Independent Neural Computation of Value from Other People's Confidence
Campbell-Meiklejohn *et al*, 2016, Journal of Neuroscience
Supplemental Material: Bayesian Framework. Version 1.

**Bayesian Inference Solution of Social Uncertainty Using Others' Choices and Confidence**

To develop an initial understanding of the problem and a solution, we created an ideal observer model that determined the probability that a marble from a mixed urn (red and green marbles) would be red. The goal of this is a *qualitative* illustration of how the computations might work, rather than a quantitative fit or a process-level mechanistic description. The model below is a generative model that represents how the subject believes different observations (marbles, agent choices, agent confidences) relate to the contents of the urn. An ideal observer inverts this generative model using Bayes' rule to infer the probable state of the urn given a set of observations, and thus predict the color of the next marble drawn from that urn.

The observations include a sample of 8 marbles from the urn (containing a number of red marbles), the choices of four agents, and each agent's confidence in his or her choice. From the choices and confidences of agents, the ideal observer infers the probability distribution of each agent's sample, and together with the observer's own sample, the likely composition of the urn.

We implemented inference in the generative model using the PyMC package in Python (script available in same space as this supplemental material), which uses Markov Chain Monte Carlo techniques to draw samples from the posterior distribution of the unobserved quantities given the generative model and the observed quantities. The following linked functions were in the generative model:

Given the probability of drawing red from the urn, $P_{red}$, the probability of the subject's sample of 8 marbles:

$$\text{Subj}_{samp} \sim \text{Binomial}(8, P_{red})$$

Similarly, the probability of each of the four agents' samples:

$$O_{samp} \sim \text{Binomial }(8, P_{red}), \text{ for each agent}$$

Given the probability of each agent's sample within $O_{samp}$, the confidence of each agent, defined as the probability assigned by each agent that at least 50% of marbles in the urn are red, drawn from $I_{0.5}$, which is the cumulative distribution function of a beta distribution evaluated at 0.5 (then adjusted to -0.5 confident green, 0.5 confident red):

$$O_{conf} = I_{0.5}(8 - O_{samp} + 1, O_{samp} + 1) - 0.5, \text{ for each agent}$$

Probability that an agent chooses red, where $k_{red}$ is a constant, for which we used 10:

$$P_{RedChoice} = \text{logit}^{-1}(k_{red} * O_{conf}), \text{ for each agent}$$

Probability that an agent responds with confidence, where $k_{conf}$ is a constant for which we used 10:

$$P_{conf} = \text{logit}^{-1}(k_{conf} * (|O_{conf}| - 0.25)), \text{ for each agent}$$

Each agent's actual choice of colour for the next marble to be drawn:

$$O_{RedChoice} \sim \text{Bernoulli}(P_{RedChoice}), \text{ for each agent}$$

Each agent's expression of confidence:

Independent Neural Computation of Value from Other People's Confidence
Campbell-Meiklejohn *et al*, 2016, Journal of Neuroscience
Supplemental Material: Bayesian Framework. Version 1.

$$O_{conf} \sim Bernoulli(P_{conf}), \text{ for each agent}$$
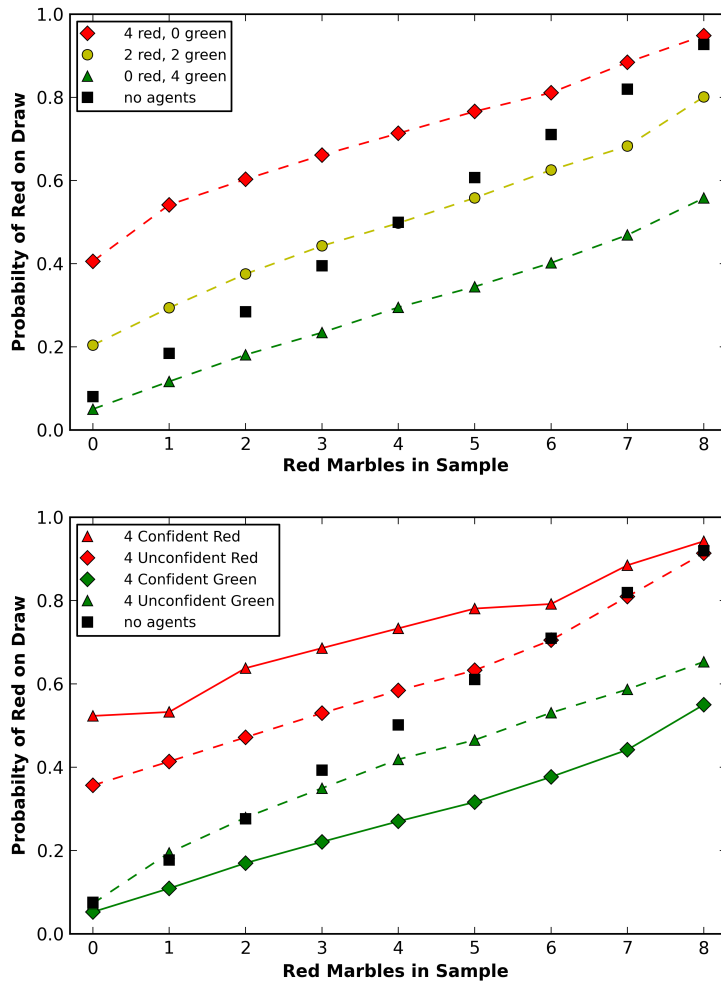
## Results of Model Simulations



**Figure S1:** Simulated estimations of probability of drawing red from a Bayesian Ideal Observer Model. We implemented inference in our generative model using Markov Chain Monte Carlo techniques to draw samples from the posterior distribution of the unobserved quantities given the generative model and the observed quantities. We chose an observer's sample size of 8 marbles and four agents as in the Urn task, and constants of $k_{red}$=10 and $k_{conf}$ = 10 which mediate the shape of the logistic function determining agent choices and confidence statements from agent beliefs about the urn. Plot points reflect the median of 25000 samples following 5000 iterations of burn-in. **Top**: Effect of agent choices of red and consistency of agents, across a range of personal samples. **Bottom**: Additional effect of agent confidences.

The Bayesian inference model of the hypothesized inferences for optimal task performance successfully demonstrated the contribution of each distinct information source to simulated subject choices (Figure 2). The key feature of this model was that the observer infers each other agent's private information about the urn (i.e., the agent's sample) from a combination of her choice and confidence, and combines this across agents and with the directly observed sample to infer the likely composition of the urn. The inferred proportion of red marbles in the urn –

Independent Neural Computation of Value from Other People's Confidence
Campbell-Meiklejohn *et al*, 2016, Journal of Neuroscience
Supplemental Material: Bayesian Framework. Version 1.

and thus the expected likelihood of success for choosing red – increases monotonically with the number of red marbles observed, and also with the number of red guesses by others. Most crucially, though, the effect of others' guesses is modulated by their confidence: the more confident an agent is in her choice, the more likely she has seen a skewed (and therefore informative) sample, and the more her choice should influence the observer's. In this way, confidence can mediate the impact of another's choice on one's own beliefs, via mentalizing.

**Bayesian Framework at Work in the Brain**

This model could also be used to generate the value of the option chosen by each subject, had they computed this value (probability that their chosen option is correct) as per Bayesian Inference. Script for these values is deposited in the same online location as this document. This value was entered into the fMRI analysis to test whether the brain's value representations reasonability conform to this model, and thus tests the viability of the Bayesian framework. See Figure 2. Increased expected value from the model predicts vMPFC activity and decreasing predicts dMPFC activity, as would be expected if it were doing an adequate job of predicting value. This shows that the Bayesian framework may be helpful, but does not suggest that it is precisely how the brain is performing the value computation during the task.
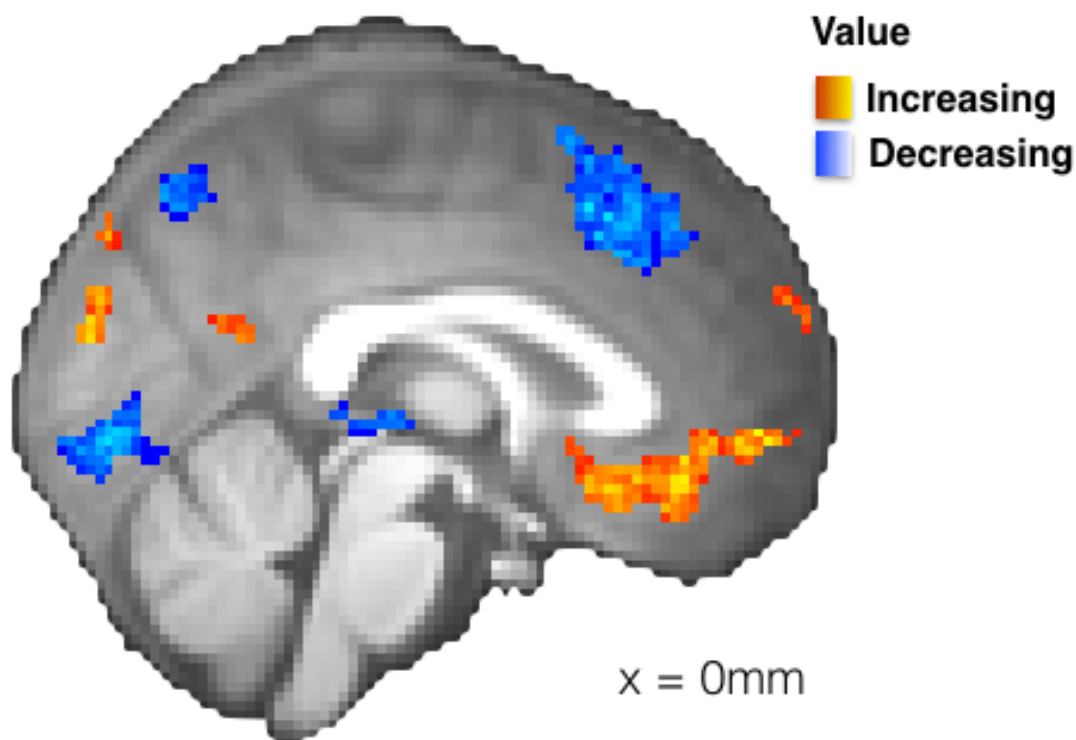


**Figure S2.** Neural correlates of chosen value as generated by Bayesian inference model. Clusters of voxels exceeding a z > 3.0 threshold, cluster significance set at p <0.05. Within vMPFC peak voxels for increasing value are [coordinates: -4 64 20, Z max = 4, 139 voxels] and [coordinates: -4 24 -12, Z max = 4.29, 759 voxels]. Peak voxel within MPFC for decreasing value are [coordinates: 6 32 42, Z max = 5.46, 695 voxels].