligning reward-specific choices to the probability
469 reversal trials. In particular, based on the probability before and after reversals, we grouped these curves into
470 incremental ($P=0.2 \rightarrow P=0.8$) and decremental ($P=0.8 \rightarrow P=0.2$, not shown) learning curves, and plotted the
471 incremental curves in **Fig. 2B**.

472

473 **Learning latency.** The learning latency was defined as the number of trials between the first behavioral change
474 point after probability reversals. The behavioral change points were identified as the significant changing
475 points of cumulative choice slopes⁴⁹, based on two-sample t-test with criteria $P < 0.05$.

476

477

478 **Probability-matching (PM) choices.** We simulated probability matching choices by first computing the
 479 relative proportions of the reward probabilities and transform them into predicted choices as follows⁵⁰:
 480

$$481 \quad \pi_i(t) = \frac{P_i(t)}{\sum_k P_k(t)}, \quad i \in k = \{LFLS, HFLS, LFHS, HFHS\}$$

$$482 \quad A_i(t) \sim B(1, \pi_i(t))$$

483 , where $\pi_i(t)$ was the probability of choosing a specific option; $P_i(t)$ denoted the reward probability of reward
 484 i on trial t , which were summed over the stimulus set as $\sum_k P_k(t)$. The reward choices $A_i(t)$ followed the
 485 binomial distribution, based on the computed probability proportions for each reward type.
 486

488 Logistic regression analysis

489 History model

490 We used multiple logistic regression (*fitglm* function, Matlab) to model choices based on recent choices and
 491 reward outcomes as follows,

$$492 \quad \text{logit}(P_L) = \beta_0 + \beta_1 \times \text{LeftFirst} + \beta_2 \times \text{FatLv} + \beta_3 \times \text{SugarLv} + \sum_{k=1}^n (\beta_{k+3} \times Cx_k) +$$

$$493 \quad \sum_{k=1}^n (\beta_{k+n+3} \times Rx_k)$$

494 , where the probability of choosing the left option (P_L) was modelled by differential choice history (Cx_n) and
 495 reward history (Rx_n) up to recent n trials while controlling the presentation sequence ($\text{LeftFirst} = 1$, if the
 496 left option was shown first; 0, if the right option was shown first) and the nutrient information cued by
 497 pretrained visual stimuli (FatLv , SugarLv = differential fat or sugar levels = 1, if left > right; 0, if left = right;
 498 -1, if left < right). Specifically, the choice history regressors Cx_n and reward history regressors Rx_n were
 499 defined as the differences between the history variables of the left and right options,
 500

$$501 \quad Cx_n = c_n^L - c_n^R, \quad c_n^i = \begin{cases} 1, & \text{if option } i \text{ was chosen } n \text{ trials earlier} \\ 0, & \text{if option } i \text{ was not chosen } n \text{ trials earlier} \end{cases}$$

$$502 \quad Rx_n = r_n^L - r_n^R, \quad r_n^i = \begin{cases} 1, & \text{if option } i \text{ was chosen and rewarded } n \text{ trials earlier} \\ 0, & \text{otherwise} \end{cases}, \quad i \in \{L_n, R_n\}$$

503 Notably, the history regressors for each option coded past trials in terms of the offered trials because the
 504 unoffered options did not carry information to influence current choices⁵¹. Therefore, the n -back trials for the
 505 left option may not be the same choice trials as those for the right option, due to the randomized offers.
 506

508 Nutrient model

509 Based on the history model, we further included nutrient-history interaction terms to characterize the
 510 influences of fat and sugar levels on the effects of recent choices and reward outcomes:
 511

$$512 \quad \text{logit}(P_L) = \beta_0 + \beta_1 \times \text{LeftFirst} + \beta_2 \times \text{FatLv} + \beta_3 \times \text{SugarLv}$$

$$513 \quad + \sum_{k=1}^n (\beta_{k+3} \times Cx_k) + \sum_{k=1}^n (\beta_{k+n+3} \times Rx_k)$$

$$514 \quad + \sum_{k=1}^n (\beta_{k+2n+3} \times FC_k) + \sum_{k=1}^n (\beta_{k+3n+3} \times FR_k)$$

$$515 \quad + \sum_{k=1}^n (\beta_{k+4n+3} \times SC_k) + \sum_{k=1}^n (\beta_{k+5n+3} \times SR_k)$$

516 , where FC_n denoted recent high-fat choices and FR_n for high-fat rewarded trials; SC_n denoted recent high-
 517 sugar choices and SR_n for high-sugar rewarded trials. The nutrient-history interaction terms were defined as
 518 follows,
 519

$$520 \quad FC_n = c_{t-n}^L \times \text{FatLv}_{t-n}^L - c_{t-n}^R \times \text{FatLv}_{t-n}^R$$

$$521 \quad SC_n = c_{t-n}^L \times \text{SugarLv}_{t-n}^L - c_{t-n}^R \times \text{SugarLv}_{t-n}^R$$

$$522 \quad FR_n = r_{t-n}^L \times \text{FatLv}_{t-n}^L - r_{t-n}^R \times \text{FatLv}_{t-n}^R$$

$$523 \quad SR_n = r_{t-n}^L \times \text{SugarLv}_{t-n}^L - r_{t-n}^R \times \text{SugarLv}_{t-n}^R$$

524 , where c_{t-n}^L and c_{t-n}^R denoted whether the left or right option was chosen n trials earlier (1, chosen; 0,
 525 unchosen); r_{t-n}^L and r_{t-n}^R denoted whether the left or right option was chosen and was rewarded (1, chosen
 526 and rewarded; 0, otherwise).
 527

530 **Reinforcement learning (RL) models**

531 Standard RL model (Q-learning)

532 We adopted a standard Q-learning algorithm that followed the Rescorla-Wagner learning rule^{1,52}. The reward
533 values (Q_t^i) were set to be 0 for all options initially ($Q_1^i = 0, \forall i \in \{LFLS, HFLS, LFHS, HFHS\}$) and were
534 updated by the reward prediction errors (RPE_t) multiplied by the learning rate $\alpha \in [0,1]$ as follows,
535

$$536 \quad RPE_t = [R_t^i - Q_{t-1}^i], \quad R_t^i = \begin{cases} 1, & \text{if rewarded} \\ 0, & \text{if otherwise} \end{cases}, \quad i \in \{LFLS, HFLS, LFHS, HFHS\}$$

$$537 \quad Q_t^i = Q_{t-1}^i + \alpha \cdot RPE_t$$

538
539 Choices were derived from transforming the value difference δ_t via the softmax function into choice
540 probability π_t^L , which was then dichotomized at 0.5 into binary choice actions A_t^L as below,
541

$$542 \quad \delta_t = Q_t^L - Q_t^R$$

$$543 \quad \pi^L(\delta)_t = \frac{1}{1 + \exp(-\beta \cdot \delta_t)} \in [0,1]$$

$$544 \quad A_t^L = \begin{cases} 1, & \text{if } \pi_t^L > 0.5 \\ Y, & \text{if } \pi_t^L = 0.5 \in \{1,0\}, Y \sim B(1,0.5) \\ 0, & \text{if } \pi_t^L < 0.5 \end{cases}$$

546
547 , where Q_t^L and Q_t^R were the reward values for the left and right option on trial t ; β was the inverse temperature,
548 the sensitivity of choice to value differences.
549

550 Alternative RL models

551 We systematically included differential learning rates and nutrient-specific learning parameters into the RL
552 models. Specifically, we examined 9 combinatorial RL models with 3 differential learning rates (*Standard*,
553 *Asym*, and *Forget*) and 3 nutrient-specific learning parameters (*Standard*, *NutVal*, *Alpha*) ($3 \times 3 = 9$ models)
554 as below.
555

556 *1. Differential learning rates (Standard, Asym, Forget)*

557 We included differential learning rates for rewarded (α^+), unrewarded (α^-), and unoffered (α^0) options to
558 update the reward values as follows,
559

$$560 \quad Q_i(t+1) = Q_i(t) + \alpha \cdot [R_i(t) - Q_i(t)], \quad \alpha = \begin{cases} \alpha^+, & \text{if rewarded} \\ \alpha^-, & \text{if unrewarded} \in [0,1] \\ \alpha^0, & \text{if unoffered} \end{cases}$$

561
562 In the *Standard* model, the agent equally updated both the rewarded and unrewarded option and kept perfect
563 memory for the unoffered option ($\alpha^+ = \alpha^-, \alpha^0 = 0$). In the *Asym* model, the agent updated the rewarded and
564 unrewarded with different learning rates, while keeping perfect memory for the unoffered rewards ($\alpha^+ \neq$
565 $\alpha^-, \alpha^0 = 0$). In contrast, in the *Forget* model, the value of the unoffered rewards decayed due to value
566 forgetting, but the rewarded and unrewarded option were updated equally ($\alpha^+ = \alpha^-, \alpha^0 > 0$).
567

568 *2. Nutrient-specific learning models (NutVal, Alpha)*

569 We examined nutrient preferences by including nutrient-specific values (*NutVal*) or nutrient-specific learning
570 rates (*Alpha*). In the *NutVal* model, the reward values depend on the reward types as follows,
571

$$572 \quad R_i(t) = \begin{cases} 1/(V_F \cdot V_S \cdot V_{FS}), & i = LFLS \\ V_F/(V_F \cdot V_S \cdot V_{FS}), & i = HFLS \\ V_S/(V_F \cdot V_S \cdot V_{FS}), & i = LFHS \\ 1, & i = HFHS \end{cases}$$

573

574 , where V_F , V_S , and V_{FS} are the values of high-fat content, high-sugar content, and their combinations,
 575 respectively, relative to the low-nutrient liquid. We normalized all reward values to $(V_F \cdot V_S \cdot V_{FS})$, to constrain
 576 all reward values between 0 and 1.

577 In the *Alpha* model, higher learning rates are used to update the values for high-nutrient rewards as follow,
 578

$$579 \quad \log \left[\frac{\alpha^+(t)}{1-\alpha^+(t)} \right] = \begin{cases} \alpha_0^+, & i_t = LFLS \\ \alpha_0^+ + \alpha_F, & i_t = HFLS \\ \alpha_0^+ + \alpha_S, & i_t = LFHS \\ \alpha_0^+ + \alpha_{FS}, & i_t = HFHS \end{cases} \in \mathbb{R}, \quad \alpha^+(t) \in [0,1], \quad \forall t \in \mathbb{N}$$

580 , where $\alpha^+(t)$ denoted the learning rate to update the value of the rewarded option on trial t , which was first
 581 transformed from $[0,1]$ to any real number and modified by the high-fat level (α_F), the high-sugar level (α_S),
 582 or their combination (α_{FS}). The logistic transformation ensured that the learning rates are always between 0
 583 and 1.
 584

585
 586 *Nutrient prediction error-RL model (NPE-RL model)*

587 In the NPE-RL model, we decomposed the nutrient-specific values $Q_i(t)$ into components of fat value $Q_i^F(t)$
 588 and sugar value $Q_i^S(t)$,
 589

$$590 \quad Q_i(t) = Q_i^F(t) \cdot Q_i^S(t) \in [0,1], \quad \forall t \in \mathbb{N}$$

591
 592 Importantly, the nutrient prediction errors were computed as the discrepancies between the subjective nutrient
 593 values and the trial-by-trial estimations of the nutrient values as follows,
 594

$$595 \quad NPE_i^F(t) = V_i^F(t) - Q_i^F(t), \quad V_i^F(t) = \begin{cases} 1/v_F, & i(t) = LFLS, LFHS \\ 1, & i(t) = HFLS, HFHS \end{cases}$$

$$596 \quad NPE_i^S(t) = V_i^S(t) - Q_i^S(t), \quad V_i^S(t) = \begin{cases} 1/v_S, & i(t) = LFLS, HFLS \\ 1, & i(t) = LFHS, HFHS \end{cases}$$

597 , where NPE_i^F and NPE_i^S denoted the fat and sugar prediction errors for the chosen reward on trial t , $i(t)$. v_i^F
 598 and v_i^S were the subjective values for fat and sugar, and Q_i^F and Q_i^S were the current values of fat and sugar
 599 components for reward i , respectively. The nutrient values were independently updated by corresponding
 600 nutrient prediction errors,
 601

$$602 \quad Q_i^F(t+1) = Q_i^F(t) + \alpha^+ \cdot NPE_i^F(t)$$

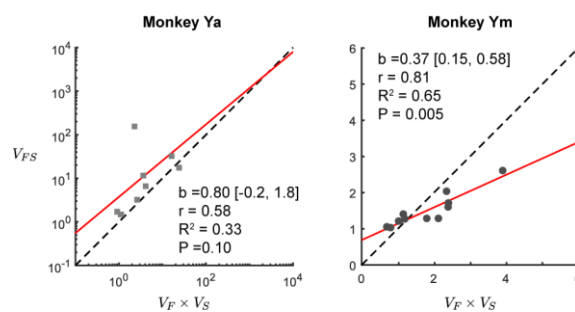
$$603 \quad Q_i^S(t+1) = Q_i^S(t) + \alpha^+ \cdot NPE_i^S(t)$$

604 , where $Q_i^F(t+1)$ and $Q_i^S(t+1)$ are the updated fat and sugar values, each was updated by the previous fat
 605 and sugar values, $Q_i^F(t)$ and $Q_i^S(t)$, by the NPEs for fat and sugar discounted by the learning rate $\alpha^+ \in [0,1]$.
 606
 607
 608
 609
 610

611 **Acknowledgements.** We thank Wolfram Schultz and his group for support; Putu Khorisantono for discussions;
 612 Christina Thompson and Aled David for animal care; Polly Taylor for anesthesia; Henri Bertrand for veterinary
 613 care. This work was funded by the Wellcome Trust and the Royal Society (Sir Henry Dale Fellowship
 614 206207/Z/17/Z to F.G.). F.-Y.H. was supported by a Fellowship from the Taiwan Ministry of Education.

615 SUPPLEMENTARY FIGURES

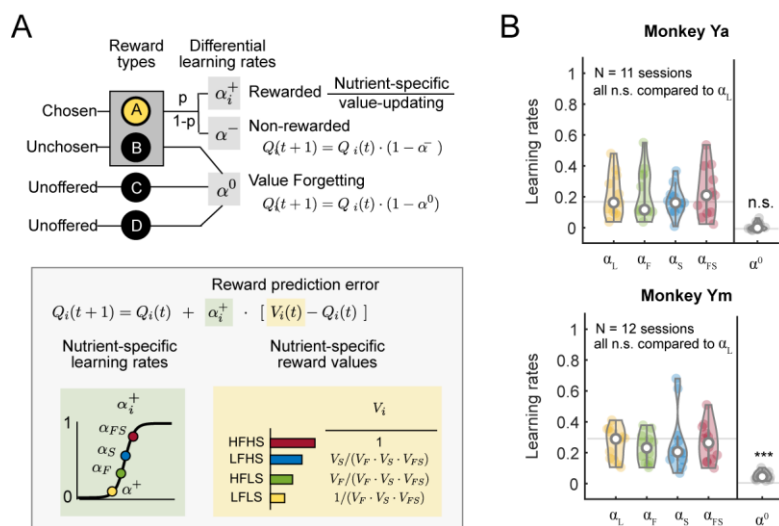
616
617



618
619

620 **Fig. S1. Fat-sugar value interactions.** The reward values of HFHS (V_{FS}) were plotted against the values
621 predicted by the multiplications of the fat values (V_F) and the sugar values (V_S) across sessions. The unity lines
622 (dashed) indicated the independence of the fat values and the sugar values, estimated by the nutrient-sensitive
623 reinforcement learning models. b = slope [95% confidence interval].

624
625
626
627
628



629
630

631 **Fig. S2. Nutrient-specific learning rates.** A) Nutrient-specific learning rate model (NutAlphaVal-Forget
632 model) architecture. The reward values were updated based on nutrient-specific learning rates (α_i^+) in addition
633 to the nutrient-specific values (V_i^t). Values of the unchosen and unoffered rewards decayed according to the
634 forgetting factor α^0 , as in the main nutrient value RL model (Figure 4A). B) Nutrient-specific learning rates
635 and forgetting factors. Learning rates for HFLS (α_F), LFHS (α_S), and HFHS (α_{FS}) were all compared to the
636 baseline learning rates for LFLS (α_L); the forgetting factors were tested against perfect value memory ($\alpha^0 =$
637 0). Wilcoxon signed-rank test.

References

- 638
639
640 1. Sutton, R.S. & Barto, A.G. *Reinforcement Learning* (MIT Press, Cambridge, MA, 1998).
641 2. Lak, A., Stauffer, W.R. & Schultz, W. Dopamine prediction error responses integrate subjective value
642 from different reward dimensions. *P Natl Acad Sci USA* **111**, 2343-2348 (2014).
643 3. Stauffer, W.R., Lak, A. & Schultz, W. Dopamine Reward Prediction Error Responses Reflect
644 Marginal Utility. *Current biology : CB* (2014).
645 4. Schultz, W., Dayan, P. & Montague, P.R. A neural substrate of prediction and reward. *Science* **275**,
646 1593-1599 (1997).
647 5. Hollerman, J.R. & Schultz, W. Dopamine neurons report an error in the temporal prediction of reward
648 during learning. *Nature neuroscience* **1**, 304-309 (1998).
649 6. Lau, B. & Glimcher, P.W. Value representations in the primate striatum during matching behavior.
650 *Neuron* **58**, 451-463 (2008).
651 7. Samejima, K., Ueda, Y., Doya, K. & Kimura, M. Representation of action-specific reward values in
652 the striatum. *Science* **310**, 1337-1340 (2005).
653 8. Grabenhorst, F., Baez-Mendoza, R., Genest, W., Deco, G. & Schultz, W. Primate Amygdala Neurons
654 Simulate Decision Processes of Social Partners. *Cell* **177**, 986-998 e915 (2019).
655 9. Costa, V.D., Mitz, A.R. & Averbeck, B.B. Subcortical Substrates of Explore-Exploit Decisions in
656 Primates. *Neuron* **103**, 533-545 e535 (2019).
657 10. Grabenhorst, F., Tsutsui, K.I., Kobayashi, S. & Schultz, W. Primate prefrontal neurons signal
658 economic risk derived from the statistics of recent reward experience. *Elife* **8** (2019).
659 11. Tsutsui, K., Grabenhorst, F., Kobayashi, S. & Schultz, W. A dynamic code for economic object
660 valuation in prefrontal cortex neurons. *Nature Communications* **7**, 12554 (2016).
661 12. Lee, D., Seo, H. & Jung, M.W. Neural basis of reinforcement learning and decision making. *Annu Rev*
662 *Neurosci* **35**, 287-308 (2012).
663 13. Seo, M., Lee, E. & Averbeck, B.B. Action selection and action value in frontal-striatal circuits. *Neuron*
664 **74**, 947-960 (2012).
665 14. Cui, Z.W., Wang, Z.L., Shao, Q., Raubenheimer, D. & Lu, J.Q. Macronutrient signature of dietary
666 generalism in an ecologically diverse primate in the wild. *Behavioral Ecology* **29**, 804-813 (2018).
667 15. Cui, Z.W., *et al.* Living near the limits: Effects of interannual variation in food availability on diet and
668 reproduction in a temperate primate, the Taihangshan macaque (*Macaca mulatta tcheliensis*). *Am J*
669 *Primatol* **82** (2020).
670 16. Yang, Y., *et al.* Cafeteria-style feeding trials provide new insights into the diet and nutritional
671 strategies of the black snub-nosed monkey (*Rhinopithecus strykeri*): Implications for conservation.
672 *Am J Primatol* **82**, e23108 (2020).
673 17. Chivers, D.J. Measuring food intake in wild animals: primates. *P Nutr Soc* **57**, 321-332 (1998).
674 18. Ma, C.Y., Liao, J.C. & Fan, P.F. Food selection in relation to nutritional chemistry of Cao Vit gibbons
675 in Jingxi, China. *Primates* **58**, 63-74 (2017).
676 19. Takahashi, M.Q., Rothman, J.M., Raubenheimer, D. & Cords, M. Dietary generalists and nutritional
677 specialists: Feeding strategies of adult female blue monkeys (*Cercopithecus mitis*) in the Kakamega
678 Forest, Kenya. *Am J Primatol* **81**, e23016 (2019).
679 20. Drewnowski, A. & Almiron-Roig, E. Human Perceptions and Preferences for Fat-Rich Foods. in *Fat*
680 *Detection: Taste, Texture, and Post Ingestive Effects* (ed. J.P. Montmayeur & J. le Coutre) (Boca Raton
681 (FL), 2010).
682 21. van der Klaauw, A.A., *et al.* Divergent effects of central melanocortin signalling on fat and sucrose
683 preference in humans. *Nature Communications* **7**, 13055 (2016).
684 22. Huang, F.-Y., Sutcliffe, M.P.F. & Grabenhorst, F. Preferences for nutrients and sensory food qualities
685 identify biological sources of economic values in monkeys. *Proc Natl Acad Sci U S A* (2021).
686 23. Pastor-Bernier, A., Volkmann, K., Stasiak, A., Grabenhorst, F. & Schultz, W. Experimentally revealed
687 stochastic preferences for multicomponent choice options. *J Exp Psychol Anim Learn Cogn* **46**, 367-
688 384 (2020).
689 24. Amato, K.R. & Garber, P.A. Nutrition and foraging strategies of the black howler monkey (*Alouatta*
690 *pigra*) in Palenque National Park, Mexico. *Am J Primatol* **76**, 774-787 (2014).
691 25. Lau, B. & Glimcher, P.W. Dynamic response-by-response models of matching behavior in rhesus
692 monkeys. *J Exp Anal Behav* **84**, 555-579 (2005).

- 693 26. Corrado, G.S., Sugrue, L.P., Seung, H.S. & Newsome, W.T. Linear-Nonlinear-Poisson models of
694 primate choice dynamics. *J Exp Anal Behav* **84**, 581-617 (2005).
- 695 27. Kennerley, S.W., Walton, M.E., Behrens, T.E., Buckley, M.J. & Rushworth, M.F. Optimal decision
696 making and the anterior cingulate cortex. *Nature neuroscience* **9**, 940-947 (2006).
- 697 28. Rothenhoefer, K.M., Hong, T., Alikaya, A. & Stauffer, W.R. Rare rewards amplify dopamine
698 responses. *Nature neuroscience* **24**, 465-469 (2021).
- 699 29. Raghuraman, A.P. & Padoa-Schioppa, C. Integration of multiple determinants in the neuronal
700 computation of economic values. *The Journal of neuroscience : the official journal of the Society for*
701 *Neuroscience* **34**, 11583-11603 (2014).
- 702 30. Pastor-Bernier, A., Plott, C.R. & Schultz, W. Monkeys choose as if maximizing utility compatible
703 with basic principles of revealed preference theory. *P Natl Acad Sci USA* **114**, E1766-E1775 (2017).
- 704 31. Padoa-Schioppa, C. & Assad, J.A. Neurons in the orbitofrontal cortex encode economic value. *Nature*
705 **441**, 223-226 (2006).
- 706 32. Yarmolinsky, D.A., Zuker, C.S. & Ryba, N.J. Common sense about taste: from mammals to insects.
707 *Cell* **139**, 234-244 (2009).
- 708 33. Carreiro, A.L., *et al.* The Macronutrients, Appetite, and Energy Intake. *Annu Rev Nutr* **36**, 73-103
709 (2016).
- 710 34. Rolls, E.T. The texture and taste of food in the brain. *J Texture Stud* **51**, 23-44 (2020).
- 711 35. Simpson, S.J. & Raubenheimer, D. The power of protein. *The American journal of clinical nutrition*
712 **112**, 6-7 (2020).
- 713 36. Rangel, A. Regulation of dietary choice by the decision-making circuitry. *Nature neuroscience* **16**,
714 1717-1724 (2013).
- 715 37. Raubenheimer, D. Toward a quantitative nutritional ecology: the right-angled mixture triangle. *Ecol*
716 *Monogr* **81**, 407-427 (2011).
- 717 38. Raubenheimer, D., Machovsky-Capuska, G.E., Chapman, C.A. & Rothman, J.M. Geometry of
718 nutrition in field studies: an illustration using wild primates. *Oecologia* **177**, 223-234 (2015).
- 719 39. Tsuji, Y. & Takatsuki, S. Interannual Variation in Nut Abundance Is Related to Agonistic Interactions
720 of Foraging Female Japanese Macaques (*Macaca fuscata*). *Int J Primatol* **33**, 489-512 (2012).
- 721 40. Hall, K.D., *et al.* Ultra-Processed Diets Cause Excess Calorie Intake and Weight Gain: An Inpatient
722 Randomized Controlled Trial of Ad Libitum Food Intake. *Cell Metab* **30**, 67-77 e63 (2019).
- 723 41. Martinez Steele, E., Raubenheimer, D., Simpson, S.J., Baraldi, L.G. & Monteiro, C.A. Ultra-processed
724 foods, protein leverage and energy intake in the USA. *Public Health Nutr* **21**, 114-124 (2018).
- 725 42. Jensen-Cody, S.O., *et al.* FGF21 Signals to Glutamatergic Neurons in the Ventromedial Hypothalamus
726 to Suppress Carbohydrate Intake. *Cell Metab* **32**, 273-286 e276 (2020).
- 727 43. Keramati, M. & Gutkin, B. Homeostatic reinforcement learning for integrating reward collection and
728 physiological stability. *Elife* **3** (2014).
- 729 44. Alonso-Alonso, M., *et al.* Food reward system: current perspectives and future research needs. *Nutr*
730 *Rev* **73**, 296-307 (2015).
- 731 45. Murray, E.A. & Rudebeck, P.H. Specializations for reward-guided decision-making in the primate
732 ventral prefrontal cortex. *Nature reviews. Neuroscience* **19**, 404-417 (2018).
- 733 46. Rolls, E.T., Mills, T., Norton, A.B., Lazidis, A. & Norton, I.T. The Neuronal Encoding of Oral Fat by
734 the Coefficient of Sliding Friction in the Cerebral Cortex and Amygdala. *Cerebral cortex* **28**, 4080-
735 4089 (2018).
- 736 47. Grabenhorst, F., Hernadi, I. & Schultz, W. Prediction of economic choice by primate amygdala
737 neurons. *Proc Natl Acad Sci U S A* **109**, 18950-18955 (2012).
- 738 48. Murray, E.A. & Rudebeck, P.H. The drive to strive: goal generation based on current needs. *Frontiers*
739 *in neuroscience* **7**, 112 (2013).
- 740 49. Gallistel, C.R., Fairhurst, S. & Balsam, P. The learning curve: implications of a quantitative analysis.
741 *P Natl Acad Sci USA* **101**, 13124-13131 (2004).
- 742 50. Herrnstein, R.J. Relative and absolute strength of response as a function of frequency of reinforcement.
743 *J Exp Anal Behav* **4**, 267-272 (1961).
- 744 51. Wittmann, M.K., *et al.* Global reward state affects learning and activity in raphe nucleus and anterior
745 insula in monkeys. *Nature Communications* **11**, 3771 (2020).
- 746 52. Rescorla, R.A. & Wagner, A.R. A theory of Pavlovian conditioning: Variations in the effectiveness of
747 reinforcement and nonreinforcement. in *Classical Conditioning II: Current Research and Theory* (ed.
748 A.H. Black & W.F. Prokasy) 64-99 (Appleton Century Crofts, New York, 1972).