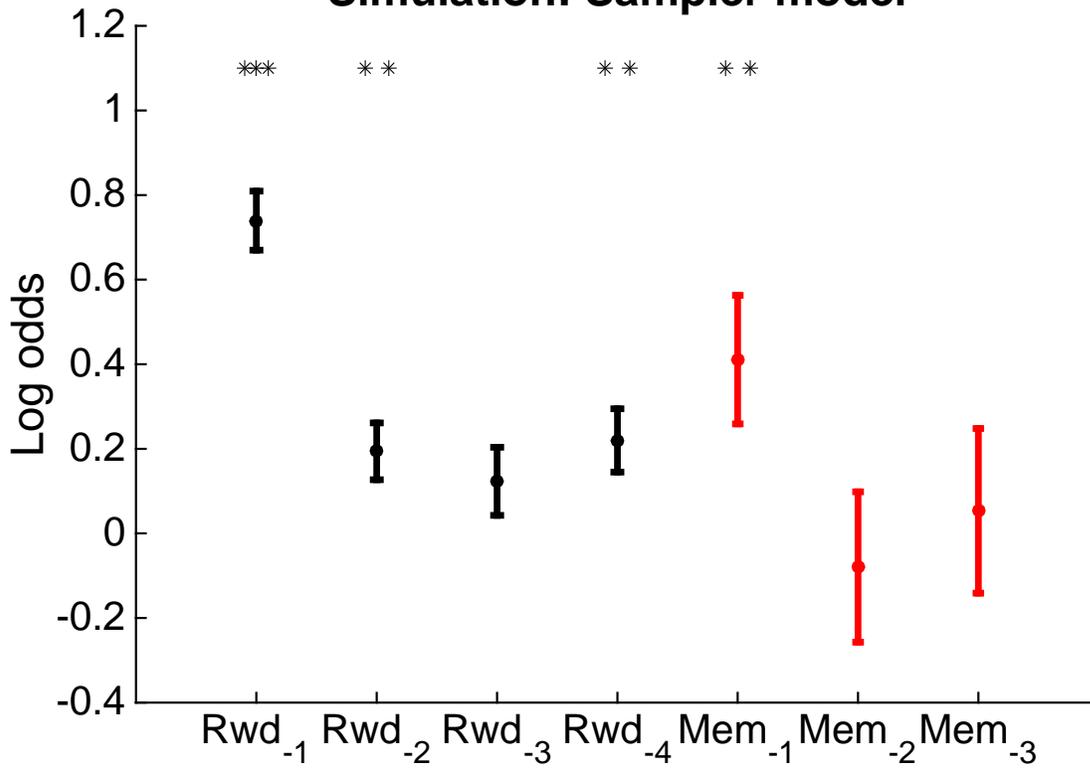Title of file for HTML: Supplementary Information
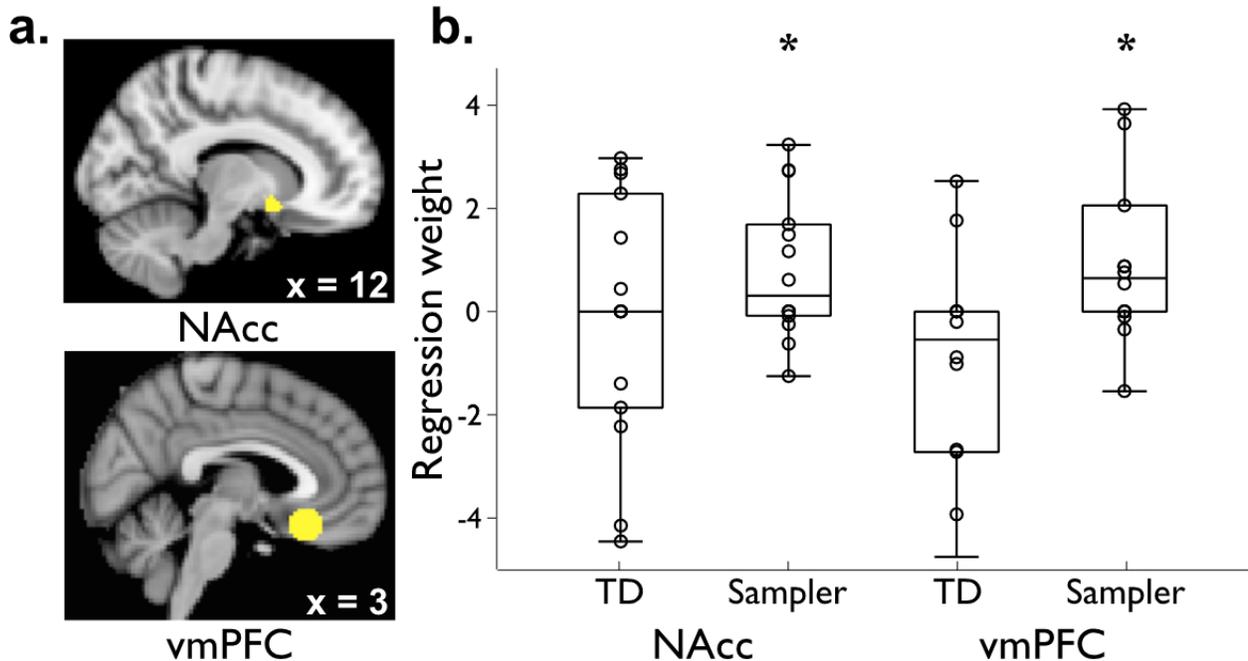Description: Supplementary Figures, Supplementary Tables, Supplementary Notes and Supplementary Reference

## Supplemental Figures

### Simulation: Sampler model



Figure S1: **Regression analysis on simulated subjects.** Average regression results for 50 populations of 20 simulated subjects each. Subjects were simulated using the Sampler model at the parameters fit to the study population. Error bars are SEM across the population means.

9

Figure S2: **Sampling model fit to neural decision variables. a. Regions of interest.** We isolated voxels of interest that corresponded to previous reports of the neural substrates for the decision variables analyzed here: Chosen Value (CV), in ventromedial prefrontal cortex (vmPFC), and Reward Prediction Error (RPE), in the nucleus accumbens (NAcc). **b. Simultaneous regression.** Candidate timeseries for each decision variable were generated according to each of the two models, and entered into a simultaneous regression against the BOLD timeseries extracted from the relevant ROI. Each plotted point represents the regression coefficient for the respective model-timeseries pair; box plots display the mean and interquartile range (* $P < 0.05$). Both regressions support the hypothesis that the Sampler model underlies neural signals (NACC-RPE: $t(13) = 2.2134$, $P = 0.0454$; vmPFC-CV: $t(13) = 2.2604$, $P = 0.0416$).

# Supplemental Tables

| Task | Simulated model | Fraction best-fit | log Bayes |
|---|---|---|---|
| Expt 1 | TD | 0.815 | 12.8726 (0.8499) |
| | Sampler | 0.897 | 8.4295 (0.7734) |
| Expt 2 | TD | 0.887 | 5.9292 (1.5065) |
| | Sampler | 0.910 | 3.5762 (0.6757) |

Table S1: **Confusion matrix for Sampler and TD models.** For each experiment and each model, we simulated 1,000 participants using the given model as ground-truth. Each individual simulated participant used a set of parameters selected at random from the parameters fit to human participants. Both models were then fit to each simulated participant's choices. Shown are the fraction of simulated participants best-fit by the ground-truth model and the mean (SEM) log Bayes factor in favor of that model.

| Model | $\alpha$ | $\alpha^{evoked}$ | $\beta$ | $\beta^c$ | log Bayes |
|---|---|---|---|---|---|
| TD | 0.5552 (0.0862) | - | 1.7551 (0.6845) | -0.0930 (0.2354) | 6.9182 (1.3227) |
| TD-evoked | 0.5269 (0.0842) | 0.2135 (0.0545) | 2.3351 (0.7223) | -0.0962 (0.2381) | 5.9167 (1.2677) |
| Sampler | 0.5393 (0.0583) | 0.4386 (0.0990) | 2.2869 (0.4943) | 0.5855 (0.3215) | - |

28 Table S2: **Fit model parameters for Experiment 2, including the TD-evoked model.** The
29 parameters shown are the mean (SEM) across subjects. The final column shows the mean (SEM)
30 of the log Bayes Factor versus the Sampler model (smaller is better).

| $\alpha^{TD}$ | $\alpha^{Sample}$ | $\beta^c$ | $\beta^{TD}$ | $\beta^{Sample}$ | $\alpha^{evoked}$ | log Bayes |
|---|---|---|---|---|---|---|
| 0.4275 | 0.5670 | 0.5281 | 0.0580 | 2.0910 | 0.6005 | 0.9700 |
| (0.0653) | (0.0521) | (0.2815) | (0.4654) | (0.5187) | (0.0871) | (1.1015) |

31 Table S3: **Fit model parameters for the Hybrid model.** The parameters shown are the mean
32 (SEM) across subjects. The final column shows the mean (SEM) of the log Bayes Factor versus
33 the Sampler model.

# Supplemental notes

## Simulated model fits

36 To demonstrate that these models are, in fact, distinguishable, we simulated the models each
37 running 1,000 instantiations of each experiment, each instance with separately initialized payoff
38 and outcome timeseries. Each model simulation was run using parameters as fit to one randomly
39 drawn subject from the respective Experiment. The Sampler model drew one sample before each
40 choice. For Experiment 2, simulated subjects responded to memory probes correctly the same
41 percentage as did our real subjects. We then fit both models to each population of 1,000
42 simulated subjects. The result of these fits is shown in Table S1.

43 For Experiment 1, subjects simulated using the TD model, 81.5% were best-fit by the TD model,
44 by an average log Bayes factor of 12.8726 (SEM 0.8499). For Experiment 1, subjects simulated
45 using the Sampler model, 89.7% were best-fit by the Sampler model, by an average log Bayes
46 factor of 8.4295 (SEM 0.7734).

47 For Experiment 2, subjects simulated using the TD model, 88.7% were best-fit by the TD model,
48 by an average log Bayes factor of 5.9292 (SEM 1.5065). For Experiment 2, subjects simulated
49 using the Sampler model, 91.3% were best-fit by the Sampler model, by an average log Bayes
50 factor of 3.5762 (SEM 0.6757).

51 In both datasets, the corresponding simulated model was a superior fit, for the bulk of the
52 population and on average at the individual level.

## Simulated regression results

54 We show that the regression results follow from the episodic sampling model. To simulate the
55 model, we generated 50 populations of 20 simulated subjects, each of whom ran a unique
56 instantiation of the task (with different payoff timeseries and outcomes), and fit the regression
57 model to each population. Simulated subjects drew one sample before each decision, used the

58  mean choice parameters as fit to the human population, and gave, on average, accurate responses
59  to memory probes at the same rate as did real subjects. Figure S1 shows the average regression
60  weights, across these populations, for each variable of interest.

## Alternative forms of choice noise

62  One potential explanation for the superior fit of the sampling model is that it simply captures
63  additional stochasticity in subjects' choices, over and above that captured by the standard
64  softmax choice function. For instance, subjects could, with some probability ε, select the highest
65  valued option, rather than selecting based on the difference in value between the two options
66  ([52]; Equation S1).

$$p_t(a = A_i) = \epsilon\left(Q_{t-1}^{TD}(A_i) = \max\left(Q_{t-1}^{TD}(\cdot)\right)\right) + (1 - \epsilon)\left(\frac{e^{\beta^C I_t^C + \beta^{TD} Q_t^{TD}(A)}}{\sum_j e^{\beta^C I_t^C + \beta^{TD} Q_t^{TD}(a_j)}}\right) \qquad (S1)$$

68  However, model comparison did not provide evidence in favor of this "ε-greedy" approach. In
69  Experiment 1 the Sampler model was favored for 15/20 subjects, by a mean Bayes Factor of
70  3.0042 (SEM 1.8137, exceedance probability > 0.99), while in Experiment 2 the Sampler was
71  favored for 19/21 subjects, by a mean Bayes Factor of 6.734 (SEM 1.7078, exceedance
72  probability > 0.99).

## Neuroimaging reanalysis

74  Given the Sampler's superior fit to behavior, we used the neuroimaging data collected alongside
75  Experiment 1 [4] to ask whether the expectation and learning variables predicted by the sampling
76  model could provide a better explanation of BOLD signal than did the corresponding variables
77  extracted from a TD model. Specifically, we tested whether the well-studied neural correlates of
78  key decision variables–Chosen Value (CV) and Reward Prediction Error (RPE)–were better
79  predicted by the sampling model than by TD. We first identified regions of interest (ROIs)
80  encompassing areas previously shown to reflect this activity: ventromedial prefrontal cortex /
81  medial orbitofrontal cortex (hereafter: vmPFC) for chosen value, and Nucleus Accumbens
82  (NAcc) for RPE (Figure S2a). Extracting the timeseries of activity within these ROIs, we next
83  performed a simultaneous regression containing the timeseries of variables predicted by both
84  models, along with several regressors of no interest. Comparing the distribution of resulting per-
85  participant regression weights against zero using a two-tailed, one-sample t-test, we evaluated
86  whether each model was a significant predictor of the target BOLD timeseries.

87  The regressors were slightly, but reliably, correlated between models (for RPE: mean $R$ =
88  0.1093, $P = 0.0215$; for CV: mean $R = 0.2701$, $P = 0.0002$). To test the exclusive contribution of
89  the Sampler-derived predictor variables, we orthogonalized the RPE and CV timeseries as
90  generated using the Sampler against their TD counterparts, and entered each set of TD and
91  Sampler predictors into a simultaneous regression on the BOLD timeseries.

92  In both cases, the predictions of the Sampler model captured additional variance in the BOLD

93  timeseries that was not modeled by TD (Figure S2b). Across participants, regression on the
94  NAcc timeseries revealed significant contribution of the RPE variable as generated by the
95  Sampler model ($t(13) = 2.2134$, $P = 0.0454$), but not TD ($t(13) = -0.1614$, $P = 0.8742$).
96  Similarly, the Chosen Value regressor generated by the Sampler model was a significant
97  predictor of the vmPFC timeseries ($t(13) = 2.2604$, $P = 0.0416$), while the TD version was not
98  ($t(13) = -1.0835$, $P = 0.2983$).

## Adding evoked trials to the TD model

100  We augmented the TD model to incorporate rewards from bandit trials evoked by valid memory
101  probes. Specifically, we added an additional parameter, $\alpha^{evoked}$, for a new, augmented TD
102  update applied to rewards $r_i^{evoked}$ during memory probe trials (Equation S2). This parameter
103  allowed the weight given to evoked bandit outcomes to vary, reflecting the fact that the sampling
104  mechanism may itself be stochastic in nature — not every probe trial will successfully trigger a
105  recall of the associated context, even those on which participants exhibit correct recognition
106  memory.

$$Q_t^{TD}(a) = Q_{t-1}^{TD}(a) + \alpha^{evoked}(r_i^{evoked} - Q_{t-1}^{TD}(a)) \tag{S2}$$

108  Table S2 expands the comparison from the main text to include this TD-evoked model. After
109  correcting for the additional parameter, the model was a slightly better fit to participants'
110  behavior than the plain TD model. It was not, however, a superior fit to the Sampler model.

## Hybrid Sampler and TD model

113  The evidence across several studies shows that multiple valuation systems contribute to choices,
114  either simultaneously or across time [6,12,35]. The current study provides evidence that one
115  component of this valuation architecture involves evaluating samples from episodic memory. (Of
116  course, we do not know to what extent this influence is coextensive with, e.g. model-based
117  learning as otherwise defined.)

118  To test the possibility that such a "hybrid" model could account for choices in this task, we
119  implemented a hybrid model combining the TD and episodic sampling models and tested it on
120  Experiment 2. The model has six parameters: a learning rate for the TD component $\alpha^{TD}$, a decay
121  rate for the Sampler component $\alpha^{sample}$, a choice stickiness term $\beta^c$, a softmax temperature for
122  the TD value $\beta^{TD}$, a softmax temperature for the $\beta^{Sample}$, and a decay rate for evoked trials
123  $\alpha^{evoked}$. Choice probability is computed using Q-values derived from each model, as specified
124  in the main text, and taken over all possible combinations of samples (following Equation 5 in
125  the main text).

126  Parameter estimates from this model show that that most of the weight is on the sampler model
127  (in the sense that it has a much higher softmax temperature, i.e. its values have a larger effect on
128  choice – indeed, the weight on the TD model is not reliably different from zero). Accordingly,
129  the addition of a TD component to the sampler model was not robustly justified in light of the

130   additional free parameters. On average, across subjects, log Bayes Factors were mildly in favor
131   of the Sampler model (mean 0.9700, SEM 1.1015, exceedance probability 0.8999), and the
132   Sampler was a better fit for 13 out of 21 subjects individually. The fit parameters and model
133   comparison results are shown in Table S3.

134   In sum, after accounting for the additional free parameters, this hybrid model was not a clearly
135   superior fit to behavior than the episodic sampling model taken alone. However, we think that–in
136   an experiment designed to distinguish these two possibilities–a more sophisticated architecture
137   (perhaps employing a common cached value representation as in DYNA [36]) could possibly
138   prove a superior explanation of behavior.

# References

141   .  [53]  Hanan Shteingart, Tal Neiman, and Yonatan Loewenstein. The Role of First Impression
142          in Operant Learning. Journal of Experimental Psychology: General, 142(2):476–488,
143          Aug 2012. ISSN 1939-2222. doi: 10.1037/a0029550. URL
144          http://www.ncbi.nlm.nih.gov/pubmed/22924882.