

Nutrient-sensitive reinforcement learning in monkeys

Fei-Yang Huang (黃飛揚)^{1,2} and Fabian Grabenhorst^{1,2*}

¹Department of Experimental Psychology, University of Oxford, Oxford OX1 3TA, UK

²Department of Physiology, Development and Neuroscience, University of Cambridge, Cambridge CB2 3DY, UK

*For correspondence: fabian.grabenhorst@psy.ox.ac.uk

Published in: *The Journal of Neuroscience*

Abbreviated title: Nutrient-sensitive reinforcement learning

8 figures

1 table

249 words in Abstract

650 words in Introduction

1,495 words in Discussion

Competing interests: The authors declare that they have no competing interests.

Acknowledgments. We thank Wolfram Schultz and his group for support and discussions; Putu Khorisantono for discussions; Christina Thompson and Aled David for animal care; Polly Taylor for anesthesia; Henri Bertrand for veterinary care. This work was funded by the Wellcome Trust and the Royal Society (Sir Henry Dale Fellowships 206207/Z/17/Z and 206207/Z/17/A to F.G.). F.-Y.H. was supported by a Fellowship from the Taiwan Ministry of Education. This research was funded in whole, or in part, by the Wellcome Trust. For the purpose of Open Access, the author has applied a CC BY public copyright license to any Author Accepted Manuscript version arising from this submission.

In Reinforcement Learning (RL), animals choose by assigning values to options and learn by updating these values from reward outcomes. This framework has been instrumental in identifying fundamental learning variables and their neuronal implementations. However, canonical RL models do not explain how reward values are constructed from biologically critical intrinsic reward components, such as nutrients. From an ecological perspective, animals should adapt their foraging choices in dynamic environments to acquire nutrients that are essential for survival. Here, to advance the biological and ecological validity of RL models, we investigated how (male) monkeys adapt their choices to obtain preferred nutrient rewards under varying reward probabilities. We found that the rewards' nutrient composition strongly influenced learning and choices. The animals' preferences for specific nutrients (sugar, fat) affected how they adapted to changing reward probabilities: the history of recent rewards influenced monkeys' choices more strongly if these rewards contained the monkey's preferred nutrients ('nutrient-specific reward history'). The monkeys also chose preferred nutrients even when they were associated with lower reward probability. A nutrient-sensitive RL model captured these processes: it updated the values of individual sugar and fat components of expected rewards based on experience and integrated them into subjective values that explained the monkeys' choices. Nutrient-specific reward prediction errors guided this value-updating process. Our results identify nutrients as important reward components that guide learning and choice by influencing the subjective value of choice options. Extending RL models with nutrient-value functions may enhance their biological validity and uncover nutrient-specific learning and decision variables.

Significance statement

Reinforcement learning (RL) is an influential framework that formalizes how animals learn from experienced rewards. Although 'reward' is a foundational concept in RL theory, canonical RL models cannot explain how learning depends on specific reward properties, such as nutrients. Intuitively, learning should be sensitive to the reward's nutrient components, to benefit health and survival. Here we show that the nutrient (fat, sugar) composition of rewards affects monkeys' choices and learning in an RL paradigm, and that key learning variables including 'reward history' and 'reward prediction error' should be modified with nutrient-specific components to account for monkeys' behavior in our task. By incorporating biologically critical nutrient rewards into the RL framework our findings help advance the ecological validity of RL models.

INTRODUCTION

According to Reinforcement Learning (RL) theory, animals choose by assigning values to options and learn by updating these values from experienced rewards (Sutton and Barto, 1998). This framework has been instrumental in identifying learning and decision variables that explain animal behavior, including object values and action values, which guide choices, and reward prediction errors, which update values from experienced outcomes. Physical implementations of these concepts have been discovered in neurons of the primate dopamine system, striatum, amygdala, and frontal cortex (Schultz et al., 1997; Samejima et al., 2005; Lau and Glimcher, 2008; So and Stuphorn, 2010; Lee et al., 2012; Seo et al., 2012; Tsutsui et al., 2016; Costa et al., 2019; Grabenhorst et al., 2019b). One important factor that limits the biological validity of canonical RL models is that they do not explain how learning and choice depend on the composition of experienced rewards. Nutrients, for example, are biologically critical intrinsic components of food rewards that engage dedicated sensory and physiological mechanisms and are essential for survival (Carreiro et al., 2016; Rolls, 2020; Simpson and Raubenheimer, 2020). Thus, investigating how nutrients influence learning and choice could enhance the biological validity of the RL framework.

From an ecological perspective, an animal's ability to adapt its food choices to changing nutrient availabilities determines its survival and long-term health (Simpson and Raubenheimer, 2012). Indeed, foraging monkeys adapt their feeding patterns to uncertainty imposed by regional and seasonal food variations (Cui et al., 2018; Cui et al., 2020), for example adjusting their diet based on the availability of nutritious foods (seeds, nuts). Primates, including humans, also exhibit subjective preferences for specific nutrients and related sensory food qualities (van der Klaauw et al., 2016; Ma et al., 2017; Takahashi et al., 2019; Huang et al., 2021). Thus, foraging animals consider both food nutrient composition and the probability of obtaining the food. However, the mechanisms underlying such adaptive, nutrient-sensitive food choices remain poorly understood.

Recent proposals in RL theory extended the reward concept to outcomes under homeostatic regulation (Keramati and Gutkin, 2014) and distinguished hedonic reward components from their post-ingestive consequences (Dayan, 2022). Generally, the importance of linking biological reward components to learning and choice theories is increasingly recognized (Rangel, 2013; Suzuki et al., 2017; Averbeck and Murray, 2020; de Araujo et al., 2020; Stuphorn, 2021). Yet, the field lacks experimental data on how nutrient-reward components influence animals' choices in formal RL paradigms.

We recently showed that monkeys in economic choice tasks exhibit individual preferences for specific nutrients and related sensory food qualities (e.g., oral texture) (Huang et al., 2021). Monkeys' preferences were well-described by subjective 'nutrient-value functions' that linked reward nutrient composition to choices. In simulations, we showed that incorporating nutrient-value functions into canonical RL models optimized model performance when preferences prioritized specific nutrients (Huang et al., 2021). Here, we extended this approach to examine the behavior of rhesus monkeys (*Macaca mulatta*) in a dynamic foraging task involving choices between rewards with different nutrient (fat, sugar) components under varying reward probabilities. Our task captured features of food choice in ecological foraging contexts (Cui et al., 2018; Cui et al., 2020): the animals could adapt their choices to changing reward availabilities to obtain preferred nutrients or choose less-preferred but more available alternatives.

Previous studies showed that macaques track changing reward probabilities based on the history of recent choices and rewards (Corrado et al., 2005; Lau and Glimcher, 2005; Kennerley et al., 2006; Lee et al., 2012; Grabenhorst et al., 2019a). The influence of reward history on choice is typically modeled by linking subjectively weighted recent rewards to current choices using logistic regression (Lau and Glimcher, 2005) and by dynamically updating expected values from experienced rewards using RL models (Sutton and Barto, 1998). We followed these approaches and examined whether, in a dynamic learning task, monkeys assigned higher value to recent rewards based on the rewards' nutrient components, and whether nutrient-sensitive RL models could account for this process.

MATERIALS AND METHODS

Animals. Two adult male rhesus macaques (*Macaca mulatta*) were trained in the study: monkey Ya (weight during the experiments: 17-19 kg, age: 6 years) and monkey Ym (12-13 kg, age 6 years).

The animals were trained and tested approximately one to two hours per day and five days per week for 6 months, interrupted by regular week-long testing breaks. Both monkeys participated in a related nutrient-choice study using the same dairy-based nutrient rewards as in this study (Huang et al., 2021). Thus, the different sensitivities of learning from the received nutrients in the probabilistic learning task can be linked to the nutrient preferences of the same monkeys in economic choice tasks without learning. The animals were on a standard diet for laboratory macaques, composed of high-protein dry pellets (% calories provided by protein: 30.36%, fat: 13.29%, carbohydrates: 56.34%), dried fruits, seeds, nuts, and fresh fruits and vegetables. We monitored the monkeys' health condition and body weights to ensure their welfare after introducing high-calorie rewards. No effects of these rewards on the animals' health were observed. Each testing day, the animals had free access to the standard diet before and after the experiments and received their main liquid intake in the laboratory. The animals' body weights increased as expected for growing animals.

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

Experimental Design

Nutrient rewards. We prepared nutrient-controlled liquids with 2x2 fat and sugar levels to examine whether fat and sugar biased learning from reward outcomes (**Fig. 1A, B**; LFLS: low-fat low-sugar; HFLS: high-fat low-sugar; LFHS: low-fat high-sugar; HFHS: high-fat high-sugar). The liquids were matched in flavor (peach or blackcurrant, varied in different sessions), temperature, protein, salt and other ingredients (see **Table 1** for detailed liquid compositions). We used commercial skimmed milk and whole milk (British skimmed milk and whole milk, Sainsbury's Supermarkets Ltd., UK) as baseline low-fat and high-fat liquids and flavored the liquids with fruit juice to increase palatability. The energy content of the liquid rewards (**Table 1**) was estimated based on metabolizable energy using Atwater general energy conversion factors (protein: 4.0 kcal/g; fat: 9.0 kcal/g; carbohydrate: 4.0 kcal/g), following the definition of 'foraging efficiency' in Optimal Foraging Theory (McNamara and Houston, 1997).

Nutrient foraging task. The four nutrient reward types were associated with four untrained visual cues in each session. When a trial started, the monkeys were first presented with two of the four visual cues, shown on a horizontally mounted touch monitor. They then made a choice between the two cues by touching one of two blue rectangle target stimuli, shown below the cues. Following the choice, they received either a large amount ('rewarded trials', 0.5 mL) or a small amount ('non-rewarded trials', 0.3 mL) of the cue-associated liquid depending on its prespecified reward probability (P) (**Fig. 1A**). When the session started, two of the rewards (LFLS/HFHS or LFHS/HFLS) were offered in high reward probabilities ($P = 0.8$), and the other two rewards in low reward probabilities ($P = 0.2$) (**Fig. 1C**, block A or block B). The reward probabilities were reversed every 100 trials ($P = 0.2 \rightarrow 0.8$; $P = 0.8 \rightarrow 0.2$) (**Fig. 1C, D**).

Data Analysis

All data were analyzed using Matlab 2017 (Mathworks).

Learning curves. Learning curves were plotted by aligning reward-specific choices to the probability reversal trials. In particular, based on the probability before and after reversals, we grouped these curves into incremental ($P=0.2 \rightarrow P=0.8$) and decremental ($P=0.8 \rightarrow P=0.2$, not shown) learning curves, and plotted the incremental curves in **Fig. 2A**. Two-sample t-test were performed to compare

choice probabilities averaged across the last 10 plotted trials after learning for each reward (70 to 80 trials).

Learning latency. The learning latency was defined as the number of trials between the first behavioral change point after probability reversals (**Fig. 2B**). The behavioral change points were identified as the significant change points of cumulative choice slopes (Gallistel et al., 2004), based on two-sample t-test with criterion $P < 0.05$.

Averaged choice probability. We compared the averaged choice probability for each reward to indicate the reward preference (**Fig. 2C**). To cancel out the influence of reward probability on choices, we truncated the final trials in unbalanced sessions and computed the averaged choice probability across the same total number of high-probability and low-probability blocks for each reward.

Logistic regression analysis

History model

We used multiple logistic regression (*fitglm* function, Matlab) to model choices based on recent choices and reward outcomes as follows (**Fig. 3A, B**),

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

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

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

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

Notably, the history regressors for each option were indexed only when the same option was offered by assuming the unoffered options in a past trial did not influence current choices (Wittmann et al., 2020). Therefore, due to the random sampling of offered options, the n -back trials for the left option may not be the same trials as those for the right option.

Nutrient model

Based on the history model, we further included nutrient-history interaction terms in the regression model to characterize the individual contribution of fat and sugar content on the effects of recent choices and reward outcomes (**Fig. 3C-G**):

$$\begin{aligned} \text{logit}(P_L) = & \beta_0 + \beta_1 \times \text{LeftFirst} + \beta_2 \times \text{FatLv} + \beta_3 \times \text{SugarLv} \\ & + \sum_{k=1}^n (\beta_{k+3} \times Cx_k) + \sum_{k=1}^n (\beta_{k+n+3} \times Rx_k) \\ & + \sum_{k=1}^n (\beta_{k+2n+3} \times FC_k) + \sum_{k=1}^n (\beta_{k+3n+3} \times FR_k) \\ & + \sum_{k=1}^n (\beta_{k+4n+3} \times SC_k) + \sum_{k=1}^n (\beta_{k+5n+3} \times SR_k) \end{aligned}$$

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

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

$$\begin{aligned}
SC_n &= c_{t-n}^L \times SugarLv_{t-n}^L - c_{t-n}^R \times SugarLv_{t-n}^R \\
FR_n &= r_{t-n}^L \times FatLv_{t-n}^L - r_{t-n}^R \times FatLv_{t-n}^R \\
SR_n &= r_{t-n}^L \times SugarLv_{t-n}^L - r_{t-n}^R \times SugarLv_{t-n}^R
\end{aligned}$$

, where c_{t-n}^L and c_{t-n}^R indicated whether the left or right option was chosen n trials earlier (1, chosen; 0, unchosen); r_{t-n}^L and r_{t-n}^R denoted whether the left or right option was chosen and was rewarded (1, chosen and rewarded; 0, otherwise). Model comparison was performed based on the Akaike Information Criterion (AIC) between the History model and the Nutrient model; these models were matched in history lengths up to 10 trials in the past (**Fig. 3E**).

Reinforcement learning (RL) models

Standard RL model (Q-learning with binary reward outcomes)

We adopted a standard Q-learning algorithm that followed the Rescorla-Wagner learning rule with binary reward outcomes (Rescorla and Wagner, 1972; Sutton and Barto, 1998). The initial reward values (Q_t^i) were set to be 0 for all options ($Q_1^i = 0, \forall i \in \{LFLS, HFSL, LFHS, HFHS\}$) and were updated by the reward prediction errors (RPE_t) scaled by the learning rate $\alpha \in [0,1]$ as follows,

$$\begin{aligned}
RPE_t &= [R_t^i - Q_{t-1}^i], \quad R_t^i = \begin{cases} 1, & \text{if rewarded} \\ 0, & \text{if otherwise} \end{cases}, \quad i \in \{LFLS, HFSL, LFHS, HFHS\} \\
Q_t^i &= Q_{t-1}^i + \alpha \cdot RPE_t
\end{aligned}$$

Choices were predicted by first transforming the left-right value difference δ_t via the softmax function into choice probability π_t^L

$$\pi^L(\delta)_t = \frac{1}{1 + \exp(-\beta \cdot \delta_t - \beta_0)} \in [0,1], \quad \delta_t = Q_t^L - Q_t^R$$

, where Q_t^L and Q_t^R were the expected values for the left and right option on trial t , β was the inverse temperature indicating the sensitivity of choice to value differences, and β_0 was the side bias independent of left-right action values. The three free parameters, learning rate (α), inverse temperature (β), and side-bias intercept (β_0) were fitted based on maximum likelihood estimation (*fminsearch* function, Matlab) using the likelihood function $\mathcal{L}(\theta)$ below ($I_t^L = 1, I_t^R = 0$ for left choices; $I_t^L = 0, I_t^R = 1$ for right choices).

$$\begin{aligned}
\mathcal{L}(\theta) &= \prod_{t \in T} l_t(\theta) = \prod_{t \in T} \left[\frac{\exp(\beta \cdot Q_t^L) \times I_t^L + \exp(\beta \cdot Q_t^R - \beta_0) \times I_t^R}{\exp(\beta \cdot Q_t^L) + \exp(\beta \cdot Q_t^R - \beta_0)} \right], \quad \theta = \{\alpha, \beta, \beta_0\} \\
\hat{\theta} &= \underset{\theta}{\operatorname{argmax}} \mathcal{L}(\theta)
\end{aligned}$$

Alternative RL models

We systematically included differential learning rates and nutrient-specific learning parameters in the RL models. Specifically, we examined nine combinatorial RL models with three model complexities (*Basic*, *Asym*, and *Forget*) and three nutrient-specific learning parameters (*Binary*, *NutVal*, *Alpha*) as described below.

1. RL model complexity (*Basic*, *Asym*, *Forget* models)

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

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

291 In the *Basic* models, the agent equally updated both the rewarded and unrewarded options and
 292 kept perfect memory for the unoffered option ($\alpha^+ = \alpha^-, \alpha^0 = 0$). In the *Asym* model, the agent
 293 updated the rewarded and unrewarded with different learning rates, while keeping perfect memory
 294 for the unoffered rewards ($\alpha^+ \neq \alpha^-, \alpha^0 = 0$). In the *Forget* model, the value of the unoffered rewards
 295 decayed due to value forgetting, but the rewarded and unrewarded options were updated equally
 296 ($\alpha^+ = \alpha^-, \alpha^0 > 0$).

297 2. Nutrient-specific learning parameters (*Binary*, *NutVal*, *NutValAlpha*)

298 We examined nutrient preferences by including either nutrient-specific values only (*NutVal* models)
 299 or additional nutrient-specific learning rates (*NutValAlpha* models). In the *NutVal* models (**Fig. 4B**),
 300 the reward values depend on the reward types as follows,
 301

$$302 \quad Q_i(t+1) = Q_i(t) + \alpha \cdot [V_i(t) \times R(t) - Q_i(t)]$$

$$303 \quad V_i(t) = \frac{V_F^{I_F(t)} \cdot V_S^{I_S(t)} \cdot V_{FS}^{I_F(t) \times I_S(t)}}{V_F \cdot V_S \cdot V_{FS}}, \quad R(t) = \begin{cases} 1 & , \text{if rewarded} \\ d & , \text{if not rewarded} \end{cases} \quad d \in [0,1]$$

$$304 \quad I_F(t) = \begin{cases} 1, & i(t) = HFHS, HFHS \\ 0, & i(t) = LFHS, LFHS \end{cases}, \quad I_S(t) = \begin{cases} 1, & i(t) = LFHS, HFHS \\ 0, & i(t) = LFHS, LFHS \end{cases}$$

305 , where V_F , V_S , and V_{FS} are fixed animal-specific values for fat, sugar, and their combinations,
 306 respectively, in addition to the low-nutrient baseline ingredients. Therefore, the experienced reward
 307 values $V_i(t)$ were computed based on subjective preferences for the fat level (I_F) and sugar levels
 308 (I_S) of the received reward. For computational simplicity, we constrained all reward values between
 309 0 and 1 by normalizing them to the value of HFHS ($V_F \cdot V_S \cdot V_{FS}$). When the animals received a small
 310 reward ('non-rewarded trials'), the reward values were discounted by a constant $R(t) = d \in [0,1]$,
 311 which scales the experienced reward values according to the large or small reward amounts. Thus,
 312 the model included free parameters for subjective values of fat, sugar, fat-sugar interaction, and for
 313 discounting of low reward amounts.

314 In the *NutValAlpha* models (**Fig. 4C**), higher learning rates are used to update the values for high-
 315 nutrient rewards as follows,
 316

$$317 \quad \log \left[\frac{\alpha(t)}{1-\alpha(t)} \right] = \alpha_0 + \alpha_F \cdot I_F(t) + \alpha_S \cdot I_S(t) + \alpha_{FS} \times [I_F(t) \cdot I_S(t)] \in \mathbb{R}, \quad \alpha(t) \in [0,1], \quad \forall t \in \mathbb{N}$$

318 , where $\alpha^+(t)$ denotes the learning rate to update the value of the rewarded option on trial t , which
 319 was first transformed from $[0,1]$ to the real domain and modified by the high-fat level (α_F), the high-
 320 sugar level (α_S), or their combination (α_{FS}), depending on the fat levels (I_F) and sugar levels (I_S). The
 321 logistic transformation ensured that the learning rates were always between 0 and 1.

322 Under these specifications, the *Standard RL model* is, therefore, equivalent to the *Basic Binary*
 323 model, and the best-fitting model in the main text refers to the *NutVal-Forget* model (**Fig. 4**).
 324

325 Energy RL model

326 In the Energy model (**Fig. 4G**), the reward values were determined solely by their energy content:
 327

$$328 \quad Q_i(t+1) = Q_i(t) + \alpha \cdot [V_i(t) \times R(t) - Q_i(t)]$$

$$329 \quad V_i(t) = \begin{cases} \frac{1}{V_E} & , i = LFHS \\ 0.5 & , i = HFHS, LFHS \\ 1 & , i = HFHS \end{cases}, \quad R(t) = \begin{cases} 1 & , \text{if rewarded} \\ d & , \text{if not rewarded} \end{cases} \quad d \in [0,1]$$

331 , where V_E denotes the subjective values for the high energy content of HFHS reward. Values for
 332 rewards with the middle energy levels, LFHS and HFHS, were both $\frac{1}{2} V_E$ and all values were
 333 normalized to V_E between 0 and 1.
 334

335 Object value RL model (ObjVal model)

336 In the *ObjVal* model (**Fig. 4G**), the animals learn from stimulus-specific values that are free
 337 parameters fixed in each session. The value function is as below, with all values normalized to V_{HFHS}
 338 to be constrained between 0 and 1.

$$339 \quad V_i(t) = \begin{cases} V_{LFLS}/V_{HFHS}, & i(t) = LFLS \\ V_{HFHS}/V_{HFHS}, & i(t) = HFHS \\ V_{LFHS}/V_{HFHS}, & i(t) = LFHS \\ 1, & i(t) = HFHS \end{cases}$$

340 Other specifications are similar to the *NutVal* models, including the learning rate, value-forgetting
 341 rate, and the discount factor.

342 Nutrient prediction error-RL model (Nut-RPE model)

343 In the *NutRPE* model (**Fig. 7**), we decomposed the nutrient-specific values $Q_i(t)$ into components of
 344 fat value $Q_i^F(t)$ and sugar value $Q_i^S(t)$,

$$347 \quad Q_i(t) = Q_i^F(t) \cdot Q_i^S(t) \in [0,1], \quad \forall t \in \mathbb{N}$$

348 Importantly, these value components were updated by nutrient prediction errors (NPEs) that were
 349 computed as the discrepancies between the experienced nutrient values and the trial-by-trial
 350 expected nutrient values as follows,
 351
 352

$$353 \quad NPE_i^F(t) = V_i^F(t) \times R(t) - Q_i^F(t), \quad V_i^F(t) = \begin{cases} \frac{1}{v_F}, & i(t) = LFLS, LFHS \\ \frac{1}{v_F}, & i(t) = HFHS, HFHS \end{cases}$$

$$354 \quad NPE_i^S(t) = V_i^S(t) \times R(t) - Q_i^S(t), \quad V_i^S(t) = \begin{cases} \frac{1}{v_S}, & i(t) = LFLS, HFHS \\ \frac{1}{v_S}, & i(t) = LFHS, HFHS \end{cases}$$

$$356 \quad R(t) = \begin{cases} 1, & \text{if rewarded} \\ d, & \text{if not rewarded} \end{cases} \quad d \in [0,1]$$

357 , where NPE_i^F and NPE_i^S denoted the fat and sugar prediction errors for the chosen reward on trial
 358 t , $i(t)$. v_i^F and v_i^S were the fixed, animal-specific values for fat and sugar, and Q_i^F and Q_i^S were the
 359 current expected values of fat and sugar components for reward i , respectively. The nutrient values
 360 were separately updated by corresponding nutrient prediction errors,
 361
 362

$$363 \quad Q_i^F(t+1) = Q_i^F(t) + \alpha \cdot NPE_i^F(t)$$

$$364 \quad Q_i^S(t+1) = Q_i^S(t) + \alpha \cdot NPE_i^S(t)$$

365 , where $Q_i^F(t+1)$ and $Q_i^S(t+1)$ are the values for fat and sugar components updated from the
 366 previous fat and sugar value estimates, $Q_i^F(t)$ and $Q_i^S(t)$, and the NPEs for fat and sugar multiplied
 367 by the learning rate $\alpha \in [0,1]$.

368 Notably, although no fat-sugar interaction term was included when fat values and sugar
 369 values were integrated into aggregated reward values, the multiplication of fat and sugar values
 370 inherently introduced supra-additivity when fat and sugar were both present. Specifically, a fixed
 371 increase in fat value from higher fat content would give rise to a greater influence on the aggregated
 372 reward value when the sugar level is higher, and vice versa. Therefore, we used this specification to
 373 test whether the monkeys were sensitive to individual fat and sugar content and possible fat-sugar
 374 interactions, similar to our nutrient-sensitive logistic regression models and the RL models described
 375 above.
 376

Model comparison

We performed model comparisons based on the Akaike Information Criterion (AIC), which evaluated model prediction while penalizing excessive free parameters and computed the relative model likelihood of model i , P_i (Burnham et al., 2011) as below,

$$P_i = \exp\left(\frac{AIC_{min} - AIC_i}{2}\right)$$

The session-averaged AIC values were compared with our best-fitting model (*NutVal-Forget* model) and the AIC differences indicated the results of the model comparison. **Fig. 4F** shows the overview of model comparisons across all RL models, with the gray box indicating the statistical threshold compared to the best-fitting *NutVal-Forget* model. The *NutVal-Forget* model was directly compared with the *Energy* model and *ObjVal* model (**Fig. 4G** inset), and with the *NutRPE* model (**Fig. 7C** inset) based on AIC differences across sessions.

Additionally, we report the number of best-fitting sessions for each model. In each session, we first selected best-fitting models from the nine RL models and calculated the number of best-fitting sessions for each model. We allowed multiple best-fitting models in each session when they were statistically indistinguishable (relative model likelihood > 0.05).

Model validation

Model-independent results

In **Fig. 4E**, we simulated choices using the best-fitting *NutVal-Forget* model and computed the conditional choice probabilities depending on the reward and choice history of the previous trial, as we did for actual choice data in **Fig. 2E**.

Model recovery analysis

To evaluate how well our model comparison could distinguish between our candidate models, we performed a model recovery analysis (**Fig. 5**) to compare the best-fitting models with the actual models used to generate simulated choices in our task (Wilson and Collins, 2019). Specifically, for the nine RL models, we simulated choices in the actual trial conditions (N = 23 sessions from two monkeys) and performed our AIC-based model comparison to select the best-fitting model for each simulated session (23 sessions × 9 models = 207 simulated sessions). We used the actual trial conditions and the fitted parameters in the 23 sessions from both monkeys to approximate conditions in our model comparison. Next, we computed the conditional probabilities of identifying the correct models given the simulated model (confusion matrix, **Fig 5A**) or the fit model (inversion matrix, **Fig. 5B**). The confusion matrix quantifies how sensitively this approach could recover the underlying model, P(fit model | simulated model); the inversion matrix indicates how reliable a best-fitting model truly identifies the underlying model, P(simulated model | fit model), compared to misclassification of other candidate models.

Parameter recovery analysis

To evaluate the reliability of the parameter estimates from the model fitting, we examined how well the fitted parameters of each model were correlated with the true simulated parameters, using the same simulation data as in the model recovery analysis (**Fig. 6A-H**). We also reported the cross-correlation between the eight free parameters in the *NutVal-Forget* model because correlations between parameters might influence the reliability of parameter estimates in the model fitting (**Fig. 6I**).

Data and materials availability: All data are available upon reasonable request to the Corresponding Author.

RESULTS

Experimental design. Two monkeys performed in a dynamic foraging task to obtain probabilistic liquid rewards with defined nutrient compositions (**Fig. 1A**). In each trial, the monkey viewed two visual cues, randomly drawn from a set of four cues, that were first presented sequentially before reappearing in a randomized left-right arrangement as choice options. The monkey was then required to choose between the two cues by touching a cue-associated target on the touch monitor. Following the touch choice, the animal received either a large amount of liquid reward ('rewarded trials') or a small amount ('non-rewarded trials') depending on a prespecified reward probability (P). The visual cues were each associated with one of four different liquid rewards. Cue-reward associations were fixed within each session and we used new, untrained visual cues in each session to avoid influences of prior experience. We designed four liquid rewards that differed only in fat and sugar levels (**Fig. 1B**; LFLS: low-fat low-sugar; HFLS: high-fat low-sugar; LFHS: low-fat high-sugar; HFHS: high-fat high-sugar) while controlling the flavor (blackcurrant or peach) and other ingredients (protein, salt, etc.; **Table 1**). The LFLS liquid was lowest in energy content; the HFLS and LFHS liquids had matched, intermediate energy content; the HFHS liquid was highest in energy. This reward set was the same as the one used in our previous study in an economic choice task (Huang et al., 2021).

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 learning and value-updating from reward outcomes. Notably, this design maintained constant availability of fat and sugar irrespective of block type as there were always two high-probability and two low-probability options for both high-fat and high-sugar rewards. The design also held the available energy content constant across blocks. The liquids were matched in flavor (blackcurrant or peach) and other ingredients (protein, salt, etc.; **Table 1**); thus, differential learning and choice patterns could only be attributed to the nutrient content of the rewards.

In the following sections, we first identify each monkey's typical learning and choice patterns in representative sessions and aggregated data across sessions. We then model these data using logistic regression based on reward history and compare the performance of different RL models in accounting for the monkeys' choices.

Nutrients bias learning and choice for probabilistic rewards: example sessions. The behavior in two example sessions (**Fig. 1D**) showed that both monkeys exhibited preferences for specific nutrients while tracking the changing reward probabilities. Monkey Ya's choices (**Fig. 1D**, left) 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 (LFHS, HFHS) frequently, even when they were associated with a lower probability of obtaining a large reward amount compared to the high-fat reward (e.g., the blue and green curves in **Fig. 1D** compare the energy-matched LFHS and HFLS rewards). In addition, choice frequencies tracked the changing reward probabilities, particularly for the high-sugar rewards (compare red and blue curves across trial blocks). In addition, monkey Ya's choices tracked the changing reward probabilities, particularly for the high-sugar rewards (**Fig. 1D**, left, compare red and blue curves across trial blocks). By contrast, monkey Ym's choices (**Fig. 1D**, right) reflected both a preference for nutrient content and a strong dependence on reward probability. Although monkey Ym also preferred high-nutrient rewards over low-nutrient rewards with matched reward probabilities (**Fig. 1D**, right, compare red and yellow curves in 200 to 300 trials, blue and green curves in 100 to 200 trials), he chose the low-nutrient reward more frequently when it was associated with a high reward probability (yellow curve, e.g., 100 to 200 trials), indicating a noticeable but weaker preference for fat and sugar than monkey Ya (compare blue and green curves, e.g., 300 to 400 trials).

Two main results emerged from these single-session data. First, the monkeys' choices were sensitive to the rewards' nutrient content: both monkeys preferred sugar over fat content, although this preference was more pronounced in monkey Ya. Second, although both monkeys tracked changing reward probabilities, they differed in the way in which they integrated reward probabilities with nutrient preferences to make choices, with monkey Ym showing a more balanced integration of

reward probability and nutrient content. These results could not be explained by canonical RL models that would only consider binary reward outcomes without also considering how different reward components that are common to different rewards (i.e., fat and sugar) may affect learning and choice.

Nutrients bias learning and choice for probabilistic rewards: aggregated data. The choice patterns observed in single sessions were also found in aggregated data across sessions. Both monkeys successfully tracked changing reward probabilities by choosing an option more frequently when its reward probability switched from low ($P = 0.2$) to high ($P = 0.8$), as evident by averaged choice probabilities around reward-probability reversal points (**Fig. 2A**). Notably, higher reward probability increased monkey Ya's choices only for the high-sugar stimuli but not for the low-sugar stimuli, which suggested that learning depended on the reward's nutrient composition. By contrast, in monkey Ym, an increase in reward probability led to an increase in choices for all stimuli. We did not find strong evidence that the monkeys adjusted their choices more quickly to changed reward probabilities for the high-sugar and high-fat rewards, indicated by learning latencies (number of trials after probability reversals before significant changes in choice behavior; **Fig. 2B**). Overall, the monkeys responded differently to probability changes for rewards that differed in fat and sugar content.

When choices stabilized after probability reversals, monkey Ya showed a strong preference for high-sugar rewards, in particular to the combination of high-sugar and high-fat content, whereas monkey Ym showed graded preferences for high-sugar over high-fat rewards, followed by the low-nutrient option (**Fig. 2A**, 60-80 trials post-reversals). These preferences were evident in averaged choice probabilities for the four reward types across truncated sessions with balanced trial types (**Fig. 2C**, see *Methods*). The choice patterns reflecting nutrient preferences and adaptations to current reward probabilities typically emerged quickly within the first 30 trials of each session when novel visual cues were introduced (**Fig. 2D**). Importantly, both monkeys preferred the high-sugar (LFHS) reward over the energy-matched high-fat (HFLS) reward. This result is important because it shows that the monkeys' preferred specific nutrients, rather than being indifferent between nutrients of the same energy content, as would be suggested by ecological energy-maximization models. The monkeys' nutrient preferences in the present learning task were broadly similar to the nutrient preferences of the same monkeys in a previously investigated choice task (Huang et al., 2021), with the exception that monkey Ya showed less preference for fat reward in the current task.

To assess whether nutrient content influenced learning from recent reward experiences, we compared the conditional choice probabilities of the four reward types based on whether they were chosen or rewarded in the preceding trial (**Fig. 2E**). Both monkeys demonstrated a baseline preference for high-sugar rewards: the monkeys were approximately two times (Ym) or three times (Ya) more likely to choose high-sugar rewards (HFHS and LFHS) than low-sugar rewards (HFLS and LFLS) (**Fig 2E**, left). Importantly, both monkeys were more likely to repeat choices for recently received large rewards compared to small rewards, except for low-sugar rewards in monkey Ya (**Fig. 2E**, conditional choice probability after small rewards ('Non-rewarded', middle) and large rewards ('Rewarded', right)). Notably, because we introduced novel visual cues for each session, these results were not due to preferences for specific visual stimuli. Similarly, these results cannot be simply explained by the preference for calories because energy-matched rewards (i.e., LFHS and HFLS) had distinct effects on subsequent choices depending on their fat and sugar content (**Fig. 2E**, green and blue curves). Thus, the monkeys learned differently from different recent rewards, according to their nutrient preferences.

These aggregated data showed that the monkeys had subjective preferences for specific high-nutrient rewards and tracked changing reward probabilities depending on the rewards' nutrient content. Because visual cues changed every session, the results were not explained by preferences for specific visual stimuli. Thus, the nutrient composition of the food rewards and the animals' individual preferences for sugar and fat biased learning and choice. We next modeled these effects in trial-by-trial data using reward-history regression models and a formal RL approach.

Nutrient-specific reward history and choice history influence monkeys' choices. To formally characterize nutrient-dependent learning mechanisms, we followed approaches from previous studies that used logistic regression to link monkeys' choices to the recent history of rewards and choices (Corrado et al., 2005; Lau and Glimcher, 2005). In this approach, current-trial choices are

modeled by the sum of (i) subjectively weighted recent rewards resulting from choices for specific options ('reward history') and (ii) subjectively weighted recent choices for specific options ('choice history'); the subjective weights are determined by fitting the logistic-regression model to an individual animal's trial-by-trial data of current choices, past choices, and past rewards. To establish a baseline reference, we first modeled the monkeys' trial-by-trial choices with a logistic regression that accounted for whether options were chosen in previous offers (choice history) and whether previous choices were rewarded (reward history), irrespective of the rewards' nutrient composition (*History model*, see *Methods*). Because options were randomly offered in our experiments and not every reward appeared on each trial, we indexed the recent rewards and choices only when the specific option was offered (see *Methods*). In this *History model* (**Fig. 3A, B**), the most recent rewards and choices were associated with the highest regression coefficients, indicating a strong effect of recent reward and recent choice on current choice, whereas regression coefficients for more remote rewards and choices declined as a function of past trials, as observed in previous studies (Corrado et al., 2005; Lau and Glimcher, 2005).

We next modeled the contribution of individual nutrient components (*Nutrient model*, see *Methods*; **Fig. 3C, D**) by including nutrient-history interaction regressors that decomposed aggregated influences of reward and choice history (V_t) into contributions from the low-nutrient baseline liquid (B_t , yellow), high-fat content (F_t , green), and high-sugar content (S_t , blue). The *Nutrient model* outperformed alternative models that encoded only binary outcomes of reward and choice history but ignored reward nutrient compositions (*History model*, see *Methods*), with model comparisons being robust for different history lengths (**Fig. 3E**). Coefficients from this *Nutrient model* showed that although more recent rewards and choices generally had stronger effects on current choices, these effects depended on the nutrient composition of recent rewards (**Fig. 3F, G**). Specifically, in monkey Ya, reward-history coefficients for recent low-nutrient rewards were small whereas reward-history coefficients for the most recent high-nutrient rewards, especially high-sugar rewards, were large and highly significant (**Fig. 3F**, left). For monkey Ym, reward-history coefficients for both the baseline and high-nutrient rewards were large and significant and declined for more remote rewards, indicating both nutrient-specific and nutrient-independent effects of reward history (**Fig. 3F**, right). In both monkeys, nutrient-specific choice history effects for fat and sugar were large and significant for recent choices and declined for more remote choices (**Fig. 3G**). Critically, because the nutrient-specific history coefficients were modeled alongside the baseline history coefficients, they explained distinct parts of the variation in the monkeys' choices. In other words, significant coefficients for fat- and sugar-specific reward and choice history could not be explained by general, nutrient-independent history effects, which were captured by the baseline regressors. Notably, nutrient-specific reward history captured effects on choices that derived from the difference between rewarded and unrewarded choices for particular nutrients. By contrast, choice history captured effects related to recent choices for particular nutrients irrespective of reward outcome (Lau and Glimcher, 2005); choice history effects, therefore, reflected the monkeys' preference for particular rewards independent of reward probability. (Accordingly, in a supplemental analysis without choice-history regressors, nutrient-specific reward-history coefficients increased to reflect monkeys' nutrient preferences as well as the effect of reward probability.)

In summary, logistic-regression analysis showed that the monkeys' choices were influenced by the history of recently obtained fat and sugar rewards, and by the history of recently made choices for fat and sugar rewards. These nutrient-specific effects of reward history and choice history were not accounted for by the separately modeled general (i.e., nutrient-independent) reward- and choice-history effects.

Reinforcement learning based on nutrient-specific values. Having established that the monkeys' choices depended on the reward history for specific nutrients, we next modeled the effect of nutrients on learning and choice within the RL framework. Learning from rewards is formalized by RL models that update trial-by-trial expected values for each option based on reward outcomes. However, our results show that the monkeys were sensitive not only to whether they received a reward but also to the reward's nutrient composition. Accordingly, RL models may require nutrient-specific parameters to account for choices in dynamic environments with varying nutrient-reward composition as tested here. We therefore developed a novel, nutrient-sensitive RL model (*NutVal-Forget model*, see *Methods*) that incorporated a subjective 'nutrient-value function' to capture how specific nutrients

(fat, sugar) may differentially influence the trial-by-trial updating of expected reward values and their influence on choice (**Fig 4A**). Instead of updating the value of the chosen option with a binary reward outcome, the *NutVal-Forget* model updated values based on animal-specific subjective values for fat, sugar, and fat-sugar interaction, as follows. (We note that this model assumes that the animals are sensitive to fat, sugar, and their interaction, and combine these nutrient valuations into integrated values for decision-making.)

$$Q_i(t+1) = Q_i(t) + \alpha \cdot [V_i(t) \times R(t) - Q_i(t)]$$

$$V_i(t) = \frac{V_F^{I_F(t)} \cdot V_S^{I_S(t)} \cdot V_{FS}^{I_F(t) \times I_S(t)}}{V_F \cdot V_S \cdot V_{FS}}, \quad R(t) = \begin{cases} 1 & , \text{if rewarded} \\ d & , \text{if not rewarded} \end{cases} \quad d \in [0,1]$$

$$I_F(t) = \begin{cases} 1, & i(t) = HFLS, HFHS \\ 0, & i(t) = LFLS, LFHS \end{cases}, \quad I_S(t) = \begin{cases} 1, & i(t) = LFHS, HFHS \\ 0, & i(t) = LFLS, HFLS \end{cases}$$

In this model, the expected value for reward i , Q_i , was updated depending on both the reward outcome on trial t , $R(t)$, (rewarded or not rewarded) and the chosen reward type on trial t , $i(t)$. The animal-specific subjective value for reward i , $V_i(t)$, was constructed by the subjective valuation of the reward's fat level $I_F(t)$ and sugar level $I_S(t)$ based on fixed subjective values for fat (V_F), sugar (V_S), and their interaction (V_{FS}). Thus, we use Q_i to refer to the trial-by-trial updated values in the RL model, following the term 'Q-value' in the RL literature. By contrast, we use V_i to refer to the animal-specific subjective values for nutrient components, which were fixed in each session. Additionally, we included a discount factor $d \in [0,1]$ as a free parameter to capture the reduced but nutrient-dependent effects of the small rewards compared to the large rewards (**Fig. 2E**). Thus, the free parameters in the nutrient-sensitive RL model included the three nutrient-value parameters—related to fat, sugar, and their interaction—the discount factor for the small reward, as well as other standard RL parameters (see *Methods*). Because subjective values were measured compared to the low-nutrient baseline reward, nutrient-value parameters (V_F , V_S) larger than 1 suggested a preference for the respective nutrient; a V_{FS} parameter larger than 1 indicated supra-additive effects of fat and sugar on subjective value. Without loss of generality, we normalized all reward values to the value of the high-fat high-sugar reward ($V_F \cdot V_S \cdot V_{FS}$), to constrain all value parameters between 0 and 1. For options that were not chosen ('unchosen') and not offered ('unoffered') on a given trial, we allowed their expected values to decay as follows,

$$Q_j(t+1) = Q_j(t) \cdot (1 - \delta), \quad \forall j \neq i(t)$$

, where the values of the unchosen and unoffered rewards, $Q_j(t)$, were discounted according to a forgetting rate (δ), which would be 0 for perfect value memory. We referred to this RL model as 'nutrient-sensitive' because it differentially learned from reward outcomes based on their nutrient content. Specifically, different from canonical RL models, our model constructed reward values for choice options based on subjective values for each reward's nutrient components. Although canonical RL models can be applied to situations with multiple outcomes, they would not estimate values based on the outcomes' nutrient (or other intrinsic) components that are common to different outcomes.

We fitted the *NutVal-Forget* model to the trial-by-trial choice data and reward outcomes in each session, separately for both monkeys (see *Methods*). The resulting best-fitting parameter values for the nutrient components indicated that across sessions, both monkeys assigned higher values to the high-sugar rewards and that monkey Ym assigned a higher value to fat but monkey Ya did not (**Fig. 4B**). Fat and sugar showed a small but significant supra-additive effect on choices in monkey Ya and a small, non-significant negative interaction in monkey Ym. In an extended model (*NutAlphaVal-Forget model*, see *Methods*), we examined whether nutrients had additional influences on the learning rates in RL models, separate from their effect on reward values, but did not find evidence for nutrient-specific learning rates in either monkey (**Fig. 4C**). (Further, learning latencies resulting from simulated choices for the actual trial sequences that the animals experienced using

the *NutVal-Forget* model showed no significant differences between the four rewards; $P > 0.08$ for all pairwise comparisons.) The nutrient values in the *NutVal-Forget* model were largely stable across testing sessions, despite some session-specific variation (**Fig. 4D**). Additionally, both monkeys' showed evidence for 'forgetting' the values of unoffered and unchosen options, according to the small but significant value-forgetting rates (**Fig. 4C**). In both monkeys, our main *NutVal-Forget* model reproduced the behavioral signatures of nutrient-specific learning from **Fig. 2E** (**Fig. 4E**). In model comparisons across sessions, the *NutVal-Forget* model outperformed alternative RL models with variations on value updating (*Basic*, *Asym*, *Forget*) and nutrient-specific parameters (*Binary*, *NutVal*, *NutValAlpha*) (**Fig. 4F**, see *Methods*). Furthermore, psychometric curves illustrated that subjective values derived from the *NutVal-Forget* model accounted for the monkeys' choices, confirmed by model-fit indicators across sessions (pseudo- R^2 : monkey Ya = 0.80 ± 0.02 ; monkey Ym = 0.48 ± 0.02 ; percentage of correctly modeled choices: monkey Ya = $94.7 \pm 1.3\%$; monkey Ym = $83.4 \pm 0.7\%$) (**Fig. 4G**, black).

To rule out the possibility that the differential learning from rewards was driven by the energy content of nutrients rather than by specific nutrients, as assumed in ecological foraging models, we compared the performance of our main *NutVal-Forget* model (**Fig. 4G**, black), which estimated separate effects of sugar and fat components on subjective values, to an RL model based on the energy content of the rewards (*Energy model*, **Fig. 4G**, red). The key difference of the *Energy model* was that it learned equally from the isocaloric HFLS and LFHS rewards, even though they differed in fat and sugar levels (**Fig. 1B**). In both monkeys, the *NutVal-Forget* model (**Fig. 4G**, black and inset) provided a significantly better fit to the monkey's choices than the *Energy model* (**Fig. 4G**, red data and inset), consistent with the monkeys' preference for high-sugar over high-fat stimuli of matched energy content (**Fig. 2C**).

For completeness, we also constructed an object-specific RL model with separate values for the four different rewards as free parameters that were fixed in each session (*ObjVal* model). Different from this object-specific model, which specified separate values for all available rewards, our main *NutVal-Forget* model specified values for the two nutrient components that were common to the different rewards. Our main *NutVal-Forget* model performed similarly to the *ObjVal* model, despite having one fewer free parameter (**Fig. 4G**, yellow). Although both models performed similarly in this data set, we note that the nutrient-based model is more flexible compared to the stimulus-specific model because it uses a value function for particular reward *components* (fat, sugar) rather than for each individual reward. Thus, the nutrient-based model could use the same value function for increasing numbers of different rewards (provided they were based on common nutrient components), whereas the object-specific model would require one additional free parameter for each reward type.

We tested the sensitivity and reliability of our model comparison in a model recovery analysis (**Fig. 5**, see *Methods*). We found that the model comparison correctly identified 82% of the simulated sessions from the *NutVal-Forget* model (**Fig. 5A**) and that the *NutVal-Forget* model was the most probable generative model given our model comparison result (**Fig. 5B**). Notably, the winning *NutVal-Forget* model also generalized well to other candidate models (**Fig. 5A**), likely due to the negligible influences of non-relevant, additional parameters, e.g. nutrient-specific learning rates (**Fig. 4C**). The *NutVal-Forget* model was also associated with relatively lower absolute confidence (red square, **Fig. 5B**), likely because it was compared with a large number of nested (i.e., overall similar) models. We also tested the identifiability of the parameters in a parameter recovery analysis, and examined the trade-offs between parameters (**Fig. 6**). Importantly, the reliability of the parameter estimates in **Fig. 4B** was supported by the successful recovery of model parameters from simulated sessions (**Fig. 6A-H**) despite some correlation between V_F and V_S (**Fig. 6I**). Further, the parameter-correlation matrix (**Fig. 6I**) indicated trade-off between specific parameters: (i) We found a positive correlation between V_F and V_S , suggesting that preferences for fat and sugar covaried across sessions. (ii) We found that the right-side bias (β_0) was positively correlated with the forgetting rate (δ) and inversely correlated with V_S and V_{FS} , indicating that nutrient-influences on choices were weaker in periods when monkeys were apparently less focused on the task (displaying side-bias and higher forgetting rate). (iii) We found that the discount rate (d) between large and small rewards (which was fixed across reward types) was inversely correlated with nutrient values; this was expected in our models because when any of the nutrient value increases, values for other reward

types decrease naturally due to value normalization (see *Methods*) and the fixed discount rate would decrease to adapt to such inadvertent influences from value normalization.

Thus, the monkeys' choices for probabilistic rewards with defined nutrient compositions were well explained by a nutrient-sensitive RL model that assigned subjective values of specific fat and sugar components to reward outcomes and choice options.

Value updating based on nutrient-specific reward prediction errors. Our nutrient-sensitive RL model suggested that the monkeys tracked fat and sugar value components of probabilistic reward outcomes and integrated them into a scalar value that guided choices. To better understand the dynamics of nutrient-specific value updating, we examined an extended nutrient prediction error-based RL model (*Nut-RPE*). In this model, the reward value on trial t , $Q_i(t)$, was decomposed into fat and sugar value components (**Fig. 7A**, see *Methods*). These value components (i) separately adapted to changes in reward probabilities based on fat and sugar reward prediction errors (RPEs) from experienced nutrient outcomes, and (ii) were flexibly combined into the integrated reward value for each option, based on the animal's nutrient preference (**Fig. 7A**). On each trial, the *Nut-RPE* model compared the expected fat and sugar value components with the obtained nutrient outcomes and computed nutrient-specific RPEs to update reward expectations. Thus, the model reveals the trial-by-trial dynamics of latent learning variables in dynamic nutrient-reward choices, including nutrient-specific values, integrated reward values, and nutrient-specific RPEs. We examined this model in part because it could serve as a tool in future studies to identify neuronal signals related to nutrient-specific learning and decision variables.

By capturing the monkeys' dynamic updating and subjective integration of fat-sugar value components, the model closely tracked evolving choice patterns for nutritionally-distinct rewards in single sessions (**Fig. 7B**). The *Nut-RPE* model provided a good fit to the monkeys' choices across sessions according to psychometric curves and model-fit indicators (**Fig. 7C**: pseudo- R^2 : monkey Ya = 0.82 ± 0.01 ; monkey Ym = 0.46 ± 0.02 ; percentage of correctly modeled choices: monkey Ya = $95.8 \pm 0.8\%$; monkey Ym = $82.3 \pm 1.0\%$). The trajectories of fat value and sugar value within single sessions characterized each monkey's idiosyncratic sensitivity to fat and sugar when updating expected reward values (**Fig. 7D**). Specifically, reward values in monkey Ya were primarily driven by the sugar value component, which was updated after the monkey received high-sugar rewards and tracked their blockwise changes in reward probabilities, whereas the fat-value component varied much less within a session (**Fig. 7D**, left). In monkey Ym, both fat and sugar value components contributed to value updating and tracked fluctuating reward probabilities throughout the session (**Fig. 7D**, right).

The *Nut-RPE* model also allowed us to examine the dynamics of nutrient-specific RPEs that separately updated the fat and sugar value components. For instance, in the first block of the single session data in **Fig. 7E**, when the LFHS reward (blue data in **Fig. 7B** and **Fig. 7E**, trials 1-100) had a high reward probability, LFHS was frequently chosen and produced mostly positive sugar RPEs. The few negative sugar RPEs were due to the occasional low reward outcomes and fluctuations of values during learning. The magnitude of sugar RPEs decreased during this block while the monkey learned the reward value more accurately. After probability reversal, LFHS was chosen less frequently and mostly produced low-magnitude rewards. However, it produced a few large positive sugar RPEs during that period (**Fig. 7E**, blue data, trials 100-200), which illustrated that even low rewards can produce large positive sugar RPE if the reward contained sugar, because reward outcomes were scaled by the nutrient-specific subjective value in the *Nut-RPE* model. By contrast, when the high-fat low-sugar stimulus was associated with high reward probability in the third block (HFLS, green data in **Fig. 7B**, **E**), it consistently produced large rewards but only very small sugar RPEs because of the rewards' low sugar content.

We note that both the *NutRPE* model and the *NutVal-Forget* model have eight free parameters. Although the *NutRPE* model does not have a free parameter for fat-sugar interaction (since it assumes multiplication of fat values and sugar values), it does have two separate learning rates for fat-value and sugar-value updating. A direct comparison between the *NutRPE* model and the *NutVal-Forget* model (**Fig. 4F**) showed that these two models performed equally well in monkey Ya and the *NutVal-Forget* model outperformed the *NutRPE* model in monkey Ym. Therefore, the *NutRPE* model revealed trial-by-trial nutrient value dynamics without compromising model performance.

Thus, subjective nutrient-value functions guided the dynamic updating and integration of reward values based on nutrient-specific reward components and related nutrient-specific RPEs. The nutrient RPEs were sensitive to both the presence of a specific nutrient in the experienced reward outcome and the reward size. These results provide a basis for hypotheses about candidate neurons encoding nutrient-specific values and RPEs (**Fig. 8**), as explained in the Discussion.

DISCUSSION

We investigated monkeys' choices for nutrient-defined rewards under varying reward probabilities and found that the rewards' nutrient composition strongly influenced choices and learning. As in our previous study involving the same animals (Huang et al., 2021), the monkeys' choices reflected individual fat and sugar preferences, consistent with the assignment of subjective values to nutrients. Importantly, in the present study we show that these nutrient preferences affected how the animals adapt their choices to changing reward probabilities during learning. Specifically, the monkeys were more likely to repeat choices that led to rewards containing preferred nutrients and frequently chose these options even under low reward probabilities. As in previous studies (Lau and Glimcher, 2005; Tsutsui et al., 2016), more recent rewards had a stronger influence on current choice. Critically, we found that this impact of reward history depended on the reward's nutrient composition: past rewards that were high in preferred sugar content influenced subsequent choices more strongly than less preferred low-nutrient or high-fat rewards. Choice history irrespective of reward outcomes also had a significant nutrient-dependent effect on choice, with stronger effects of past choices for preferred nutrients. We developed a nutrient-sensitive RL model that incorporated these influences of preferred nutrients on learning and choice. The model updated the values of sugar and fat components of expected rewards trial by trial, based on recently experienced rewards, and integrated these nutrient-specific values into scalar reward values that explained the monkeys' choices. Value updating in response to particular nutrient outcomes was governed by nutrient-specific RPEs. Thus, different from canonical RL models that learn from binary reward outcomes (Sutton and Barto, 1998), our nutrient-sensitive RL model learned based on the rewards' nutrient components, which influenced the subjective values of choice options and reward outcomes. Nutrient effects on learning and choice were not explained by energy content or other reward properties (flavor, salt, protein), which we controlled, or preferences for visual cues, which we changed every session. Our results suggest that nutrients constitute important reward components that influence subjective valuation, learning and choice, and that RL models can be usefully extended to incorporate nutrient-specific effects on behavior.

Previous studies in macaques revealed important factors influencing learning and choice, including reward and choice history (Corrado et al., 2005; Lau and Glimcher, 2005; Samejima et al., 2005; Kennerley et al., 2006; Lee et al., 2012; Seo et al., 2012; Tsutsui et al., 2016), the variance of recent rewards (Grabenhorst et al., 2019a), novelty and rarity of choice objects (Costa et al., 2019; Rothenhoefer et al., 2021), and social observations (Grabenhorst et al., 2019b). Because these studies did not vary the composition of reward outcomes, they could not test how specific nutrients affect learning and choice. Other studies demonstrated that macaques have sophisticated preferences for different reward types that comply with principles of economic choice theory (Padoa-Schioppa and Assad, 2006; Lak et al., 2014; Raghuraman and Padoa-Schioppa, 2014; Pastor-Bernier et al., 2017) but did not examine how nutrient rewards affect RL. Here, we reasoned that nutrients represent biologically critical, intrinsic reward components that are essential for survival and thus should affect learning and choice in dynamic reward environments, to ensure stable nutrient intake. By systematically manipulating the sugar and fat content of liquid rewards, we showed that monkeys' have subjective preferences for specific nutrients and that these preferences influence how monkeys adapt their choices to changing reward probabilities. Similar to one previous study (Raghuraman and Padoa-Schioppa, 2014), our monkeys based their choices on both subjective reward preferences and reward probabilities. However, different from that study, our monkeys were required to learn reward probabilities from experience rather than from explicit pre-trained cues; we also systematically varied the rewards' nutrient composition, rather than varying reward type irrespective of nutrient content. Thus, beyond the general effects of reward type, we determined that specific nutrients differentially affect learning and choice. We do not suggest that nutrients influence learning independently of subjective value. Rather, our data and model suggest that nutrients

influence learning by affecting the subjective values that animals assigned to reward outcomes and choice options.

In our study, monkeys adapted their choices to changing reward probabilities to obtain large amounts of preferred nutrient-defined rewards; when probabilities changed, the monkeys could switch their choices to less-preferred but more available alternatives. The animals were not required to learn the nutrient composition of novel foods as the four rewards used were highly familiar to the animals. This situation resembles natural scenarios in which monkeys regularly forage for a range of familiar food rewards and must determine which foods are currently available (Cui et al., 2018; Cui et al., 2019; Cui et al., 2020). Wild macaques solve such foraging problems by adapting their choices to short-term changes in reward availability (e.g., switching between foraging grounds) and to longer-term (e.g., seasonal) changes. For example, preferred foods, such as fat- and protein-rich seeds may only be available in autumn and winter. Accordingly, wild macaques adapt their foraging patterns to pursue different foods in other seasons, such as leaves and herbs, while maintaining a stable nutrient balance (Cui et al., 2019). Our simple repeated-choice paradigm captured these features of macaques' foraging environments in the laboratory. Our findings suggest that the computations underlying nutrient-adaptive foraging choices involve the assignment of subjective values to specific nutrients, the integration and updating of nutrient values from recently experienced reward outcomes, and the comparisons of values to make choices.

Flavor-nutrient conditioning is an important process in feeding and foraging behavior (de Araujo et al., 2020; Dayan, 2022); however, it is unlikely to have played a major role in our study. The animals were highly familiar with the four different nutrient-defined rewards through task training over many months. Thus, flavor-nutrient conditioning effects would have stabilized by the time the present experiments were conducted.

We extended a canonical RL model to account for the observed nutrient sensitivity of the monkeys' choices. Different from RL models that learn from binary reward outcomes, our nutrient-sensitive RL model evaluated rewards using an animal-specific 'nutrient-value function'. This nutrient-value function detected the fat and sugar content of the reward in addition to reward size, and weighted the reward according to individual monkeys' preferences for fat, sugar, and their interaction. As a consequence, the model learned expected values of choice options based on the reward's nutrient composition and monkeys' nutrient preferences, rather than based on reward frequency irrespective of nutrient content as in standard RL models. An extension of this model decomposed subjective values into separate fat and sugar value components, and updated these value components with nutrient-specific prediction errors that considered the difference between experienced and expected values for fat and sugar components. Although not investigated here, our experimental paradigm and model could be extended to incorporate internal, physiological set-points as proposed in recent homeostatic RL (Keramati and Gutkin, 2014). Our model could also separately model short-term influences of sensory, hedonic food components on learning and choice, including recently demonstrated fat-related oral-texture effects (Huang et al., 2021), and longer-term effects related to post-ingestive processing (Dayan, 2022).

Our findings have implications for understanding the neural mechanisms underlying RL. We described a nutrient-specific learning process that updates value estimates for separate fat and sugar reward components and integrates this information into scalar values that guide adaptive choices. At the neural level, this mechanism would require neurons that encode individual nutrient values and dynamically update them via nutrient-specific RPEs (**Fig. 8**). By decomposing trial-by-trial reward values that guide RL into nutrient-value components, our model specifies computational signals that could help identify such neurons (**Fig. 7**). Midbrain dopamine neurons, orbitofrontal cortex, and amygdala participate in value-based decision-making, RL, and food evaluation (Padoa-Schioppa and Assad, 2006; Grabenhorst et al., 2010; Grabenhorst et al., 2012; Murray and Rudebeck, 2013; Lak et al., 2014; Stauffer et al., 2014; Suzuki et al., 2017; Murray and Rudebeck, 2018; Rolls et al., 2018; Grabenhorst et al., 2019b; Averbeck and Murray, 2020) and thus constitute targets for testing these hypotheses. For example, it will be interesting to determine whether primate dopamine neurons encode RPEs for specific nutrients. Recent studies demonstrated that dopamine neurons encode RPEs in terms of economic utility, thereby integrating influences of reward type, quantity, probability, and risk (Lak et al., 2014; Stauffer et al., 2014). It remains to be tested whether the presently observed nutrient-specific effects on learning and choice are incorporated into the dopamine neurons' RPE signal, or whether nutrient-specific reward components are processed in

other reward structures such as the orbitofrontal cortex and amygdala, where single neurons encode nutrients and sensory reward properties (Kadohisa et al., 2005; Huang and Grabenhorst, 2022).

In summary, our results identify nutrients as important reward components that influence monkeys' subjective valuations, choices, and learning. These processes were well described by a nutrient-sensitive RL model that updated the value of expected rewards based on their sugar and fat components using nutrient-specific RPEs. Our data and nutrient-sensitive RL model can serve as tools to guide future studies that aim to uncover nutrient-specific learning and decision computations and their neurophysiological implementations.

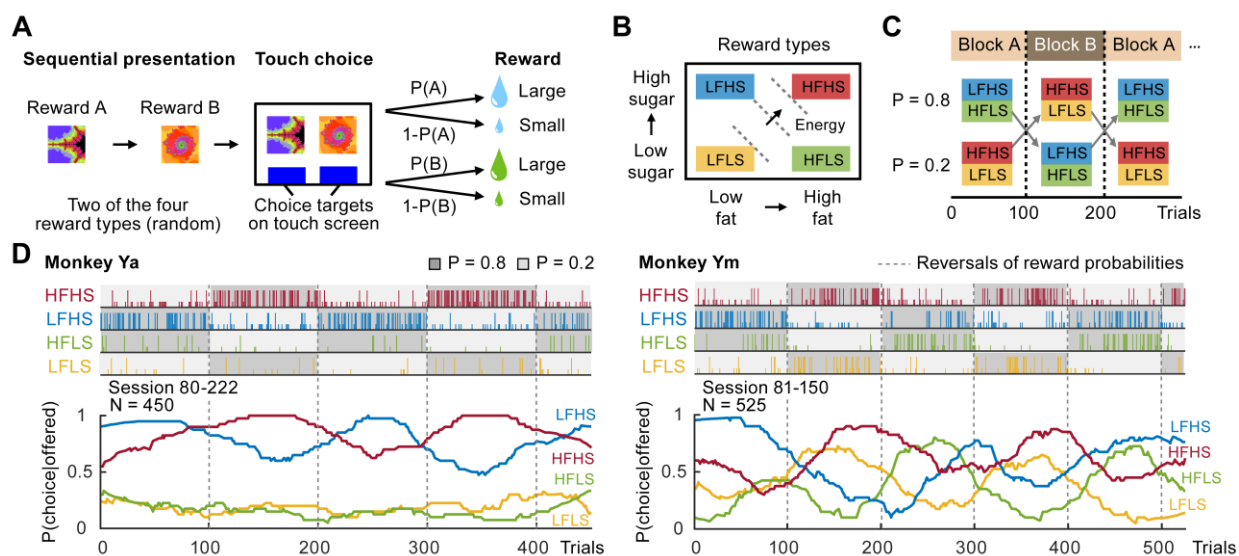


Fig. 1. Dynamic foraging task with nutrient-defined rewards. (A) Task structure. In each trial, the monkeys were first sequentially presented with two visual cues randomly drawn from a set of four cues and then made a left or right touch choice between these two cues when they were simultaneously presented. Following the touch choice, the animals received a large amount (0.5 mL) or a small amount (0.3 mL) of the associated liquid reward depending on a prespecified reward probability (P). (B) Reward design. Four types of liquids with 2×2 factorial fat and sugar levels were offered in each session: the low-fat low-sugar liquid (LFLS, yellow), the high-fat low-sugar liquid (HFLS, green), the low-fat high-sugar liquid (LFHS, blue), and the high-fat high-sugar liquid (HFHS, red). The LFHS and HFLS liquids were isocaloric and all rewards were matched in flavor (blackcurrant or peach) and other ingredients (e.g., protein, salt, etc., see **Table 1**). (C) Reward-probability schedule. The probabilities of receiving large rewards ('reward probability') were assigned in two block types. In block A, LFHS and HFLS were associated with a high probability ($P = 0.8$); LFLS and HFHS were associated with a low reward probability ($P = 0.2$). All reward probabilities were reversed in block B. Each session started with either block A or block B and the reward probabilities changed between the two block types every 100 trials with typically three to five alternations in each session. (D) Choices and reward outcomes in single sessions for monkey Ya (left) and monkey Ym (right). Tick marks represent choices for specific rewards: long marks indicate large-reward outcomes ('rewarded trials'); short marks indicate small-reward outcomes ('non-rewarded trials'). Reward types in dark-gray blocks were associated with high reward probability ($P = 0.8$) and light-gray blocks were associated with low reward probability ($P = 0.2$). Choice-probability curves show nine-trial running averages of choices for each reward (N: trials).

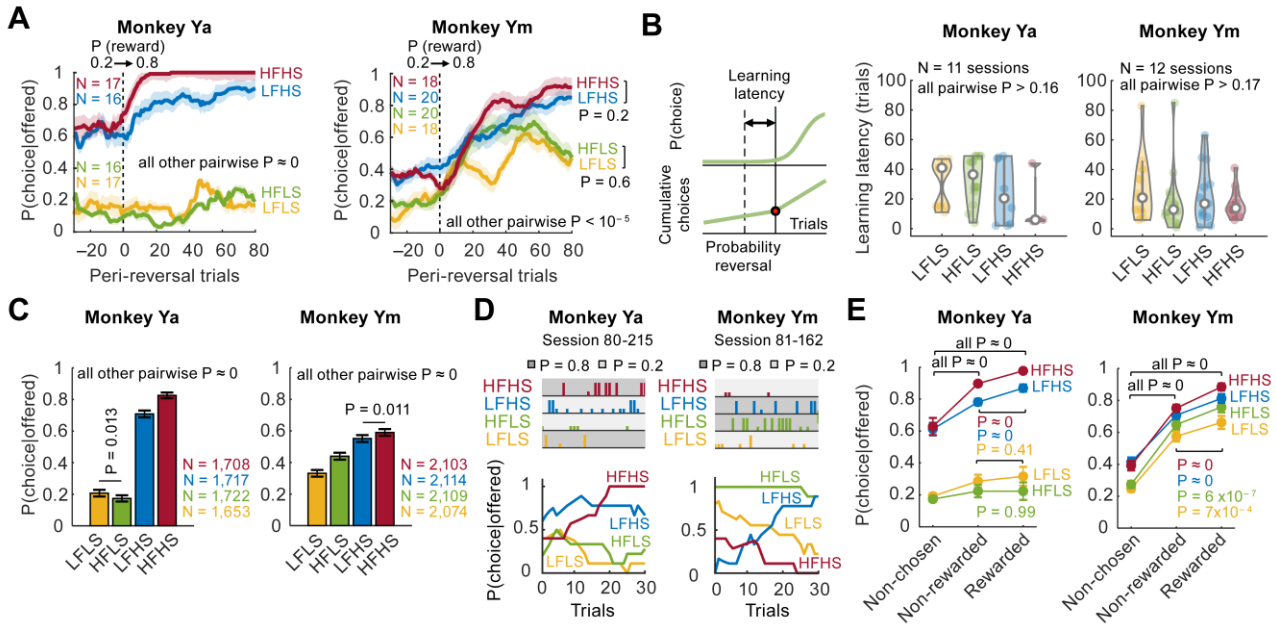


Fig. 2. Nutrient-specific learning and choice patterns across sessions. (A) Learning curves. Mean choice frequencies (nine-trial running averages \pm s.e.m.) aligned to reward probability reversals (dashed line, $P = 0.2 \rightarrow 0.8$) indicate how choices responded to changes in reward probabilities depending on reward nutrient content. N: number of low to high reward probability reversals. Two-sample t-test on the choice probability averaged across the 70th and 80th trials after reversals. **(B)** Learning latency was defined as trial intervals between probability reversals and the first significant change of choice patterns (see *Methods*). We included only the first probability reversal across sessions to avoid influences of different pre-reversal choice probabilities on learning latency. Median \pm 95 % confidence interval. P-values: pairwise two-sided Wilcoxon rank sum test. **(C)** Reward preference. Averaged choice frequencies (mean \pm s.e.m.) indicate preferences for the four reward types. The choice frequencies were aggregated across sessions that were truncated to have balanced block types (see *Methods*). P-values: pairwise two-sample binomial test. N: trial numbers. **(D)** Initial learning from novel visual cues. Choices and reward outcomes in the initial 30 trials of two example sessions for monkey Ya (left, block type A) and monkey Ym (right, block type B) show how the monkeys differentially associated novel visual cues to the reward types while learning from reward outcomes. Tick marks represent choices for specific rewards: long marks indicate large-reward outcomes ('rewarded trials'); short marks indicate small-reward outcomes ('non-rewarded trials'). Reward types in dark-gray blocks were associated with high reward probability ($P = 0.8$) and light-gray blocks were associated with low reward probability ($P = 0.2$). Choice-probability curves show nine-trial running averages of choices for each reward. **(E)** History-dependent choice probabilities. The probability of choosing each reward (mean \pm s.e.m.) increased after choosing and receiving the specific reward. Such influence of reward history depended on the reward's fat and sugar content. P-values: pairwise two-sample binomial test. The conditional probability of choosing each reward was computed based on the outcomes when the reward was last offered. Left: the reward was not chosen ('Non-chosen'). Trial numbers [LFLS; HFLS; LFHS, HFHS] = [1,775; 1,909; 622; 342] (Ya), [1,677; 1,507; 1,117; 1,042] (Ym); Middle: the reward was chosen but only a small reward was delivered ('Non-rewarded'). Trial numbers [LFLS; HFLS; LFHS, HFHS] = [472; 425; 1,707; 1,996] (Ya), [967; 1,171; 1,545; 1,616] (Ym); Right: the reward was chosen and a large reward was delivered ('Rewarded'). Trial numbers [LFLS; HFLS; LFHS, HFHS] = [231; 206; 971; 1,099] (Ya), [563; 704; 859; 953] (Ym).

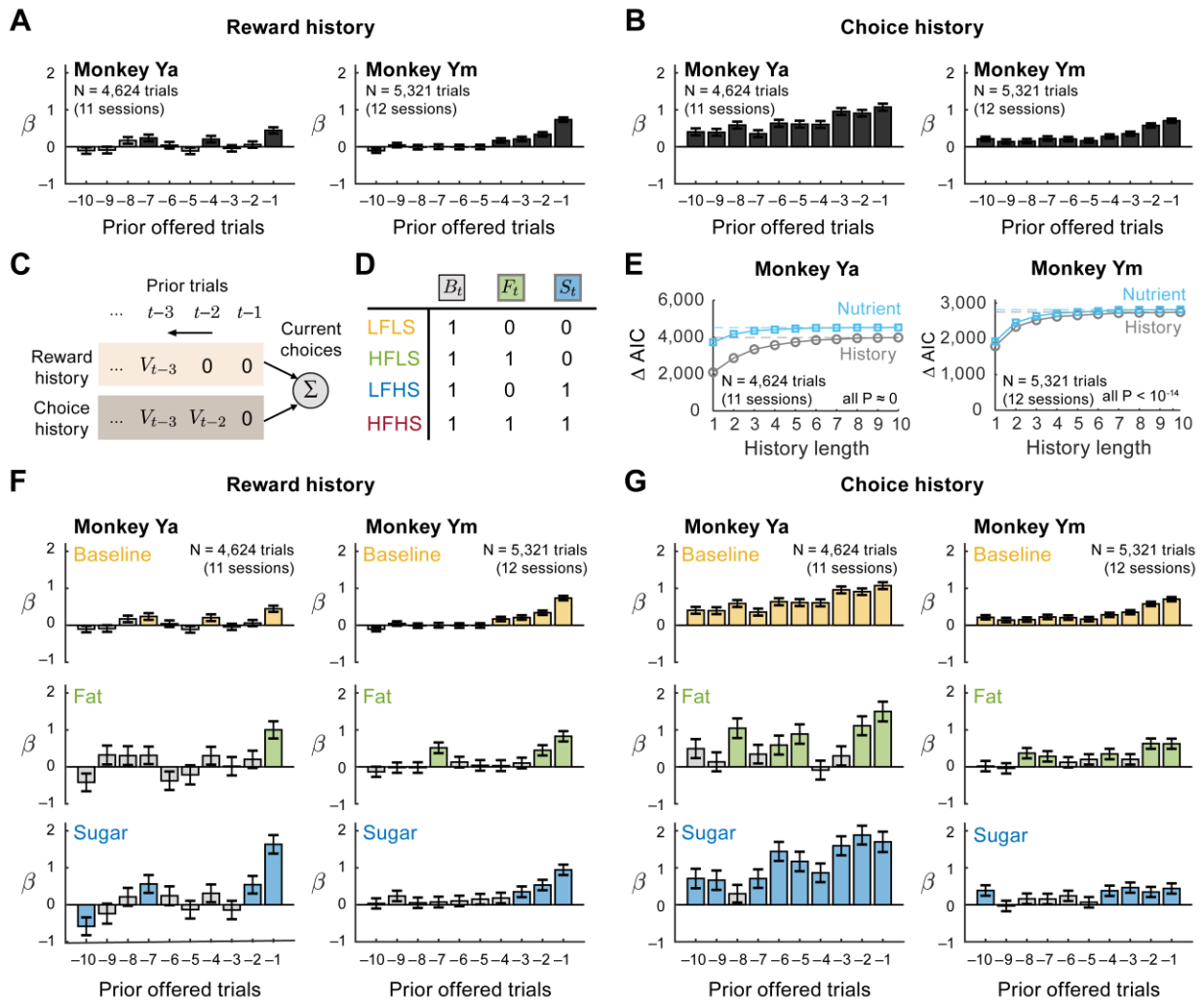


Fig. 3. Nutrient-dependent influences of reward and choice history on choice. (A-B) General reward- and choice-history effects in logistic regression. **(A)** Reward-history effects in logistic regression. Regression coefficients (\pm s.e.m.) for reward history, modeled across all stimuli, reveal baseline effects of recent large vs. small rewards on choice in the last ten trials preceding current-trial choice. Filled bars: $P < 0.05$; gray bars: non-significant. **(B)** Choice history effects in logistic regression. Regression coefficients for choice history, obtained from the same logistic regression model as in (A), reveal the baseline effect of recent choice on current-trial choice irrespective of reward outcome. **(C-G)** Nutrient-dependent logistic regression model. **(C)** Schematic of reward and choice history logistic regression model. Prior trials were indexed only when the same type of reward was offered. Current-trial choices were explained by recently rewarded and chosen rewards. **(D)** Nutrient-specific history regressors. In the nutrient-specific logistic regression model (*Nutrient model*), each reward outcome and choice on trial t was decomposed into effects of baseline low-nutrient ingredients (B_t , gray), additional fat content (F_t , green) and additional sugar content (S_t , blue). **(E)** Model comparison across history trial lengths. Performance of the *Nutrient model* (blue) and *History model* (gray) matched in history trial lengths up to the past 10 trials. Models were compared based on $\Delta AIC = AIC(\text{history trial length} = 0) - AIC(\text{history trial length} = i, i = 1, 2, \dots, 10)$ and confirmed by the loglikelihood test (P-value). AIC = Akaike Information Criterion. **(F-G)** Nutrient-specific reward- and choice-history effects. **(F)** Effects of nutrient-specific reward history. Regression coefficients for recent low-nutrient baseline rewards (yellow), fat-containing rewards (green), and sugar-containing rewards (blue) on current-trial choice. Nutrient-specific effects were estimated in the same model as low-nutrient baseline effects; thus, effects of fat and sugar reward history were not accounted for by general effects of reward history. **(G)** Effects of nutrient-specific choice history. Regression coefficients for recent low-nutrient baseline choice (yellow), the choice for fat-containing rewards (green), and the choice for sugar-containing rewards (blue) on current-trial choice. Nutrient-

960 specific effects were estimated in the same model as low-nutrient baseline effects; thus, effects of
961 fat and sugar choice history were not accounted for by general effects of choice history.

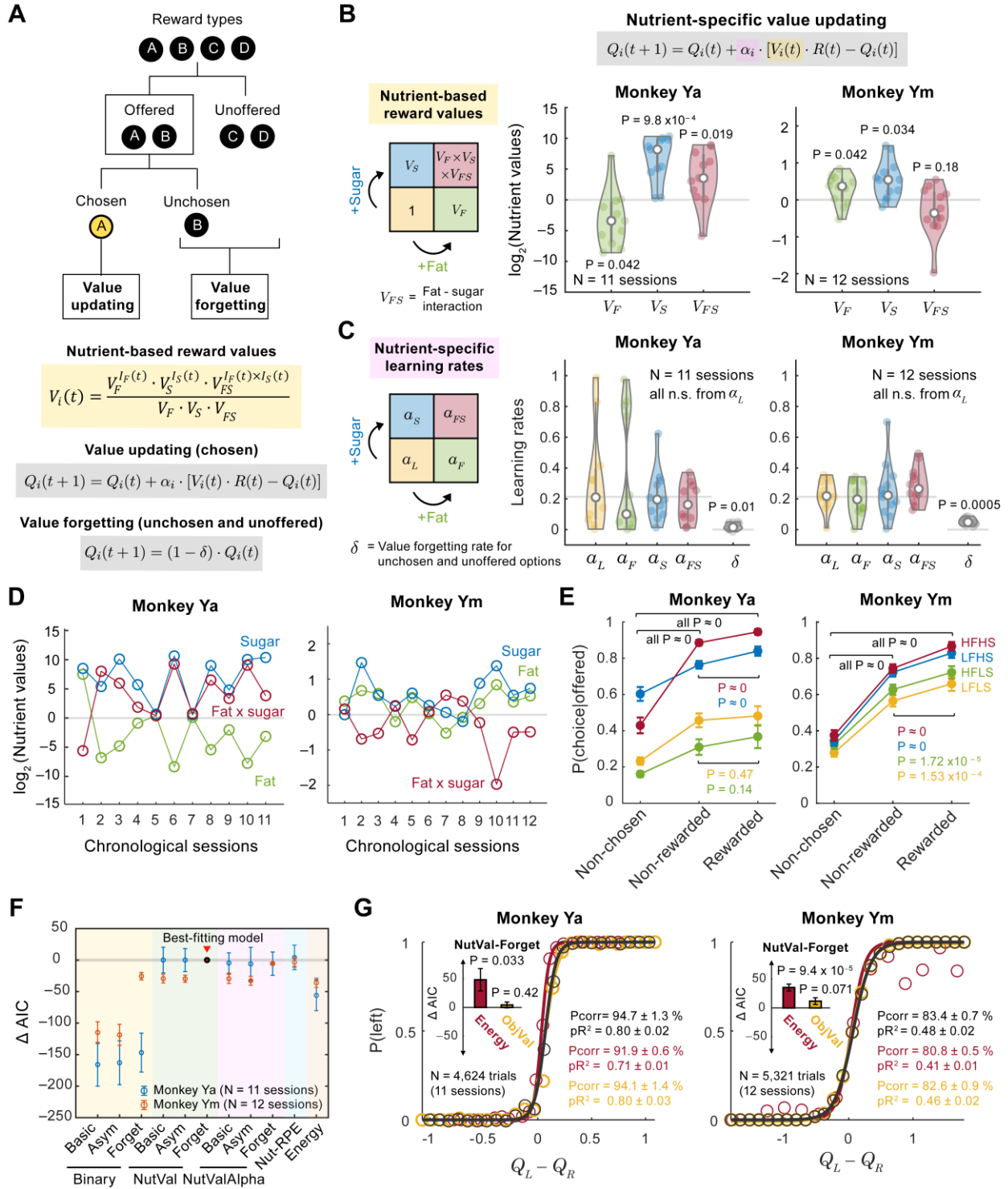


Fig. 4. Nutrient-sensitive reinforcement learning. (A) Nutrient-sensitive RL model (*NutVal-Forget* model, see *Methods*). Expected values for each option (Q_i) were iteratively updated based on subjective reward values, V_i , constructed by animal-specific values for high-fat content (V_F), high-sugar content (V_S), and fat-sugar interaction (V_{FS}), compared to the low-nutrient baseline reward ($V = 1$). $I_F(t)$ and $I_S(t)$ indicate the fat and sugar levels of the reward, respectively. Q_i refers to the value for reward i , which was updated depending both on the reward outcome on trial t , $R(t)$, and the reward type received on trial t , $i(t)$. All reward values were normalized to the highest reward value ($V_F \cdot V_S \cdot V_{FS}$) to be constrained between 0 and 1. Values for unrewarded and unoffered options decayed by a factor of $(1 - \delta)$. δ : forgetting rate $\in [0,1]$. (B) Subjective nutrient values. Fitted parameters indicating subjective values for fat, sugar, and fat-sugar interaction in each session fitted by the *NutVal-Forget* model. Data were log-transformed (base 2) for visualization. P-value: Wilcoxon

975 signed-rank test. Gray lines indicate reference values as null hypotheses for each parameter. **(C)**
 976 Nutrient-specific learning rates (compared to low-nutrient baseline α_L) and the value-forgetting rate
 977 δ (compared to $\delta = 0$, perfect value memory), fitted in the *NutValAlpha-Forget* model. P-value:
 978 Wilcoxon signed-rank test. Gray lines indicate reference values as null hypotheses for each
 979 parameter. **(D)** Temporal dynamics of nutrient values. Nutrient values were plotted across
 980 chronological sessions for both monkeys. **(E)** Nutrient-based learning behavior of the *NutVal-Forget*
 981 model. Simulated choices based on the *NutVal-Forget* model reproduced nutrient-specific learning
 982 observed in **Fig. 2E**: the probability of choosing each reward (mean \pm s.e.m.) increased with previous
 983 reward outcomes but to different extents depending on reward fat and sugar content. P-values:
 984 pairwise two-sample binomial test. **(F)** Model comparison across RL models including the nine
 985 combinatorial RL models, the *Nut-RPE* model (**Fig. 5**), and the *Energy* model (**Fig. 4G**). The model
 986 comparison was conducted based on the Akaike Information Criterion (AIC) across testing sessions
 987 (mean \pm s.e.m.) for monkey Ym (blue) and monkey Ym (orange). The gray line indicates the
 988 statistical decision threshold (relative likelihood of a given model < 0.05 , see *Methods*) compared to
 989 the best-fitting *NutVal-Forget* model (red arrowhead). Comparisons between any two of the other
 990 models can also be performed by taking their AIC differences. **(G)** Psychometric curves relating
 991 model-derived reward values to choice probability. Model-fitted reward values from the *NutVal-*
 992 *Forget* model (black) outperformed the *Energy* model (*Energy RL*, red) and performed equally well
 993 as the *ObjVal* model in explaining monkeys' choices (see *Methods*). Inset: $\Delta AIC = AIC_{Energy} -$
 994 $AIC_{Nutrient}$. Pcorr: percentage correctly modelled choices \pm s.e.m. pR^2 : pseudo- $R^2 \pm$ s.e.m.. *: $P <$
 995 0.05 ; **: $P < 0.01$; ***: $P < 0.001$.

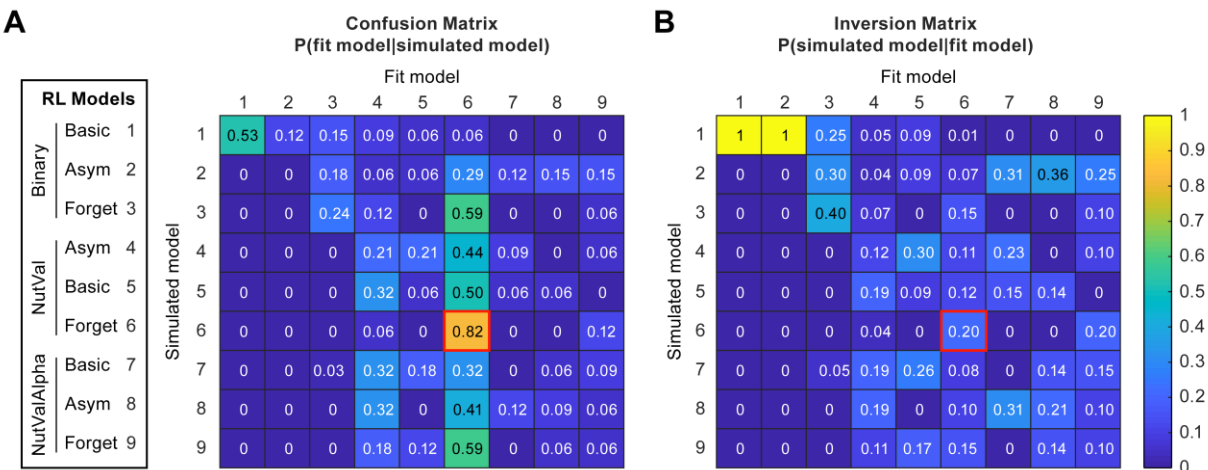


Fig. 5. Model recovery analysis. We evaluated the sensitivity and reliability of our model comparison based on **(A)** the proportion of correctly recovered models from a specified simulated model (confusion matrix) and **(B)** the confidence of best-model predictions using our approach (inversion matrix) (see *Methods*). Colors indicate the conditional probabilities stated above each matrix. The nine combinatorial RL models are listed in the left box. Red boxes in the matrices indicate data for our best-fitting *NutVal-Forget* model. The model comparison correctly identified 82% of the simulated sessions from the *NutVal-Forget* model (compare the conditional probabilities in the sixth row of **Fig. 5A**). Additionally, the *NutVal-Forget* model was the most probable generative model when predicted as the best-fitting model (compare the posterior probabilities in the sixth column of **Fig. 5B**).

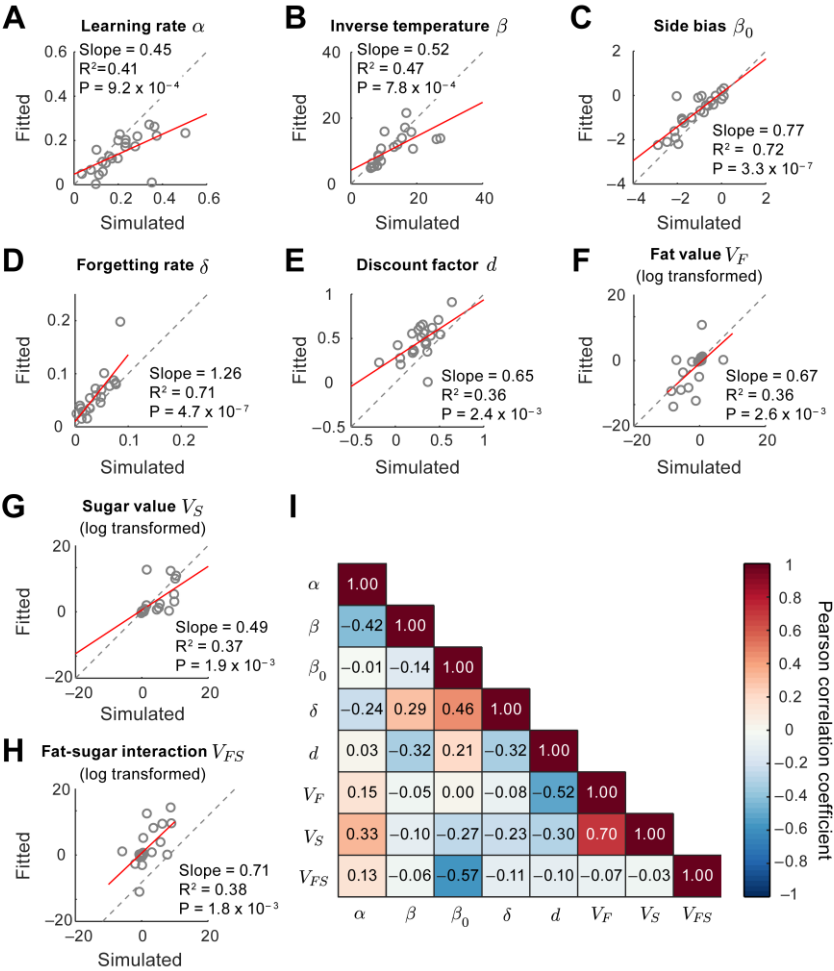


Fig 6. Parameter recovery analysis. (A-H) Correlations between simulated and fitted parameters for each of the eight free parameters in the *NutVal-Forget* model, including (A) learning rate, α ; (B) inverse temperature, β ; (C) side bias, β_0 ; (D) value-forgetting rate, δ ; (E) discount factor, d ; (F) fat value, V_F ; (G) sugar value, V_S ; and (H) fat-sugar interaction, V_{FS} . The simulation was performed based on session-specific fitted parameters from each monkey to approximate the valid range of parameters (see *Methods*). The value parameters in (F-H) were log-transformed (base 2). Dashed line: unity line. Red line: least-square regression line. Inset, slope: the fitted slope of the regression line; P-value was estimated based on the least-square linear fit of the data points. (I) Parameter trade-off. Cross-correlation between the eight fitted parameters in (A-H) identified mutual dependence between free parameters in the *NutVal-Forget* model.

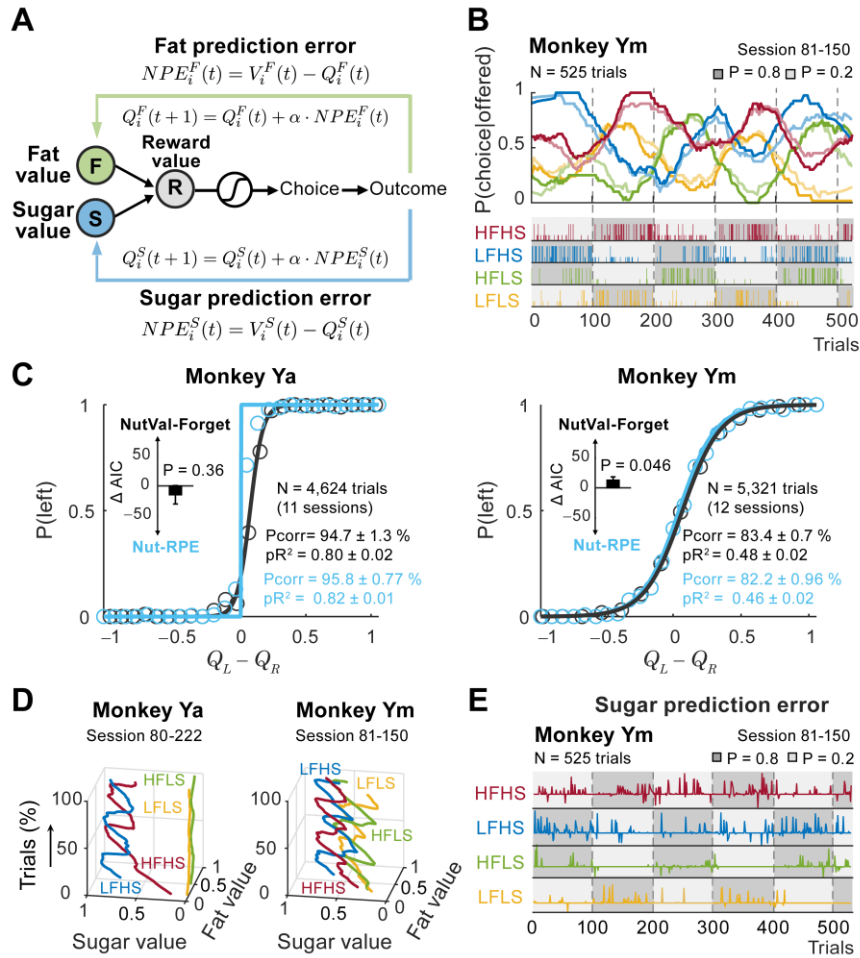


Fig. 7. Learning with nutrient-specific reward prediction errors. (A) Nutrient prediction error-based RL model (*Nut-RPE* model). Subjective values for fat (top, green) and sugar (bottom, blue) were updated based on differences between reward outcomes and expected nutrient values (nutrient prediction errors, NPEs); fat values and sugar values were multiplied into integrated reward values to guide choices. **(B)** Single-session data for choices, rewards, and modeled choice probabilities based on the *Nut-RPE* model (monkey Ym). Top: Model-predicted choices (faint lines) tracked monkey Ym's actual choices (thick lines) for nutrient-defined rewards. Bottom: the trial-by-trial record of choices and reward outcomes. Tick marks represent choices for each reward: long marks indicate large-reward outcomes ('rewarded trials'); short marks indicate small-reward outcomes ('non-rewarded trials'). **(C)** Psychometric curves based on value estimates from the *Nut-RPE* model accounted for the monkeys' choices with good model-fit indicators. Pcorr: percentage correctly modeled choices ± s.e.m. pR²: pseudo-R² ± s.e.m. **(D)** Single-session dynamics of nutrient values. Reward values in monkey Ya were dominated by the sugar values. In the *Nut-RPE* model, these sugar values were updated after choices of high-sugar liquids and tracked blockwise changes in reward probability. In monkey Ym, both fat values and sugar values contributed to value updating and tracked fluctuating reward probabilities. **(E)** Sugar prediction errors in the single session shown in panel (B). Prediction errors were sensitive to sugar content and reward size, determined by model-derived nutrient value parameters fitted to the monkey's choices. Single session data in (Fig. 7B, D, E) are the same sessions as in the right panel of Fig. 1D.

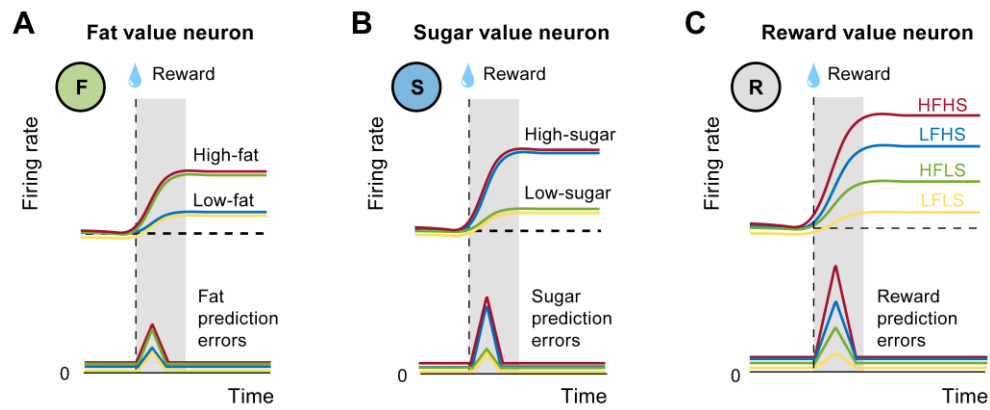


Fig. 8. Hypothesized neuron types encoding nutrient-specific learning and decision variables. (A) Fat value neurons signal the trial-by-trial fat-specific value component and update their activity based on fat-specific reward prediction errors. (B) A similar process operates for sugar-value neurons. (C) Inputs from fat and sugar value neurons may converge onto reward value neurons that signal integrated, scalar value, depending on the subjective nutrient preferences, and integrated reward prediction errors to guide learning and choice.

1049 **Table 1. Nutrient content of the liquid food rewards**

	Nutrient rewards	Peach flavor				Blackcurrant flavor				Water
		LFLS	HFLS	LFHS	HFHS	LFLS	HFLS	LFHS	HFHS	
Recipe	Peach juice (mL)	30	30	30	30	0	0	0	0	NA
	Blackcurrant juice (mL)	0	0	0	0	30	30	30	30	
	Skimmed milk (mL)	270	0	270	0	270	0	270	0	
	Whole milk (mL)	0	270	0	270	0	270	0	270	
	Caster sugar (g)	0	0.81†	20.4	21.21	0	0.81	20.4	21.21	
	Total (mL)	300	300	300	300	300	300	300	300	
Nutrient content (per 100 mL)	Calorie (kcal)	33.5	60.7	60.7	87.9	43.9	71.1	71.1	98.3	0
	Fat (g)	0.45	3.33	0.45	3.33	0.45	3.33	0.45	3.33	0
	Sugar (g)	4.50	4.50	11.30	11.30	6.80	6.80	13.60	13.60	0
	Protein (g)	3.24	3.24	3.24	3.24	3.24	3.24	3.24	3.24	0
	Salt (g)	0.11	0.10	0.11	0.10	0.14	0.13	0.14	0.13	0
	Sugar/fat (kcal/kcal)	4.444	0.601	11.161	1.508	6.716	0.908	13.432	1.815	NA

† We added 0.81 g of caster sugar for the HFLS reward to compensate for the slightly lower sugar content in the commercial whole milk compared to the skimmed milk.

1050
1051

References

- Averbeck BB, Murray EA (2020) Hypothalamic interactions with large-scale neural circuits underlying reinforcement learning and motivated behavior. *Trends in Neurosciences* 43:681-694.
- Burnham KP, Anderson DR, Huyvaert KP (2011) AIC model selection and multimodel inference in behavioral ecology: some background, observations, and comparisons. *Behavioral Ecology and Sociobiology* 65:23-35.
- Carreiro AL, Dhillon J, Gordon S, Higgins KA, Jacobs AG, McArthur BM, Redan BW, Rivera RL, Schmidt LR, Mattes RD (2016) The Macronutrients, Appetite, and Energy Intake. *Annual Review of Nutrition* 36:73-103.
- Corrado GS, Sugrue LP, Seung HS, Newsome WT (2005) Linear-Nonlinear-Poisson models of primate choice dynamics. *Journal of the Experimental Analysis of Behavior* 84:581-617.
- Costa VD, Mitz AR, Averbeck BB (2019) Subcortical substrates of explore-exploit decisions in primates. *Neuron* 103:533-545 e535.
- Cui Z, Shao Q, Grueter CC, Wang Z, Lu J, Raubenheimer D (2019) Dietary diversity of an ecological and macronutritional generalist primate in a harsh high-latitude habitat, the Taihangshan macaque (*Macaca mulatta tcheliensis*). *American Journal of Primatology* 81:e22965.
- Cui ZW, Wang ZL, Shao Q, Raubenheimer D, Lu JQ (2018) Macronutrient signature of dietary generalism in an ecologically diverse primate in the wild. *Behavioral Ecology* 29:804-813.
- Cui ZW, Wang ZL, Zhang SQ, Wang BS, Lu JQ, Raubenheimer D (2020) Living near the limits: Effects of interannual variation in food availability on diet and reproduction in a temperate primate, the Taihangshan macaque (*Macaca mulatta tcheliensis*). *American Journal of Primatology* 82.
- Dayan P (2022) "Liking" as an early and editable draft of long-run affective value. *PLoS Biology* 20:e3001476.
- de Araujo IE, Schatzker M, Small DM (2020) Rethinking Food Reward. *Annual Review of Psychology* 71:139-164.
- Gallistel CR, Fairhurst S, Balsam P (2004) The learning curve: implications of a quantitative analysis. *Proceedings of the National Academy of Sciences of the United States of America* 101:13124-13131.
- Grabenhorst F, Hernadi I, Schultz W (2012) Prediction of economic choice by primate amygdala neurons. *Proceedings of the National Academy of Sciences of the United States of America* 109:18950-18955.
- Grabenhorst F, Rolls ET, Parris BA, D'Souza A (2010) How the brain represents the reward value of fat in the mouth. *Cerebral Cortex* 20:1082-1091.
- Grabenhorst F, Tsutsui KI, Kobayashi S, Schultz W (2019a) Primate prefrontal neurons signal economic risk derived from the statistics of recent reward experience. *eLife* 8.
- Grabenhorst F, Baez-Mendoza R, Genest W, Deco G, Schultz W (2019b) Primate Amygdala Neurons Simulate Decision Processes of Social Partners. *Cell* 177:986-998 e915.
- Huang F-Y, Grabenhorst F (2022) Nutrient and sensory coding of anticipated food reward in primate amygdala neurons. *Society for Neuroscience Abstracts* 318.03.
- Huang F-Y, Sutcliffe MPF, Grabenhorst F (2021) Preferences for nutrients and sensory food qualities identify biological sources of economic values in monkeys. *Proceedings of the National Academy of Sciences of the United States of America* 118:e2101954118.
- Kadohisa M, Rolls ET, Verhagen JV (2005) Neuronal representations of stimuli in the mouth: the primate insular taste cortex, orbitofrontal cortex, and amygdala. *Chemical Senses* 30:401-419.
- Kennerley SW, Walton ME, Behrens TE, Buckley MJ, Rushworth MF (2006) Optimal decision making and the anterior cingulate cortex. *Nature Neuroscience* 9:940-947.
- Keramati M, Gutkin B (2014) Homeostatic reinforcement learning for integrating reward collection and physiological stability. *eLife* 3.
- Lak A, Stauffer WR, Schultz W (2014) Dopamine prediction error responses integrate subjective value from different reward dimensions. *Proceedings of the National Academy of Sciences of the United States of America* 111:2343-2348.
- Lau B, Glimcher PW (2005) Dynamic response-by-response models of matching behavior in rhesus monkeys. *Journal of the Experimental Analysis of Behavior* 84:555-579.

- 1108 Lau B, Glimcher PW (2008) Value representations in the primate striatum during matching behavior.
 1109 Neuron 58:451-463.
- 1110 Lee D, Seo H, Jung MW (2012) Neural basis of reinforcement learning and decision making. Annual
 1111 Review in Neuroscience 35:287-308.
- 1112 Ma CY, Liao JC, Fan PF (2017) Food selection in relation to nutritional chemistry of Cao Vit gibbons
 1113 in Jingxi, China. Primates 58:63-74.
- 1114 McNamara JM, Houston AI (1997) Currencies for foraging based on energetic gain. The American
 1115 Naturalist 150:603-617.
- 1116 Murray EA, Rudebeck PH (2013) The drive to strive: goal generation based on current needs.
 1117 Frontiers in Neuroscience 7:112.
- 1118 Murray EA, Rudebeck PH (2018) Specializations for reward-guided decision-making in the primate
 1119 ventral prefrontal cortex. Nature Reviews Neuroscience 19:404-417.
- 1120 Padoa-Schioppa C, Assad JA (2006) Neurons in the orbitofrontal cortex encode economic value.
 1121 Nature 441:223-226.
- 1122 Pastor-Bernier A, Plott CR, Schultz W (2017) Monkeys choose as if maximizing utility compatible
 1123 with basic principles of revealed preference theory. Proceedings of the National Academy of
 1124 Sciences of the United States of America 114:E1766-E1775.
- 1125 Raghuraman AP, Padoa-Schioppa C (2014) Integration of multiple determinants in the neuronal
 1126 computation of economic values. The Journal of Neuroscience 34:11583-11603.
- 1127 Rangel A (2013) Regulation of dietary choice by the decision-making circuitry. Nature Neuroscience
 1128 16:1717-1724.
- 1129 Rescorla RA, Wagner AR (1972) A theory of Pavlovian conditioning: Variations in the effectiveness
 1130 of reinforcement and nonreinforcement. In: Classical Conditioning II: Current Research and
 1131 Theory (Black AH, Prokasy WF, eds), pp 64-99. New York: Appleton Century Crofts.
- 1132 Rolls ET (2020) The texture and taste of food in the brain. Journal of Texture Studies 51:23-44.
- 1133 Rolls ET, Mills T, Norton AB, Lazidis A, Norton IT (2018) The neuronal encoding of oral fat by the
 1134 coefficient of sliding friction in the cerebral cortex and amygdala. Cerebral Cortex 28:4080-
 1135 4089.
- 1136 Rothenhoefer KM, Hong T, Alikaya A, Stauffer WR (2021) Rare rewards amplify dopamine
 1137 responses. Nature Neuroscience 24:465-469.
- 1138 Samejima K, Ueda Y, Doya K, Kimura M (2005) Representation of action-specific reward values in
 1139 the striatum. Science 310:1337-1340.
- 1140 Schultz W, Dayan P, Montague PR (1997) A neural substrate of prediction and reward. Science
 1141 275:1593-1599.
- 1142 Seo M, Lee E, Averbeck BB (2012) Action selection and action value in frontal-striatal circuits.
 1143 Neuron 74:947-960.
- 1144 Simpson SJ, Raubenheimer D (2012) The Nature of Nutrition: A Unifying Framework from Animal
 1145 Adaptations to Human Obesity: Princeton University Press.
- 1146 Simpson SJ, Raubenheimer D (2020) The power of protein. The American Journal of Clinical
 1147 Nutrition 112:6-7.
- 1148 So NY, Stuphorn V (2010) Supplementary eye field encodes option and action value for saccades
 1149 with variable reward. Journal of Neurophysiology 104:2634-2653.
- 1150 Stauffer WR, Lak A, Schultz W (2014) Dopamine reward prediction error responses reflect marginal
 1151 utility. Current Biology : CB.
- 1152 Stuphorn V (2021) Food reward derives from nutrient content and sensory qualities. Proceedings of
 1153 the National Academy of Sciences of the United States of America 118.
- 1154 Sutton RS, Barto AG (1998) Reinforcement Learning. Cambridge, MA: MIT Press.
- 1155 Suzuki S, Cross L, O'Doherty JP (2017) Elucidating the underlying components of food valuation in
 1156 the human orbitofrontal cortex. Nature Neuroscience 20:1780-1786.
- 1157 Takahashi MQ, Rothman JM, Raubenheimer D, Cords M (2019) Dietary generalists and nutritional
 1158 specialists: Feeding strategies of adult female blue monkeys (*Cercopithecus mitis*) in the
 1159 Kakamega Forest, Kenya. American Journal of Primatology 81:e23016.
- 1160 Tsutsui K, Grabenhorst F, Kobayashi S, Schultz W (2016) A dynamic code for economic object
 1161 valuation in prefrontal cortex neurons. Nature Communications 7:12554.

- 1162 van der Klaauw AA, Keogh JM, Henning E, Stephenson C, Kelway S, Trowse VM, Subramanian N,
1163 O'Rahilly S, Fletcher PC, Farooqi IS (2016) Divergent effects of central melanocortin
1164 signalling on fat and sucrose preference in humans. *Nature Communications* 7:13055.
1165 Wilson RC, Collins AG (2019) Ten simple rules for the computational modeling of behavioral data.
1166 *elife* 8.
1167 Wittmann MK, Fouragnan E, Folloni D, Klein-Flugge MC, Chau BKH, Khamassi M, Rushworth MFS
1168 (2020) Global reward state affects learning and activity in raphe nucleus and anterior insula
1169 in monkeys. *Nature Communications* 11:3771.
- 1170