

# A Bayesian Account of Reconstructive Memory

**Pernille Hemmer (phemmer@uci.edu)**

Department of Cognitive Sciences, University of California, Irvine  
Irvine, CA, 92697-5100

**Mark Steyvers (msteyver@uci.edu)**

Department of Cognitive Sciences, University of California, Irvine  
Irvine, CA, 92697-5100

## Abstract

It is well established that prior knowledge influences reconstruction from memory, but the specific interactions of memory and knowledge are unclear. Extending work by Huttenlocher et al. (1991, 2000) we propose a hierarchical Bayesian model of reconstructive memory in which prior knowledge interacts with episodic memory at multiple levels of abstraction. The combination of prior knowledge and noisy memory representations is dependent on familiarity. We present empirical evidence of the hierarchical influences of prior knowledge, showing that the reconstruction of familiar objects is influenced toward the specific prior for that object, while unfamiliar objects are influenced toward the overall category.

**Keywords:** Long term memory; Prior knowledge; Bayesian models; Reconstructive memory

## Introduction

Knowledge is essential for our interactions with the environment. We learn more easily by using what we know to relate to new information and associations for objects are learned over a lifetime. The challenge, however, is to understand how this knowledge interacts with memory.

Bartlett (1932) showed that memories are guided by schemas that help to fill in the details of memories. For example, providing labels can activate schemas that guide the interpretation of the stimulus and serves as an aid to memory. Carmichael, Hogan, & Walter (1932) showed that providing labels can facilitate and influence later reconstruction. They had subjects study a simple line drawing, e.g., two circles and a line (o-o), along with a label. Subjects who were given the label ‘eyeglasses’ later tended to reconstruct the drawing with a curve rather than a line connecting the circles (o^o), representing the nosepiece on a pair of glasses. Subjects who were given the label ‘dumbbell’ tended to reconstruct a thicker line (o=o) similar to the handle on a dumbbell.

Biases need not be from labels provided by the experimenter, but may arise from internal sources as well. Bartlett showed that the participants themselves bring certain biases to the task. In both temporal and serial reproduction he demonstrated how a person’s cultural and social experiences influence their reconstruction to conform to their idiosyncratic biases. Kalish, Griffiths, and Lewandowsky (2007) formalized Bartlett’s serial reproduction task using iterated learning with Bayesian and

human agents. They showed that Bayesian and human learners revert to their prior when inferring the underlying function of a set of coordinates.

While serial reproduction is about the evolution from iteration to iteration, the approach presented here will focus the retrieval from memory based on a single specific event. Previous work by Huttenlocher and colleagues (Crawford, Huttenlocher, & Engebretson, 2000; Huttenlocher, Hedges, & Duncan, 1991; Huttenlocher, Hedges, & Vevea, 2000) has shown that that prior knowledge exerts strong influences on reconstruction from memory. Huttenlocher, Hedges, & Duncan (1991) presented a Bayesian model of category effects positing that reconstruction from memory is a weighted average of specific memory traces and category information. This weighted average ‘cleans up’ noisy memory traces and prevents large errors in reconstruction.

In this paper, we first present the basic approach of the model presented by Huttenlocher and colleagues and then introduce a series of extensions to this model. We assume that the observer is presented with an object during study and is instructed to retrieve from memory a feature of that object at a later time. In the experiment reported in this paper, we test memory for one-dimensional stimulus values, such as the size of an object. In this context, the goal for the observer is to reconstruct the original size  $\mu$  of an object using noisy samples  $y$  that are retrieved from memory. Bayes’ rule gives us a principled way of combining prior knowledge and evidence from memory:

$$p(\mu | y) \propto p(y | \mu)p(\mu) \quad (1)$$

The posterior probability  $p(\mu | y)$  gives the likely stimulus values  $\mu$  given the noisy memory contents  $y$ . This posterior probability is based on a combination of  $p(\mu)$ , the prior knowledge of the likely sizes of the object and  $p(y|\mu)$ , the likelihood of obtaining evidence  $y$  from memory. This Bayesian approach gives a principled account of how prior knowledge of the world is combined with memory contents to recall information about events.

For example, suppose the feature values of objects are Gaussian distributed,  $\mu \sim N(\mu_0, \sigma_0^2)$ , where  $\mu_0$  and  $\sigma_0^2$  are the prior mean and variance of the feature values. Furthermore, when a specific object value  $\mu_s$  is studied, suppose this leads to samples  $y$  drawn from episodic memory with the samples having a Gaussian noise distribution centered on the original studied value,  $y \sim N(\mu_s, \sigma_m^2)$ . The variance of the noise process,  $\sigma_m^2$ , controls the degree to which the

stored episodic representations resemble the original studied object features. The exact source of the noise is not modeled in this account but this could be related to decay or interference with other events entering memory. Standard Bayesian techniques can now be used to calculate the posterior distribution in Eq.1. The conditional probability of recalled stimulus value  $\mu_r$  given the contents of memory  $y$  is given by a Gaussian distribution with mean  $\mu_n$ ,

$$\mu_n = w\mu_o + (1-w)\bar{y} \quad (2)$$

where  $w=(1/\sigma_o^2)/[(1/\sigma_o^2)+(n/\sigma_m^2)]$  and  $n$  is the number of samples taken from episodic memory. Note that the mean of the recalled stimulus values is a weighted linear combination of the prior mean  $\mu_o$  and the mean of memory content  $\bar{y}$ . The prior mean  $\mu_o$  is weighted more heavily in recall when the prior has a higher precision ( $1/\sigma_o^2$ ) and when the memory noise increases. This corresponds to the intuitive notion that if the prior is strong, it will have a strong influence on recall. Similarly, if memory contents are very noisy, the prior will also exert a strong influence on recall.

This model predicts systematic biases toward the category center, or prior category mean, at reconstruction. Figure 1 illustrates these biases and the effect of the strength of the prior. The small vertical lines represent the small noisy samples (around 0.2) drawn from episodic memory at the time of test. In the left panels, we simulate drawing a single sample ( $n=1$ ) from memory. The dashed lines represent prior knowledge, and the solid lines represent the posterior distribution that forms the basis for recall. Using classical statistical inference (top panel) with an uninformative prior, the posterior is centered on the mean of the memory sample -- there is no effect of the prior. Using Bayesian inference (bottom two panels), we specified a prior with mean  $\mu_o = 0.4$ . We simulated a relative vague prior with precision  $1/\sigma_o^2 = 200$ . Using the Bayesian inference procedure as described above, the resulting posterior is slightly shifted toward the prior. For a relative precise prior with precision  $1/\sigma_o^2 = 2000$ , the result is a posterior that is shifted much more away from the data and toward the prior. The right panels show the results when four samples ( $n=4$ ) are drawn from memory. In this case, the evidence from memory is stronger which decreases the influence of the prior. Subsequently, the posterior distribution is less influenced by the prior.

### Extending the Basic Approach

The approach sketched above formed the basis for the theory by Huttenlocher and colleagues. We propose a hierarchical extension to this theory where prior knowledge can come from multiple sources. We will conduct a behavioral experiment using natural objects such as fruits and vegetables for which participants have pre-experimental knowledge at multiple levels. For example, we expect that

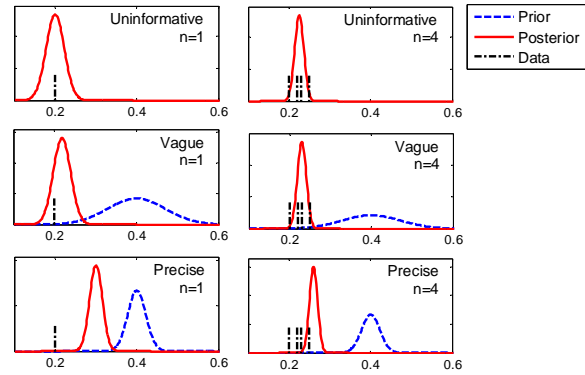


Figure 1. Illustrations of a Bayesian account for the systematic biases in reconstructive memory due to prior knowledge. See text for details.

participants not only have prior knowledge at the category level (e.g. “I expect fruits to be roughly of this size”) but also at the object level (e.g. “I expect an apple to be of this size”). We predict that the influence of the object and category prior knowledge depends on an individual’s familiarity with the object and category. If a participant studies an object with which they are familiar, e.g., a chayote (a type of gourd), then they can use their knowledge about the common size of this object to aid their reconstruction and correct an otherwise noisy memory trace at test. Another participant that studies the same chayote, who does not know this object might be able to recognize it as a vegetable and can use his general knowledge at the category level to guide reconstruction. In the experiment, we will test some of the predictions from this extended theory, and focus on the role of multiple, hierarchical sources of prior knowledge in reconstructive memory.

### Experiment

In the following experiment we first measured the perceived size of common natural objects such as fruits and vegetables as an estimate of participants’ prior knowledge for these objects. In the second phase of the experiment we assessed recall memory for size. We used the observed size ranges from the norming phase as the foundation for the study sizes in order to encourage the use of prior knowledge.

We predict that the effect of prior knowledge at the category and object level will be observed by systematic biases towards the mean of the object and category prior at reconstruction. At the category level, this means that small objects (e.g. raspberries) will be overestimated while large objects (e.g. pineapple) will be underestimated. At the object level, this means that objects presented at relatively small sizes (e.g., a small apple relative to all apples) will be overestimated while large objects (e.g., a large apple) will be underestimated.

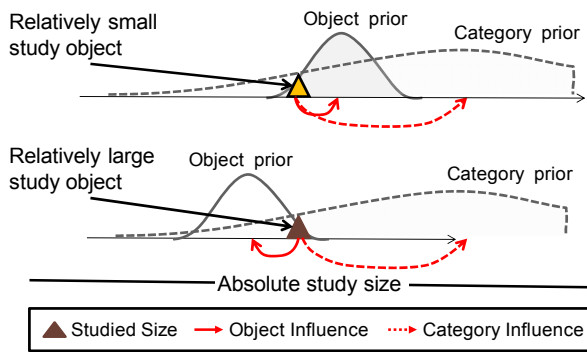


Figure 2. Predicted influences of category and object level priors for two objects studied at the same size.

Separating out the contributions from category and object level is difficult because in many cases, the effects might operate in the same direction. This is illustrated in Figure 2, top panel. If an object is studied at a size that is small relative to both the category and object prior (e.g., a small apple), both of these priors will result in a positive bias. They are both operating in the same direction, toward the category center. The clearest demonstration of independent contributions of object and category level prior knowledge is provided when the effects go in opposite directions. For example, in Figure 2, bottom panel, the object (e.g., a large strawberry) is studied at a size that is large relative to the object prior but small relative to the category prior. In this case, the category effect leads to an overestimation while the object effect will lead to an underestimation of the object at reconstruction. The crucial comparison is between the top and bottom panel. In both cases, the objects are shown at exactly the same size during study. However, we predict that the reconstructed sizes will be different because of the difference in relative object sizes.

## Methods

Participants were undergraduate students at the University of California, Irvine. There were 18 participants in the norming phase and 25 participants in the test phase.

For the norming phase there were 37 images in 2 categories: fruits and vegetables. For the test phase, 24 of the objects from each of the norming categories were used. See Figure 4 for examples. Another class of stimuli was also developed: abstract shapes created by drawing outlines of objects and filling with blue. See Figure 3 for examples.

**Norming Phase.** All materials were presented on two computer screens. A reference object was presented on the



Figure 3. Examples from the shapes category created by drawing outlines of objects filled in blue.

left screen and the object of interest was presented on the right screen. Participants were asked to make three size judgments for each object: “What is the smallest (or average or largest) size of an object like this?” Participants manipulated the size of the object using a slider. Responses were measured on a scale from zero to one. Images were presented in random order and at one of four initial sizes relative to the overall screen size: 0.2, 0.4, 0.6 or 0.8.

**Memory Phase.** The study sizes of the images were sampled from the size ranges collected in the norming phase. For sampling, we used a truncated Gaussian distribution between the min and the max of the individual object range. The objects were never shown outside of the min-max range. The shapes category was yoked to the vegetable category for size and orientation on the screen. The specific study size for each shape was the same as that of its yoked vegetable. Participants were shown a continuous random sequence of study and test images. Each study image was presented a total of three times during the experiment, and there was always a related intervening test trial between presentations. Each participant completed three blocks of 72 study and test images. Study images were presented for two seconds. At test participants were asked to make two memory judgments. They were first asked to make a recognition decision about whether they remembered seeing the object at study. Second they were asked to make a recall judgment about the size of the object at study using the slider on the screen to manipulate the size of the object. Responses were measured on a scale from zero to one.

## Results

**Norming Phase.** Figure 4 depicts the 24 objects from the vegetable category. The top panel indicates the range of the size judgments for individual objects averaged over

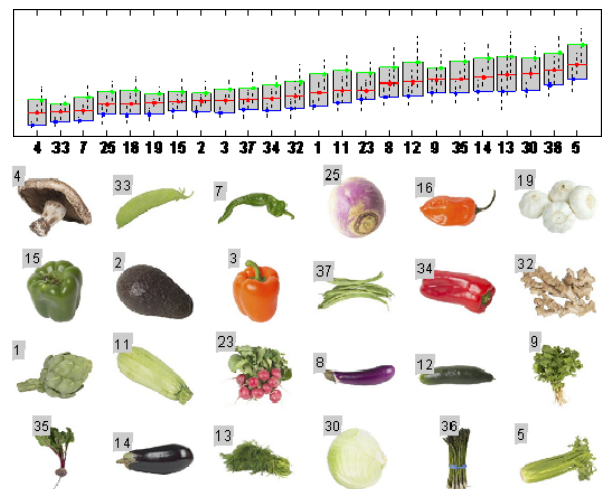


Figure 4. Norming results for the vegetable category. Bars show the range of size judgments. The center horizontal lines show the mean of the ‘average’ judgments.

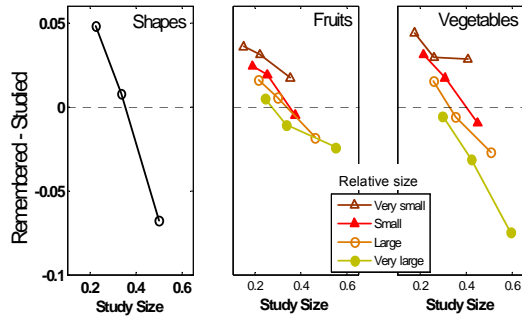


Figure 5. Reconstruction error by category. Positive and negative error indicates over- and underestimation respectively.

participants. The results follow a natural order: the mean ‘average’ size judgment for mushrooms is smaller than for bell peppers, which are all smaller than celery, and so on. Participants expressed a large degree of agreement, although variability does increase with the magnitude of the objects.

**Memory Phase.** Reconstruction error (reconstructed size – studied size) was used to measure performance in each category. Reconstruction error as a function of category and object class is plotted in Figure 5. Positive reconstruction error indicates overestimation while negative reconstruction error indicates underestimation. The observed pattern of correction toward the category center as indicated by negative slopes for all categories supports the prediction of category effects. To assess the influence of object priors, we divided the study objects into four classes based on the study sizes relative to the minimum and maximum acceptable sizes as assessed in the norming experiment. We divided the range between the minimum and maximum in four equal ranges and named those ranges “very small”, “small”, “large”, and “very large”. These classes therefore give the sizes of objects relative to the mean of the object – e.g. a “very large” object might be an apple that is studied at close to the maximum size (relative to all apples). The results show a regular pattern for different object classes (very small < small < large < very large). This difference in intercepts by relative study size supports the prediction of object prior effects. To measure the effects of prior knowledge on reconstruction memory at both the category and object level, a regression model was fitted to each subject assuming a fixed slope and separate intercepts for each relative object size (except for the shapes category where we use one intercept). Average slopes and intercepts are reported in Table 1.

The slope for each category was significantly different from zero (fruits:  $t(24)=-4.714$ ,  $p=0.000$ , vegetables:  $t(24)=-5.657$ ,  $p=0.000$ , shapes:  $t(24)=-10.754$ ,  $p=0.000$ ). This is consistent with a category level influence of prior knowledge. A significant trend was observed within each

Table 1. Average slopes and intercepts by category.

|            | Fruits |      |    | Vegetables |      |    | Shapes |      |    |
|------------|--------|------|----|------------|------|----|--------|------|----|
|            | Mean   | SD   | N  | Mean       | SD   | N  | Mean   | SD   | N  |
| Slope      | -.120  | .127 | 25 | -.191      | .169 | 25 | -.415  | .193 | 25 |
| Intercepts |        |      |    |            |      |    | .142   | .069 | 25 |
| Very small | .057   | .049 | 25 | .089       | .062 | 25 |        |      |    |
| Small      | .048   | .045 | 25 | .075       | .058 | 25 |        |      |    |
| Large      | .040   | .044 | 25 | .065       | .059 | 25 |        |      |    |
| Very large | .035   | .054 | 25 | .047       | .065 | 25 |        |      |    |

category such that intercepts for very small object were larger than those of very large objects (fruits:  $t(24)=2.569$ ,  $p=0.016$ , vegetables:  $t(24)=3.991$ ,  $p=0.001$ ). These differences are consistent with an object level influence of prior knowledge.

## Model

The results showed that natural stimuli such as fruits and vegetables are associated with multiple levels of pre-experimental prior knowledge, each exerting an influence on reconstructive memory. In our first extension of the basic model by Huttenlocher and colleagues, we propose that prior knowledge can be represented at multiple levels of abstraction which can independently influence reconstruction from memory. We propose a simple mixture model where the prior mean and variance ( $\mu_o$ ,  $\sigma_o^2$ ) is a combination of category and object level priors,

$$\mu_o = z\mu_i + (1-z)\mu_c \quad (3)$$

$$\sigma_o^2 = z\sigma_i^2 + (1-z)\sigma_c^2 \quad (4)$$

where ( $\mu_i$ ,  $\sigma_i^2$ ) represents the object prior associated with object  $i$  and ( $\mu_c$ ,  $\sigma_c^2$ ) represents the category prior. The variable  $z$  weights the contribution of the object prior relative to the category prior. We assume that this weighting is determined by

$$z \sim \text{Bernoulli}(\theta_i) \quad (5)$$

where  $\theta_i$  is a constant that represents the familiarity of an object. In this model, familiar objects lead to a prior that is more dependent on the object rather than the category. Similarly, this implements the intuitive notion that for unfamiliar objects, it is unlikely that the object prior is reliable and inference instead reverts to a higher-level prior based on categorical knowledge.

As before, we assume that the computational goal for the participant is to invert the forward memory model and reconstruct the original event given the noisy memory contents and prior knowledge about the study event. The solution to this computational problem was described in Eq. 2.

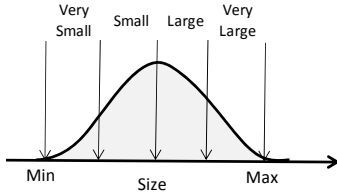


Figure 6. Example of the Gaussian distribution for a simulated object prior. The four regions label the study sizes relative to the object prior.

We applied the model in Eqs. 2-5 to our experimental setup and aimed for qualitative fits to the data as opposed to detailed quantitative fits. We used three values of  $\theta$  corresponding to no object familiarity ( $\theta=0$ ), medium familiarity ( $\theta=.4$ ) and high familiarity ( $\theta=.7$ ). In this model, the category prior can represent a combination of a priori knowledge about the category as well as knowledge accumulated during the experiment (such as the distribution learned for the shapes category). Here, we will not distinguish between these two sources for the category prior and use a single prior with  $\mu_c=0.5$  and  $1/\sigma_c^2=20$ . This is a relatively vague prior that is centered near the mean of study sizes we used in the experiment. For the object priors, we simulated Gaussians with means centered across the range  $[0,1]$  and precision  $1/\sigma_o^2=200$ . This implements a relatively precise object prior compared to the category prior. For the study sizes, we drew samples from the object priors, and rejected samples outside the  $[0,1]$  range. For the purpose of data analysis, we categorized the study sizes into four classes: “very small”, “small”, “large” and “very large”. These size indications are relative to the object prior. Figure 6 illustrates this discretization process. Just as in the experiment, the label “very small” refers to an object that was presented at study at a value close to the minimum size for that particular object. This size is *not* related to the absolute study size. For example, we can simulate a very small pineapple that is still larger than most other fruits. Finally, we ran the simulation with a memory precision of  $1/\sigma_m^2=50$ .

Figure 7A shows the model predictions. Overall, the results show effects of both the category and object prior. Objects that were studied at small sizes with respect to the category and the object prior are overestimated while large study sizes relative to the object and category prior are underestimated. Also, as expected, variations in familiarity can modulate the influence of the object prior. For  $\theta=0$ , there is no influence of the prior. This situation is comparable to the experimental results for the shapes category for which participants did not have any pre-experimental knowledge specific to the object.

When comparing the slopes in this simulation and the experimental results in Figure 5, an important discrepancy arises. In the experimental data, the effect of the category prior is stronger for the shapes compared to the fruits and vegetables (see Table 1 for the difference in estimated slopes across categories). In the simulation in Figure 7A, the

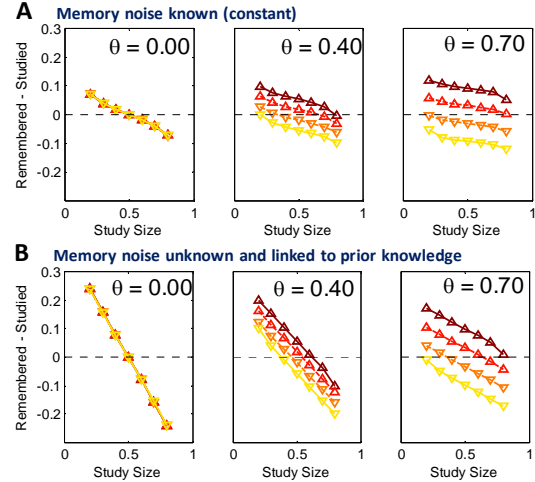


Figure 7: (A) model predictions when memory noise is a constant parameter. (B) model predictions when memory noise is an unknown variable.

effect of the category prior is more or less constant across the levels of familiarity. The difference in category prior cannot be explained due to differences in the study size distributions because the shapes category was yoked to the vegetables category and exactly the same sizes were presented across the two categories.

These experimental results raise an interesting issue about the relative effects of the priors. For objects that have presumably very little prior knowledge, we see relative strong effects of the priors, exactly opposite to what the basic Bayesian approach would predict. This suggests that an additional change to the theory is needed to fully explain the data. We will now describe a change to the noise process that governs the sampling of memory representations. This additional extension will lead to a model that is able to qualitatively describe our findings.

In the basic approach, the memory noise  $\sigma_m^2$  is treated as a constant parameter and the theory does not explain how this parameter is set or varies across experimental conditions. Moreover, this approach assumes that the observer knows the memory noise parameter during the inference process. However, it seems unlikely that the observer has access to such knowledge. We propose that the memory noise is itself an unobserved variable which needs to be estimated from the data (i.e., the memory samples) and prior knowledge. From a statistical point of view, we propose a system where the goal is to make inferences about data with an unknown mean and unknown variance. In the statistics literature, several solutions exist for this problem, and we follow a standard approach (e.g. Gelman et al., 2003) that allows an analytic solution.

As before, we assume that noisy memory samples  $y$  are drawn from episodic memory with a Gaussian noise distribution  $y \sim N(\mu_s, \sigma_m^2)$  that is centered around the original studied value  $\mu_s$ . Instead of assuming a constant noise variance  $\sigma_m^2$ , the noise variance is sampled from an inverse- $\chi^2$  distribution:

$$\sigma_m^2 \sim \text{Inv-}\chi^2(v_0, \sigma_o^2) \quad (6)$$

and the mean of the stimulus values is assumed to be conditionally dependent on the noise variance:

$$\mu_s | \sigma_m^2 \sim N(\mu_o, \sigma_m^2 / \kappa_0) \quad (7)$$

The constants  $v_0$  and  $\kappa_0$  represent the prior degrees of freedom and the prior sample size respectively. The goal for the observer is to calculate the conditional probability of recalling size  $\mu_r$  given the contents  $y$  in memory. This leads to the following solution:

$$\mu_r | y \sim t_{v_0+n}(\mu_n, \sigma_n^2) \quad (8)$$

$$\mu_n = \frac{k_o}{k_o+n} \mu_o + \frac{n}{k_o+n} \bar{y} \quad (9)$$

$$\sigma_n^2 = \frac{v_o \sigma_o^2 + (n-1)s^2 + (k_o n / k_o + n)(\bar{y} - \mu_o)^2}{(v_o + n)(k_o + n)} \quad (10)$$

Note the similarity of Eq. 9 to Eq. 2. In both cases, the mean of the recall distribution is a linear combination of the prior mean and the mean of the observed memory samples. Figure 8 shows a graphical representation of the complete model. Shaded nodes represent observed variables while nodes without shading represent unobserved variables. The arrows indicate the conditional dependencies between the variables.

Note also that memory noise is modeled as an unobserved variable and that there is a coupling between the memory noise variance and the prior variance of the study event. This corresponds to intuitive notions about memory; if we encode objects with which we not very familiar and have little associated prior knowledge, it is more difficult to store accurate representations in memory for that object. In contrast to the previous model where memory noise was left as an unexplained parameter, this model explains memory noise as a variable dependent on prior knowledge. We simulated this model in the same manner as the previous model. We used the same category and object priors and set

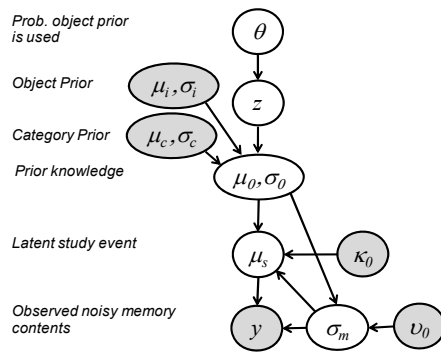


Figure 8. The graphical model representation for the hierarchical Bayesian model for reconstructive memory.

$v_0=5$  and  $\kappa_0=5$ . The results are shown in Figure 7B. Note that the model predicts that the category prior is relatively strong for the low familiarity conditions. This somewhat paradoxical effect falls out of the model because of the coupling between memory and noise and prior knowledge. Objects with weak priors (e.g., shapes) are associated with relatively noisy samples from memory. The result is that the prior exerts a stronger influence to reduce the effects of the memory noise. On the other hand, objects with strong priors (e.g. fruits and vegetables) are associated with relative precise samples from memory leading to a reduced influence of the prior overall.

## Conclusion

We have given a hierarchical Bayesian account of reconstructive memory, where reconstruction of the size of the original study event is influenced by prior knowledge at multiple levels. Unfamiliar objects lead to inferences that are more influenced by the category center, whereas familiar objects lead to inferences that are more influenced by the object prior. A novel assumption of the model is that memory noise is unknown to the observer and becomes part of the inference process. This assumption is different from the basic approach as described by Huttenlocher et al. (1991, 2000) but is consistent with empirical data showing that category effects exert a greater influence when the observer has no pre-experimental knowledge for the object, i.e., the object is unfamiliar. While it seems counter intuitive that a vague prior exerts a stronger influence in reconstructive memory, this is to be expected if we couple the memory noise process to the prior.

## References

- Bartlett, F. C. (1932). *Remembering: A study in experimental and social psychology*. Cambridge, England: Cambridge University Press.
- Carmichael, L. C., Hogan, H. P., & Walter, A. A. (1932). An experimental study of the effect of language on the reproduction of visually perceived form. *Journal of Experimental Psychology*, 15, 73-86.
- Crawford, E., Huttenlocher, J., & Engebretson, P. H. (2000). Category effects on estimates of stimuli: Perception or reconstruction? *Psychological Science*, 11, 280-284.
- Gelman, A., Carlin, J.B., Stern, H.S., and Rubin, D.B. (2003). *Bayesian Data Analysis*. Chapman & Hall.
- Huttenlocher, J., Hedges, L. V., & Duncan, S. (1991). Categories and particulars: Prototype effects in establishing spatial location. *Psychological Review*, 98, 352-376.
- Huttenlocher, J., Hedges, L. V., & Vevea, J. (2000). Why do categories affect stimulus judgments? *Journal of Experimental Psychology: General*, 129, 220-241.
- Kalish, M. L., Griffiths, T. L., & Lewandowsky, S. (2007). Iterated learning: Intergenerational knowledge transmission reveals inductive biases. *Psychonomic Bulletin & Review*, 149, 288-294.