# Supplemental materials for "Putting value in context: A role for context memory in decisions for reward".

Aaron M. Bornstein <sup>1\*</sup>, Kenneth A. Norman <sup>1,2</sup>

- 1. Neuroscience Institute, Princeton University, Princeton, NJ, USA.
- 2. Department of Psychology, Princeton University, Princeton, NJ, USA.
- \* To whom correspondence should be addressed: aaronmb@princeton.edu.

## Supplemental Methods

#### Context-aware sampling model

To investigate how single-trial and context reward trade off with each other as the number of past episodes sampled increases, we simulated the task as performed by a context-aware episodic sampling model. In this simulation, all choices are made using episodic sampling alone (no influence of model-free values), so as to clearly isolate the influence of changing the number of samples. In episodic sampling, option values are estimated using the values encountered on one or more past episodes, with the likelihood of sampling a given episode diminishing exponentially with its recency. The context-aware episodic sampling model augments this idea, by positing that additional samples after the first are (with some probability) selected uniformly from the same context as the preceding sample.

The model maintains a cache of episodes representing each experienced trial. When a subject responds correctly to a valid memory probes, the model with some probability "reinstates" the episode by copying the reminded trial to the front of the cache—thus, making it more likely to be drawn when evaluating options during the next choice. If the subject's correct response to the memory probe is of high confidence, then the context of the probed trial is included in the episode copied to the front of the cache.

We used this model to simulate subjects who used different numbers of samples from episodic memory to make decisions. The model had four parameters which were fixed across all simulations:  $\alpha_{direct}$ , or the decay rate on temporal recency;  $\alpha_{evoked}$ , or the probability of reinstating evoked trials as a result of memory probes;  $\beta$ , the softmax temperature;  $\beta_p$ , the choice perseveration term; and  $\pi$ , the likelihood of drawing sample k from the same context as sample k-1 (as opposed to based on temporal recency).

A final parameter, the number of samples drawn to make a decision, was varied between 1 and 15. For each fixed number of samples, we simulated 1,600 subjects each performing an instantiation of the task. The simulated subject's parameters  $\alpha_{direct}$ ,  $\beta$ , and  $\beta_p$  were set to those fit to the real subjects with the number of samples equal to 1, and  $\pi$  was set to 1.

For simplicity, in the first simulation  $\alpha_{evoked}$  was set to 1. This assumption was relaxed for our second set of simulations, during which we varied  $\alpha_{evoked}$  between 0 and 1 to illustrate the impact of changing this parameter. The subjects were programmed to make, on average, the same proportion of low confidence and incorrect responses as did the subjects from Experiments 1 and 2.

#### Fitting incremental learning models to our data

To compare our context reward model to the others, we generated the timeseries of reinstated context reward values that would be learned according to three different specifications of model-free RL.

The first variant followed the traditional method, learning the value of each card deck without regard to changes in room context. Specifically, at each step, the value for the chosen card deck,  $B_t$ , was updated according to the reward received on that trial,  $R_t$ , and the learning rate  $\alpha$ :

$$Q_B = Q_B + \alpha * (R_t - Q_B) \tag{S1}$$

We refer to this as the *standard* model.

The second variant reset the value of each card deck when the room changed, giving separate values  $Q_{B,C}$  for each room-context, following Equation S2.

$$Q_{B,C} = Q_{B,C} + \alpha * (R - Q_{B,C}) \tag{S2}$$

We refer to this as the *room-reset* model.

A third variant also reset action values, but this time it did so at a variable trial number within each room (e.g. 1, 2, 3, or more trials after room change). We designed this model to account for the possibility that participants took note of payoff reversals and discounted information prior to the switch [1]. We refer to this model as the *reversal* model.

The three incremental learning models fit to choices in rooms one through six, using maximum likelihood estimation of parameters. To compare the model likelihoods we first transformed them using the Bayesian Information Criterion (BIC; Schwarz 2), which penalized the third model for its additional parameter.

The best-fitting model was used to generate a replacement for the context-reward regressor as specified in main text, *Methods* subsection *Regression analysis*. The resulting regression model weights were compared to the original (shown in Figure 5 in the main text).

### Supplemental Results

### Sampling model simulations

We demonstrated the context-aware sampling model (described in *Supplemental Methods*, subsection *Context-aware sampling model*) predicts that the effect of reminded context

should be greater than the effect of reminded trials (as in Figure 5 of the main text). To show this we simulated the context-aware sampling model 1,600 times, and split the sampled subjects into 50 populations of 32 each. (Figure S2 shows the results for one simulated population of 32 subjects each sampling 12 episodes in the time between probe and choice.) As the number of samples to make a decision increased, so did the ratio of the effect of context reward to that of the reminded trial (correlation between context and single-trial effects: R = -0.9436, p = 1.3125e - 07; Figure S1).

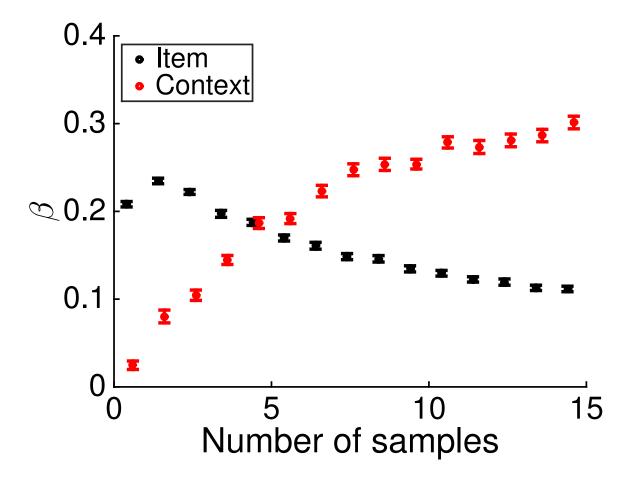


Figure S1: Simulations show that the influence of context reward should increase with greater numbers of episodic samples. We simulated the context-aware sampling model, holding fixed all parameters except the number of samples drawn in support of each decision. We then performed the regression analysis shown in Figure 5 of the main text on this simulated data, and plot here the regression weights for single-trial (item) rewards and context rewards. As the number of samples drawn increases, the effect of context reward increases, while the effect of item reward decreases.

The effects of both reminded trials and reminded context are lower in Experiment 2 than Experiment 1, relative to the effect of recent reinforcement. In Experiment 1, the

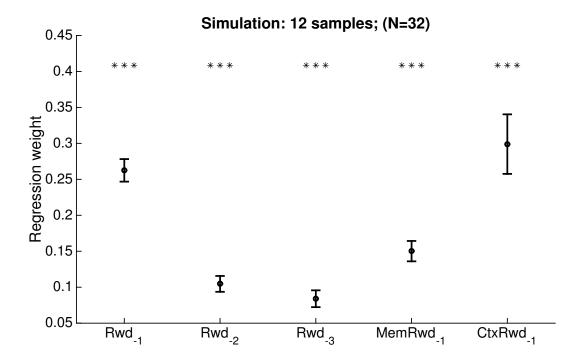


Figure S2: Full regression model fit to a simulated population of 32 subjects that sampled 12 episodes between each probe and the ensuing choice.

context reward effect is not significantly different from the effect of the most recent reward (p=0.6601) nor the second most recent (p=0.2449), but is significantly larger the effect of the third most recent reward (p=0.0128). In Experiment 2, the context reward effect is lower than the effect of the most recent reward (p=0.0027), but not distinguishable from the second most recent (p=0.3148) nor the third most recent reward (p=0.8191). In other words, subjects in Experiment 2 showed a lower influence of rewards from probe-evoked information, and thus were relatively more likely to make decisions on the basis of recent reinforcement.

We used the simulation to illustrate how our mechanism can eliminate the effect of reminded trials while sparing an effect of recent rewards and context. We followed the simulation procedure above, this time generating populations of 1,600 subjects, with each population using a different value of  $\alpha_{evoked}$ . All other model parameters were fixed. Figure S3 shows the results for one population of 32 simulated subjects with  $\alpha_{evoked}$  set to 0.3.

Figure S4 shows how, as the value of  $\alpha_{evoked}$  decreases, so do both memory effects. Crucially, though, because the reminded context has a larger impact on behavior, relative to the reminded trial (for reasons discussed above), it remains impactful, while the item effect drops to nearly zero.

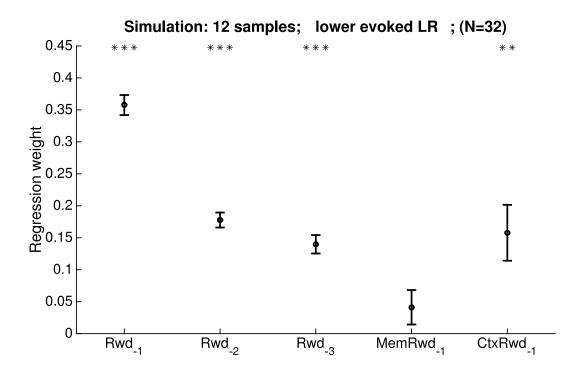


Figure S3: Full regression model fit to a simulated population of 32 subjects that sampled 12 episodes before each choice, with  $\alpha_{evoked}$  set to 0.3.

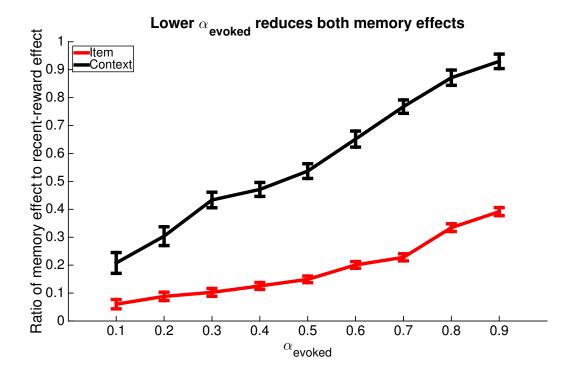


Figure S4: As the probability of reinstating reward information decreases, so does the contribution of both item and context memory to decisions. We ran the experiment for 9 populations of 1,600 simulated subjects. Each population used a different value of the  $\alpha_{evoked}$  parameter, the probability that a reminder probe would result in reward information being reinstated from the reminded trial and context. We then measured the regression weights for the reminded trial and the reminded context, plotted here by their ratio to the regression weight for reward experienced one trial ago. As  $\alpha_{evoked}$  decreased, so did the ratio of each type of memory-based effect to the effect of recent rewards. While the effect of reminded context is preserved at even very small values for  $\alpha_{evoked}$ , the effect of the reminded trial drops to near zero. This pattern matches the reduced effect of both reminded trial and context in Experiment 2.

#### Fitting incremental learning models to our data

First, we compared the fit to behavior of three alternative RL models. Each model learned cached action values for the three card decks; two of these models reset those values when context changed—in one model, the context shifts/value resets were at the time that the room changed, in the other model, the context shifts were at a variable trial after the start of each room (to account for context boundaries being drawn at the time the payoffs changed).

The model that reset action values when the room changed was the best fit to behavior. By BIC versus the second-best model, in Experiment 1, the room-reset model was superior for 18/20 subjects, mean difference in BIC: 5.6981; In Experiment 2, 26/32 subjects, 4.6406.

The fit parameters for the room-reset model-learning rate  $\alpha$ , softmax temperature  $\beta$ ,

choice stickiness  $\beta_p$ —were: Expt. 1  $\alpha$  mean 0.4802 SEM 0.0617,  $\beta$  mean 0.2333 SEM 0.2829,  $\beta_p$  mean 0.4348 SEM 0.1392; Expt. 2  $\alpha$  mean 0.5738 SEM 0.0375,  $\beta$  mean 0.4623 SEM 0.0313,  $\beta_p$  mean 0.2231 SEM 0.0738. These values were consistent across the two experiments (by two-sample t-test:  $\alpha p = 0.1745$ ;  $\beta p = 0.3156$ ,  $\beta_p p = 0.1480$ ).

We then used this room-reset model to generate a replacement context reward regressor, used in place of context reward the final action-values learned for each card deck in each room. The mean correlation between this new regressor and the corresponding original regressor was R=0.0564 (p=0.3926) in Expt. 1 and R=0.1714 (p=5.4728e-06) in Expt. 2.

# Context reward effects are specifically modulated by scene evidence in PPA

To confirm that context reward was specifically modulated by scene reinstatement, we repeated the quartiles analysis of Figure 5 in the main text to look for an effect of univariate activity in PPA (as opposed to scene evidence). Across quartiles and subjects, neither reminded trial (mean slope= -0.0422, SEM 0.0397, p = 0.3154), nor reminded context (mean slope= -0.0162, SEM 0.0184, p = 0.3875) showed a relationship with PPA activity, nor did we observe an interaction between these relationships (p = 0.7559).

We next defined a region of interest that was differentially responsive to the "scrambled" scenes that were used as a control in our localizer task (section Methods, subsection ROI definition). Across quartiles and subjects, neither reminded trial (mean slope = 0.0194, SEM 0.0479, p = 0.6956), nor reminded context (mean slope = -0.0214, SEM 0.0212, p = 0.3210) showed a relationship with univariate activity levels in this ROI, nor did we observe an interaction between these relationships (p = 0.5047). We next trained a classifier (analogously to how we trained the scene evidence classifier in the main analysis) to decode both "scrambled scene" evidence and scene evidence from this ROI, and explored how these classifier evidence values related to trial and context effects in our regression model. Across quartiles and subjects, neither reminded trial (mean slope= -0.0033, SEM 0.0160, p = 0.8391), nor reminded context (mean slope= 0.0314, SEM 0.0200, p = 0.1264) showed a relationship with "scrambled scene" evidence in this ROI, nor did we observe an interaction between these effects (p = 0.1467). Across quartiles and subjects, neither reminded trial (mean slope= 0.0041, SEM 0.0132, p = 0.7602), nor reminded context (mean slope= 0.0185, SEM 0.0158, p = 0.2518) showed a relationship with scene evidence in this ROI, nor did we observe an interaction between these effects (p = 0.5046).

As described in the main text, we also tested whether each regressor of interest was modulated by scene evidence, by repeating the quartiles analysis of Figure 5 for each regressor. We found that scene evidence did not reliably modulate any of the other regressors of interest (all p > 0.14; Figure S5).

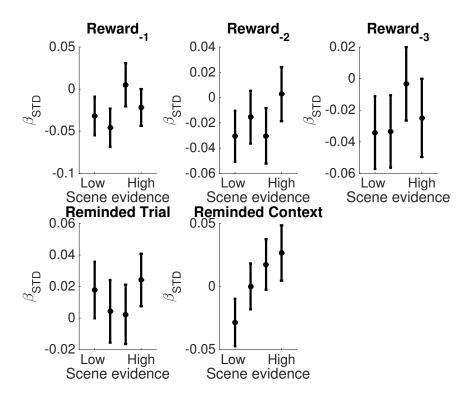


Figure S5: Each regressor of interest plotted as a function of scene reinstatement evidence.

#### Hippocampal activity during decision-making

The hippocampus has long been understood to be critical to episodic memory encoding and retrieval. We investigated whether our regression weights were correlated with univariate activity levels, or classifier evidence for scenes, in this structure (Figure S6).

First, using the FreeSurfer atlas, we defined for each subject an anatomical mask of bilateral hippocampus proper.

Next, at each correct, valid probe trial, we took the average activity level in hippocampus across our timepoints of interest (Figure 6b of the main text). We then divided trials into quartiles based on this measure, and recomputed the regression in Figure 6c of the main text for each bin. We repeated the regression analysis within each of these quartiles.

While the relationship between hippocampal activity levels and single-trial reward effect was numerically negative (average slope = -0.0231 SEM 0.0288, p = 0.4444), and the context reward effect was numerically positive (average slope = 0.0277 SEM 0.0285, p = 0.3422), neither were statistically significant, not was there an interaction between the effects (p = 0.2928).

We next tested for a relationship between scene reinstatement evidence in hippocampus and the regression effects. For this analysis, we trained a scene classifier on multi-voxel patterns in the hippocampus, following a procedure analogous to that used for PPA in the main text. We then divided trials into quartiles based on this measure, and recomputed the regression analysis (Figure 6c of the main text) for each quartile. We did not observe a reliable relationship between hippocampal scene evidence and the effect of either reminded trial or reminded context (both p > 0.12).

Lastly, we investigated the hypothesis that hippocampal activity could reflect retrieval of memories in support of decisions [3, 4]. We used the scene-specific reinstatement weights to investigate activity related to memory retrieval. In previous studies, we observed that hippocampal activity increases along with uncertainty about an action's outcome, both for simple sequential responses and goal-directed planning decisions [5, 6]. We interpreted these findings as consistent with hippocampus' known role in memory retrieval. In the context of action evaluation, memory retrievals could constitute evidence about the outcome of the actions under consideration (similar to how forward trajectories are "replayed" as rodents make navigation decisions [7]). In this task, greater uncertainty about the associated context should lead to a wider range of next-step outcomes to evaluate. We therefore reasoned that activity in hippocampus at choice might scale with uncertainty about the probed item's context. To test this, we computed the entropy of the distribution of scene-specific weights on each trial. We then measured, for each participant, the correlation between this trialby-trial entropy and univariate activity in hippocampus, defined broadly using a bilateral anatomical mask. Consistent with a role for hippocampus in retrieving memories that are used to evaluate outcomes, this correlation was reliably positive across participants: mean R = 0.0670, SEM 0.0261,  $t_{29} = 2.5733$ , p = 0.0155 (Figure S7).

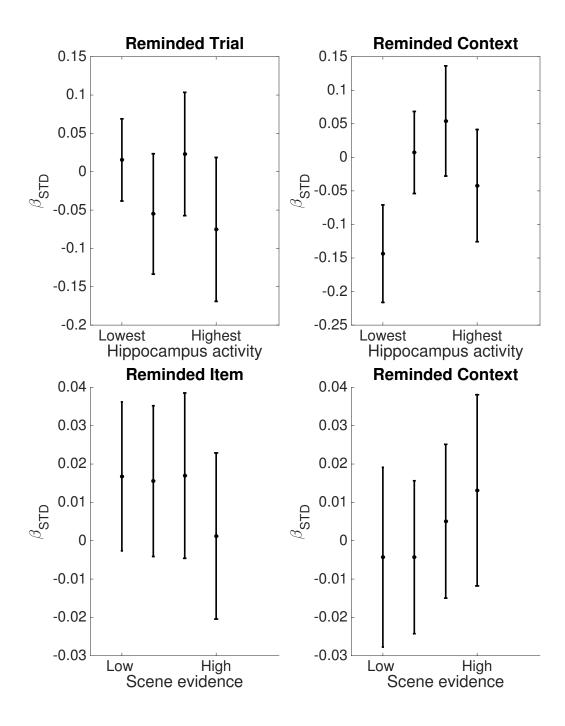


Figure S6: **Hippocampal activity does not predict behavioral effects.** The effect of reminded trial and of reminded context reward as a function of univariate activity and scene evidence in hippocampus.

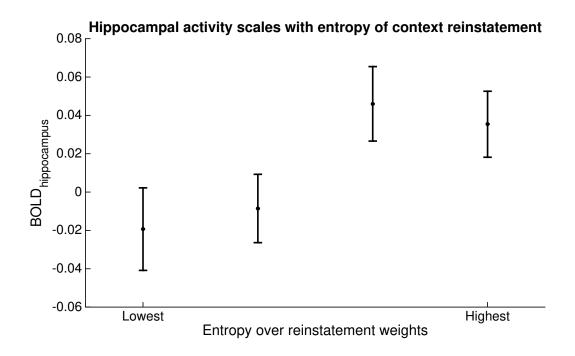


Figure S7: Activity in bilateral hippocampus scales with entropy over scene reinstatement evidence. We computed the evidence that participants reinstated each context image on each trial, by taking the correlation between per-scene template patterns and PPA activity on that trial. The resulting six numbers were then normalized to create a probability distribution for that trial, reflecting the relative likelihood that each scene was being remembered. The entropy over this distribution can thus be considered the uncertainty over the reinstated context. This entropy value is reliably correlated with hippocampal activity at the time of choice (mean R = 0.0670, SEM 0.0261,  $t_{29} = 2.5733$ , p = 0.0155). Shown here are the mean hippocampal activity levels for each quartile of entropy values; quartiles are computed within-participant, and the means taken across the population.

#### References

- [1] Timothy E J Behrens, Mark W Woolrich, Mark E Walton, and Matthew F S Rushworth. Learning the value of information in an uncertain world. *Nature Neuroscience*, 10(9):1214–21, sep 2007. ISSN 1097-6256. doi: 10.1038/nn1954. URL http://www.ncbi.nlm.nih.gov/pubmed/17676057.
- [2] Gideon Schwarz. Estimating the Dimension of a Model. *Annals of Statistics*, 6(2):461–464, 1978.
- [3] Aaron M Bornstein, Mel W Khaw, Daphna Shohamy, and Nathaniel D. Daw. What's past is present: Reminders of past choices bias decisions for reward in humans. *bioRxiv*, 2015. doi: 10.1101/033910.

- [4] Michael Ν Shadlen Daphna Shohamy. Decision Making and Sequential and Sampling from Memory. Neuron, 90(5):927-939, 2016. ISSN 0896-6273. doi: 10.1016/j.neuron.2016.04.036. URL http://dx.doi.org/10.1016/j.neuron.2016.04.036.
- [5] Aaron M. Bornstein and Nathaniel D. Daw. Dissociating hippocampal and striatal contributions to sequential prediction learning. European Journal of Neuroscience, 35(7): 1011–1023, apr 2012. ISSN 0953816X. doi: 10.1111/j.1460-9568.2011.07920.x. URL http://doi.wiley.com/10.1111/j.1460-9568.2011.07920.x.
- [6] Aaron M. Bornstein and Nathaniel D. Daw. Cortical and Hippocampal Correlates of Deliberation During Model-Based Decisions for Rewards in Humans. *PLoS Computational Biology*, 9(12):e1003387, dec 2013. ISSN 1553-7358. doi: 10.1371/journal.pcbi.1003387. URL http://dx.plos.org/10.1371/journal.pcbi.1003387.
- [7] Adam Johnson and A. David Redish. Neural ensembles in CA3 transiently encode paths forward of the animal at a decision point. *Journal of Neuroscience*, 27(45): 12176–89, nov 2007. ISSN 1529-2401. doi: 10.1523/JNEUROSCI.3761-07.2007. URL http://www.ncbi.nlm.nih.gov/pubmed/17989284.