# The Importance of Long-term Memory in Infant Perceptual Categorization

**Martial Mermillod**
**Robert M. French**
Quantitative Psychology and
Cognitive Science
Psychology, U. of Liège,Belgium
{rfrench, mmermillod}@ulg.ac.be

**Paul C. Quinn**
Psychology,
University of Delaware
Newark, DE, USA
pquinn@psych.udel.edu

**Denis Mareschal**
Psychology,
Birkbeck College
London, UK
d.mareschal@bbk.ac.uk

## Abstract

Quinn and Eimas (1998) reported that young infants include non-human animals (i.e., cats, horses, and fish) in their category representation for humans. To account for this surprising result, it was proposed that the representation of humans by infants functions as an attractor for non-human animals and is based on infants' previous experience with humans. We report three simulations that provide a computational basis for this proposal. These simulations show that that a "dual-network" connectionist model that incorporates both bottom-up (i.e., short-term memory) and top-down (i.e., long-term memory) processing is sufficient to account for the empirical results obtained with the infants.

## Introduction

During the last decade, an increasing amount of computational research, in particular, connectionist modeling, has been devoted to the basic mechanisms underlying human categorization (e.g., Anderson & Fincham, 1996; Kruschke, 1992). Our own research has focused on developing a computational model of early infant categorization and testing that model empirically (French, Mermillod, Quinn, & Mareschal, 2001; Mareschal, & French, 1997; Mareschal, French, & Quinn, 2000; Mareschal, Quinn, & French, 2002).

Quinn, Eimas, and Rosenkrantz (1993) observed a surprising categorization asymmetry in young infants between 3 and 4 months of age. After being exposed to a series of photos of cats, the infants showed greater interest in an image of a novel dog compared to a novel cat. However, after exposure to a series of dogs, infants of the same age showed no significantly different interest in either a new dog or a new cat.

We hypothesized that this categorization asymmetry was due to the greater perceptual variability of dogs and to the fact that the ranges of perceptual features of cats were largely included in those of dogs. In short, when familiarized on dogs, a new cat was perceived as something very much like what had already been seen. But, when familiarized on cats, a new dog was generally outside of what the infants had been familiarized on (i.e., cats). The explanation required that very young infant categorization of these animals be essentially a bottom-up process.

For reasons that are given in detail elsewhere (see, Mareschal & French, 1997; Mareschal et al., 2000) we used a three-layer, non-linear autoencoder to model this categorization asymmetry. The model predicted a reversal of this categorization asymmetry when the original variances and inclusion relationship between the two sets of stimuli was reversed by selecting a highly varied set of cats and a set of dogs with low variability. This prediction was subsequently verified experimentally with young infants (French et al., 2001). The model also predicted a disappearance of this categorization asymmetry when the inclusion relationship was removed by careful selection of cat and dog breeds for the stimuli. Again, we were able to empirically verify that the asymmetry did, in fact, disappear (French, Mareschal, Mermillod, & Quinn, 2003). This work strongly supports the view that categorization by young infants of certain types of objects (cats, dogs, horses, cars, etc.) is almost exclusively a bottom-up, statistically driven process with no contribution from prior conceptual knowledge.

### "Perceptual attractors"

Recently, however, Quinn and Eimas (1998) reported a very interesting effect that suggests that this picture has to be modified when *human* perceptual features are involved. The essence of their experiment is as follows. Using an experimental design identical to that used in Quinn et al. (1993), they showed 3- and 4-month-old infants images of a series of pairs of horses, followed by a pair of test images consisting of a novel horse and a human (or a fish or a car). As expected, the infants looked longer at the novel category (humans, fish, and cars) than the new exemplar from the familiarization category (horse). However, when the infants were familiarized on twelve images of humans, and then were presented with an image of a novel human or a horse (or a fish or a car), *there was no significant increase in looking time for the exemplar from the novel animal categories, although there was a significant increase in looking time for the car.*

In other words, infants do not seem to be able to recognize an animal exemplar from a novel category after being familiarized with humans. The result was initially attributed to a lack of power of the experiment. However, two replications were done with a large number of subjects and the effect remained. Further, control experiments show that there is no discrimination bias among the exemplars used in the experiments and no spontaneous preference for exemplars of humans over instances of non-human animals. It was also

suggested that the possibility of broader variance of the human category, combined with overlapping distributions of various perceptual features of the images, might have produced categorization asymmetries similar to those in Quinn et al. (1993) and French et al. (2001). However, these hypotheses did not stand up to closer scrutiny. By testing the typicality of the pictures of humans, horses, and fish by naïve observers, Quinn and Eimas (1998) found that the human category actually seems to be the least variable of the three. (This was subsequently verified by an analysis of the Gabor filtered images.)

To explain the asymmetrical categorization of humans and non-human animals, Quinn and Eimas (1998) suggested that the early exposure of the infants to human visual stimuli might generate a global category that included other animals. The human visual stimuli might act as a powerful "perceptual attractor" for stimuli sharing even a small number of common perceptual attributes with humans.

The simple autoencoder model of Mareschal and French (1997), Mareschal et al. (2001), and French et al. (2001), cannot model the categorization asymmetry observed by Quinn and Eimas (1998). We hypothesize the need for a shift from the purely bottom-up perceptual categorization paradigm of the simple autoencoder to a model that includes a long-term memory capacity that is able to influence purely bottom-up categorization.

The remainder of this paper is organized as follows. We first show that the standard autoencoder model fails on the Quinn and Eimas (1998) data and discuss why this occurred. We then present a dual-network memory system (Ans & Rousset, 1997, 2000; French, 1997) that involves a continual interaction between two networks — one designed to process new input, which we might loosely designate as the STM network, the other designed for long-term storage, which we call the LTM network. We first code the images in terms of neurobiologically plausible spatial-frequencies (Archambault, Gosselin, & Schyns, 2000; French, Mermillod, Mareschal & Quinn, 2002). We then show that, if this dual-network system has previously stored in its LTM network prior perceptual information about human images, and if this information is re-introduced into STM when it is processing new input, the STM network reproduces the asymmetric categorization results of Quinn and Eimas (1998). Finally, we examine the hypothesis that, because the effect could be due to a simple enlargement of the "human image" attractor basin, it might be possible to obtain the same effect by simply enlarging this basin by adding noise to the input when training the network on perceptual images of humans. Our initial investigation of this issue shows that the addition of noise is not sufficient, implying that this process really does require the co-mingling with the new data to be learned by the network of previously-stored information of the perceptual images of humans.

## Contribution of LTM in young infants

Our hypothesis raises the question of early long-term memory storage and consolidation of perceptual stimuli. Previous research has shown evidence of early long-term memory storage under certain circumstances. Rovee-Collier, Evancio, and Earley (1995) found that 3-month-old infants are capable of long-term storage as long as the stimulus is "refreshed" within a certain time window. They reported long-term retention for reinforcement learning if the infants had a reminder within 2 to 3 days after initial learning. In related research, Merriman, Rovee-Collier, and Wilk (1997) showed that long-term retention could influence a categorization task by 3-month-old infants. They showed that infants exposed to the stimuli during 3 daily sessions are capable of some degree of long-term retention. In our study, we assume that young infants are exposed to humans sufficiently often to allow consolidation in long-term memory of this particularly important perceptual stimulus class.

## Dual-network memory systems

Near the end of the 1980's a serious problem with many connectionist models came to light — namely, the problem of catastrophic forgetting (McCloskey & Cohen, 1989), where new learning completely destroys previously learned information. McClelland, McNaughton, and O'Reilly (1995) suggested that the brain's way of avoiding this problem was the development of two complementary learning systems, the hippocampus and the neocortex. New information was learned in the hippocampus and old information was stored out of harm's way in the neocortex. At about the same time, French (1997) and Ans & Rousset (1997, 2000) suggested dual-network connectionist architectures to overcome this problem. These were coupled networks, continually exchanging information by means of *pseudopatterns* (Robins, 1995). One network served as a long-term storage network (LTM); the other (STM) was used to learn new information. When new information was to be learned by the STM, a number of LTM pseudopatterns (each of which reflected the contents of LTM) were produced by sending noise through the LTM network and associating this noise with its output. A series of LTM input/output patterns generated in this way were then mixed with the new patterns to be learned by the STM network. Catastrophic forgetting of previously learned information was thereby effectively overcome.

Results show that for some categories young infants do form (and presumably use) long-term memory traces. This makes particularly appropriate the use of this dual-network architecture to simulate the asymmetric categorization results of Quinn and Eimas (1998). We will see that this type of dual-network connectionist model does, indeed, reproduce these results.

## Simulation of the perceptual system

In order to simulate infant learning, we need a neurobiologically plausible means of encoding perceptual information from the visual environment. We chose an encoding scheme that mimics the neural processes from the retina to the V1 visual pathway when in the presence of an image. This scheme involves decomposing each image into a spatial frequency map (Acerra, Burnod, & de Schonen, 2002; Archambault et al., 2000; French et al., 2002). The neurobiological plausibility of this encoding derives from the fact that different columns in V1 are sensitive to different ranges of spatial frequencies (De Valois & De Valois, 1988; Tootell, Silverman, & De Valois, 1981). We were able to characterize each image as a unique vector of 26 real numbers each of which corresponded to an "energy" value for a Gabor filter, simulating the activity of V1 complex cells (Sakaï & Tanaka, 1999). Space constraints do not allow us to present the details of this encoding here; they can be found in French et al. (2002).

## Overview of the three simulations

We present three simulations. The first shows that the autoencoder model originally used by Mareschal and French (1997) and Mareschal et al. (2000), implementing the infant habituation theories of Sokolov (1963), cannot simulate the empirical data in Quinn and Eimas (1998). In Simulation 2, we show that the dual-network model described above in which the LTM network has stored encodings of visual images of humans does correctly simulate the data in Quinn and Eimas (1998). Finally, in Simulation 3, if the LTM network has no prior learning of human images, the dual-network model does not reproduce the results of Quinn and Eimas (1998).

## Simulation 1: Failure of the autoencoder model to simulate Quinn & Eimas (1998)

As discussed in the Introduction, the bottom-up autoencoder model of Mareschal and French (1997) and Mareschal et al. (2000) has been remarkably successful in simulating (and predicting) categorization performance for certain types of categories in young infants. Autoencoders are connectionist networks that learn to produce on output what is presented on input. The original model used measurements of explicit features to encode the images seen by the infants. This encoding was made more neurobiologically plausible by using spatial-frequency data to characterize the inputs (French et al., 2002). This suggests that, at least for the categories used (dog, cat, fish, horses, etc.), young infant categorization was a bottom-up process driven by the statistical distributions of the perceptual features of the stimuli.

However, Quinn and Eimas (1998) used this procedure with humans and horses (and cats and fish, as well). To reiterate, they found that when familiarized on horses, infants show, as expected, a significantly higher interest thereafter when shown the image of a human compared to that of a novel horse. However, when

exposed first to images of humans, there is subsequently no significantly higher interest in horses (or the other nonhuman animal species)!

In the following "dual-network" simulations, we consider the LTM network to be the "top-down, knowledge-based" network. This addition contrasts with the purely "bottom-up" statistical learning of the patterns in the environment by the simple autoencoder without a LTM network.

### Network

We used a standard 26-20-26 feedforward backpropagation autoencoder network (learning rate: 0.1, momentum: 0.9). We chose a 26-20-26 architecture to resemble, in terms of the input-hidden unit compression, the architecture used in previous simulations on perceptual categorization (French et al., 2002; Mareschal et al., 1997; Mareschal et al., 2000).

### Stimuli

Using a spatial frequency encoding of the stimuli (French et al., 2002), we simulated the visual acuity of the 3- to 4-month-old infants (4 cycles/degree) for the data used in Quinn and Eimas (1998). The vectors were normalized between 0 and 1, filter by filter, across all of the 36 items comprising the stimuli. For each run of the program, the network was trained on 12 stimuli from one category (either Horses or Humans), and then tested on the 6 remaining stimuli from the training category and 6 randomly chosen stimuli from the 18 stimuli in the remaining category.

### Procedure

As in the original simulations (French et al., 2002; Mareschal & French, 1997; Mareschal et al., 2000), the autoencoder was trained on 12 randomly selected stimuli from one of the two categories (each category had 18 stimuli total). The stimuli were presented to the network in pairs (to simulate presenting the infants with pairs of images) for a fixed duration of 250 epochs (corresponding to the 15-second presentation for each pair of images shown to the infants).



Figure 1. Network error produced by the autoencoder after training on Human and Horse categories. Exemplars of the non-training category produce significant increases in error compared to novel exemplars from the training category.

Upon completion of the training phase, the 6 remaining test vectors from the training category were presented to the network, along with the 6 randomly

chosen vectors from the other category. The observed output of the network was compared to the original input in order to give an error value that measured how well the network was able to autoassociate each of the test patterns. All results were averaged over 50 runs.

Results

The autoencoder produced a significant increase in error when trained on images from the Human category and tested on novel humans compared to horse exemplars ($\underline{F}(1, 98) = 337.2$, $\underline{p}<0.001$). When the network was trained on the category Horse, it also produced a significant increase in error ($\underline{F}(1, 98) = 111.74$, $\underline{p}<.001$). (Figure 1).

Discussion

The model's largely symmetric increase in error was not observed by Quinn and Eimas (1998). They found that when familiarized with images of humans, the preference scores for horses and novel humans were not significantly different from chance, whereas when familiarized with images of horses and tested on humans, the preference scores for humans were significantly above chance.

**Simulation 2: LTM storage of human images**

Overview of the simulation

In order to examine the influence of prior learning and storage in LTM of the human category, we used the dual-network memory model proposed by French (1997), consisting of a long-term storage network (LTM network), where previously learned information is stored, and a short-term storage network (STM network), where new information is learned. We first trained the LTM network on 18 exemplars of humans in different postures and positions. Once this was completed, we then compared the categorization performance of the STM network in two situations. In the first, we trained it on 12 images of humans (randomly selected from a second set of 18 and not the same images as those used to train the LTM network) as in Simulation 1. During learning of the human images, the STM network also received input from the LTM network. After the completion of this familiarization phase on human images, the STM network was tested on the 6 remaining images from the human image set and 6 randomly selected images from the horse image set. The network's categorization performance on these two sets of test images was compared.

We hypothesized that the influence of the representations of humans in LTM on learning in the STM network would produce the categorization asymmetry observed in Quinn and Eimas (1998).

Material

The dual-network memory model is composed of two neural networks (Figure 2) called the STM network and LTM network. Although this is not a requirement of the model, in this simulation each of these networks

is a 26-20-26 feedforward backpropagation autoencoder network identical to the one used in Simulation 1. The LTM network was first trained on a set of images of humans. Then the STM network was simultaneously trained on the new stimuli from the environment and pseudopatterns generated by the LTM network. All parameters of the STM network were identical to those of the LTM network (learning rate of 0.1, momentum: 0.9 and a Fahlman offset of 0.1).



*Figure 2*. The dual-network memory model.

Stimuli

The human-image stimuli used to train the LTM network were 18 pictures of different humans in various positions as might be seen by a 3-month-old infant in different situations. Each of these images was uniquely encoded as a 26-element vector, each of whose values represented the energy value of a particular Gabor filter.

Procedure

We first created a Human category representation in the LTM network based on the learning of 18 exemplars of humans taken from real-life settings. Each stimulus was filtered with an average acuity of 2 cycles per degree for the category learned by the LTM network (to simulate the visual acuity of infants before 3 to 4 months of age when, presumably, they would have acquired this category). In the test phase (simulating 3- to 4-month-old infants) this visual acuity was increased to 4 cycles per degree. The number of training epochs the LTM network was set at 1000 in order to create a reasonably reliable representation of the Human category in this network.

We then tested the influence of that LTM representation on category learning in the STM network. Each time a set of patterns (in the present simulation each set contains two patterns) was presented to the STM network, 4 new pseudopatterns were generated by the LTM memory. Feedforward-backpropagation weight changes were then made for patterns to be learned, as well as for each of the four pseudopatterns. For each learning epoch, 4 new LTM pseudopatterns were generated. In this way, a reflection of the contents of LTM is learned by the STM network, along with the new patterns. The maximum number of

training epochs was raised from 250 to 2000 epochs in order to allow the STM network to develop reliable internal representations of the new patterns from the environment combined with the contents from LTM memory. The ratio of pseudopatterns to real patterns is 2:1 in order to ensure the STM network is provided with a relatively good reflection of the contents of LTM.



*Figure 3*. Neural network error produced by the STM autoencoder after training on Humans and Horses with input from the LTM network previously trained on exemplars from the Human category.

Results

The STM network was trained, as in Simulation 1, first on images from the Human category (while also receiving pseudopattern input from the LTM network). It was then tested on novel images from the Human category and images from the Horse category. As in Quinn and Eimas (1998), now there was no significant increase in error for the test exemplars in the Horse category compared to novel exemplars from the Human category ($\underline{F}(1, 98) = .854$, $\underline{p}>0.358$). The STM network was then re-initialized and trained on images from the Horse category (again, while receiving pseudopattern input from the LTM network) and, after this familiarization phase, was tested on novel images of horses and images of humans. In this case there was a significant increase in error for the human images compared to the novel horse images, as in Quinn and Eimas (1998) ($\underline{F}(1, 98) = 86.42$, $\underline{p}<.001$). See Figure 3.

Discussion

These results, using a dual-network model of memory (French, 1997), support the hypothesis that the asymmetric categorization observed in Quinn & Eimas (1998), which we could not simulate with a simple autoencoder, could be due to the influence on STM of a representation of the Human category in LTM.

## Simulation 3: The contents of LTM

Overview of the simulation

Our hypothesis is that the dual-network memory model was able to reproduce the results of Quinn and Eimas (1998) because the LTM network contained a representation of Humans that influenced processing in the STM network. In short, this LTM information was increasing the attractor basin of Humans, causing it to

largely include Horses, thereby giving rise to the asymmetry reported by Quinn and Eimas (1998). However, it might be possible that the contribution of pseudopatterns from the LTM network alone, without this network necessarily having learned anything, could be enough to increase the Human attractor basin, thereby giving rise to the observed categorization asymmetry. This would be equivalent to adding noise to the patterns to be learned by the STM.

To test this we ran the dual-network model *without the LTM network having first learned the Human category*, but with it nonetheless contributing pseudopatterns when the STM network was learning new patterns.

Material and Procedure

The dual-network was identical in all respects to the one run in Simulation 2. The only difference is that the LTM network was left completely untrained. The training and testing procedures were identical to those in Simulation 2.

Results

The results (Figure 4) show that the network returns to the symmetric categorization situation of Simulation 1 in which a simple autoencoder was used. In other words, the content of the LTM network is, indeed, influencing learning in the STM network as it learns new patterns.



*Figure 4*. When the LTM-network is "empty" and generates pseudopatterns that are simply noise, the STM categorization performance returns to the performance of the autoencoder model (see Figure 1).

The autoencoder produced a significant increase in error when trained on the Human category and tested on novel humans compared to horse exemplars ($\underline{F}(1, 98) = 142.87$, $\underline{p}<0.001$). When the network was trained on the Horse category, it also produced a significant increase in error ($\underline{F}(1, 98) = 43.09$, $\underline{p}<.001$).

## Predictions

There are a number of implications of this work on the categorization processes of infants as they grow older and their long-term memory capacity develops. The most important of these is that we should see the disappearance, or at least a significant attenuation, of the purely bottom-up categorization asymmetries observed in Quinn et al. (1993) and French et al. (2001).

It is also perhaps reasonable to assume that there is nothing special about the Human category that was stored in the LTM network in our model. The prediction is that *any* category to which young infants are exposed repeatedly would also serve as an attractor. Presumably, this hypothesis could be tested by artificially exposing young infants repeatedly to a particular category.

## Conclusions

This work represents a first step in the study of the transition from the largely bottom-up processing of category information by very young infants to the categorization mechanisms that are integrated with the long-term memory capacities of the developing infant. Many questions remain about how this change takes place, but we have shown the important contribution of concepts stored in long-term memory to the otherwise largely bottom-up learning of young infants.

## Acknowledgments

## Bibliography

Acerra, F., Burnod, Y., & de Schonen, S. (2002). Modelling aspects of face processing in early infancy, *Developmental Science, 5*, 98–117.

Anderson J.R., & Fincham J.M. (1996). Categorization and sensitivity to correlation. *Journal of experimental psychology:* LMC, *22,* 259-277.

Ans, B., & Rousset, S. (2000). Neural networks with a self-refreshing memory: Knowledge transfer in sequential learning tasks without catastrophic forgetting. *Connection Science, 12,* 1-19.

Ans, B., & Rousset, S. (1997). Avoiding catastrophic forgetting by coupling two reverberating neural networks. *Comptes-Rendus de l'Académie des Sciences,* Série III, 320, 989-997.

Archambault, A., Gosselin, F., & Schyns, P. (2000). A natural bias for the basic level? *Proc. of the 22nd Annual Cognitive Science Conference* (pp. 585-590). NJ: LEA.

De Valois, R.L., & De Valois K.K. (1988). *Spatial vision.* Oxford University Press. New York.

French R. M., Mermillod M., Quinn P. C., & Mareschal D. (2001). Reversing category exclusivities in infant perceptual categorization: simulations and data. *Proc. of the 23th Annual Cog. Sci. Society Conference* (pp. 307-312). NJ: LEA.

French R. M., Mermillod M., Quinn P. C., Chauvin A.., & Mareschal, D. (2002). The importance of starting blurry: Simulating improved basic-level category learning in infants due to weak visual acuity. *Proc. of the 24th Annual Cog. Sci. Society Conference* (pp. 322-327). NJ: LEA.

French, R. M. (1997). Pseudo-recurrent connectionist networks: An approach to the "sensitivity-stability" dilemma. *Connection Science, 9,* 353-379.

French, R. M., Mareschal, D., Mermillod, M., Quinn, P. C. (2003). The role of bottom-up processing in perceptual categorization by 3- to 4-month-old infants: Simulations and data. Manuscript submitted for publication.

Kruschke, J.K. (1992). ALCOVE: An exemplar-based connectionist model of category learning. *Psychological Review, 99,* 22-44.

Mareschal, D., & French, R. (1997). A connectionist account of interference effects in early infant memory and categorization. *Proc. of the 19th Annual Cognitive Science Society Conference* (pp. 484-489). NJ: LEA.

Mareschal, D., French, R., & Quinn, P. (2000). A connectionist account of asymmetric category learning in early infancy. *Developmental Psych., 36,* 635-645.

Mareschal, D., Quinn, P. C., & French, R. M. (2002) Asymmetric interference in 3- to 4-month-olds' sequential category learning. *Cognitive Science, 26,* 377-389.

McClelland J.L., McNaughton B.L., & O'Reilly R.C. (1995). Why there are complementary learning systems in the hippocampus and neocortex: insights from the successes and failures of connectionist models of learning and memory. *Psychological Review, 102,* 419-457.

McCloskey, M., & Cohen, N. (1989)..Catastrophic interference in connectionist networks: The sequential learning problem. In G. Bower (Ed.), *The Psychology of Learning and Motivation*, V.24 (pp. 109-165). NY: Academic Press.

Merriman J., Rovee-Collier C., & Wilk A. (1997). Exemplar spacing and infants' memory for category information, *Infant Behavior and Development, 20,* 219-232.

Quinn, P. C., & Eimas, P. D. (1998). Evidence for a global categorical representation of humans by young infants. *Journal of Experimental Child Psychology, 69,* 151-174.

Quinn, P. C., Eimas, P. D., & Rosenkrantz, S. L. (1993). Evidence for representations of perceptually similar natural categories by 3- and 4-month-old infants. *Perception, 22,* 463-475.

Robins, A.V. (1995) Catastrophic forgetting, rehearsal and pseudorehearsal. *Connection Science*, 7, 123–146.

Rovee-Collier, C., Evancio, S., & Earley, L.A. (1995). The time window hypothesis: Spacing effects. *Infant Behavior and Development, 18 (1),* 69-78.

Sakaï, K., & Tanaka, S. (1999). Spatial pooling in the second-order spatial structure of cortical complex cells. *Vision Research, 40,* 855-871.

Sokolov, E. N. (1963). *Perception and the conditioned reflex.* Hillsdale, NJ: LEA.

Tootell, R. B., Silverman, M. S., & De Valois, R. L. (1981). Spatial frequency columns in primary visual cortex, *Science,* 214, 813-815.

Wilson, H. (1988). Development of spatiotemporal mechanisms in infant vision, *Vision Research, 28,* 611-628.