

Read et al. (2023). *Hippocampus*

**TITLE: Computational models can distinguish the contribution from different mechanisms to familiarity recognition**

John Read<sup>1</sup>, Emma Delhayé<sup>1,2</sup>, & Jacques Sougné<sup>2,3</sup>

<sup>1</sup> GIGA Centre de Recherche du Cyclotron In Vivo Imaging, University of Liège, Liège, Belgium

<sup>2</sup> Psychology and Cognitive Neuroscience Research Unit, University of Liège, Liège, Belgium

<sup>3</sup> UDI-FPLSE, University of Liège, Liège, Belgium

**Corresponding author:** John Read

Email: john.read08@gmail.com

ORCID for the corresponding author: 0009-0009-3269-4362

*This is the peer reviewed version of the following article:*

**“Read, J., Delhayé, E., & Sougné, J. (2023). Computational models can distinguish the contribution from different mechanisms to familiarity recognition. *Hippocampus*, 1–15”,**

*which has been published in final form at <https://doi.org/10.1002/hipo.23588>. This article may be used for non-commercial purposes in accordance with Wiley Terms and Conditions for Use of Self-Archived Versions. This article may not be enhanced, enriched or otherwise transformed into a derivative work, without express permission from Wiley or by statutory rights under applicable legislation. Copyright notices must not be removed, obscured or modified. The article must be linked to Wiley’s version of record on Wiley Online Library and any embedding, framing or otherwise making available the article or pages thereof by third parties from platforms, services and websites other than Wiley Online Library must be prohibited.*

## **ACKNOWLEDGMENTS**

The authors thank Stephan Defraire for his invaluable support throughout the coding process at the early stages of the study as well as Dr. Christine Bastin for the inspiring discussion about the behavioral part of this work. The authors also want to give a very special thanks to Dr. Daniel Defays for his help in the comprehension of the mathematics behind the algorithms, his insightful comments on the manuscript and his support throughout this entire project.

## **ABSTRACT**

Familiarity is the strange feeling of knowing that something has already been seen in our past. Over the past decades, several attempts have been made to model familiarity using artificial neural networks. Recently, two learning algorithms successfully reproduced the functioning of the perirhinal cortex, a key structure involved during familiarity: Hebbian and anti-Hebbian learning. However, performance of these learning rules is very different from one to another thus raising the question of their complementarity. In this work, we designed two distinct computational models that combined Deep Learning and a Hebbian learning rule to reproduce familiarity on natural images, the Hebbian model and the anti-Hebbian model respectively. We compared the performance of both models during different simulations to highlight the inner functioning of both learning rules. We showed that the anti-Hebbian model fits human behavioral data whereas the Hebbian model fails to fit the data under large training set sizes. Besides, we observed that only our Hebbian model is highly sensitive to homogeneity between images. Taken together, we interpreted these results considering the distinction between absolute and relative familiarity. With our framework, we proposed a novel way to distinguish the contribution of these familiarity mechanisms to the overall feeling of familiarity. By viewing them as complementary, our two models allow us to make new testable predictions that could be of interest to shed light on the familiarity phenomenon.

**KEY WORDS:** Neural Networks (Computer), Recognition (Psychology), Perirhinal Cortex, Algorithms

## INTRODUCTION

Recognition memory has been described as the ability to determine if one has already encountered or not an event such as an object or a person (see Besson, Ceccaldi, & Barbeau, 2012 for a review article on the subject). Although it was highly debated among the scientific community, it is now commonly accepted that two retrieval processes can occur during recognition (Jacoby, 1991; Mandler, 1980; Tulving, 1985; Yonelinas, Ramey, & Riddell, 2022). Familiarity-based recognition is the feeling of knowing that something – or someone – has already been seen in the past, without recall of the context in which it has been encountered (Tulving, 1985; Yonelinas, Aly, Wang, & Koen, 2010). By contrast, recollection-based recognition refers to the experience of consciously remembering an event (Tulving, 1985; Yonelinas et al., 2010). Over the past decades, *Dual Process theories* proposed that recollection and familiarity work as two functionally and anatomically independent processes (see Diana, Reder, Arndt, & Park, 2006; Eichenbaum, Yonelinas, & Ranganath, 2007; Yonelinas, 2002 for reviews).

Recent studies suggest that familiarity emerges through the implication of an anterior-temporal network including several brain regions (Bastin et al., 2019; Merkow, Burke, & Kahana, 2015; Ritchey, Libby, & Ranganath, 2015; Scalici, Caltagirone, & Carlesimo, 2017; Yonelinas, Otten, Shaw, & Rugg, 2005). Previous works have also shown that the perirhinal cortex (PrC) is crucial during familiarity detection (Aggleton & Brown, 2006; Bowles et al., 2010; Eichenbaum et al., 2007; Montaldi & Mayes, 2010). For example, a study showed that during a recognition task, patients with specific lesions in the PrC present impaired familiarity without recollection dysfunction (Brandt, Eysenck, Nielsen, & von Oertzen, 2016). These works were also supported by Wolk, Dunfee, Dickerson, Aizenstein, & Dekosky (2011), who showed an anatomic double dissociation between familiarity associated with the PrC and recollection associated with the hippocampus.

Looking more closely at patterns of neural firing during familiarity-based recognition, electrophysiological studies in monkeys showed that a small fraction of PrC neurons – called *novelty neurons* – respond in a stronger manner when new stimuli are presented (Brown & Xiang, 1998; Xiang & Brown, 1998). More importantly, this pattern of high activation tends to decrease when the same stimuli are presented again (Brown & Aggleton, 2001). In other words, when a stimulus is new, novelty neurons in the PrC will respond with a higher firing rate. But, when the same stimulus becomes familiar, its activity in the PrC is reduced compared to a novel stimulus. This phenomenon known as repetition suppression has also been observed in the human brain. This is notably the case in the inferotemporal cortex, a region which is adjacent to the PrC and is involved in visual perception (Grill-Spector, Henson, & Martin, 2006; Meyer & Rust, 2018).

Several works in computational modeling are grounded around Dual Process frameworks and the implication of the PrC in familiarity detection (see Cowell, 2012 for a review). For example, a neurocomputational model brought evidence that human must resort to two complementary learning systems to adequately capture the mechanisms of recognition memory (Norman & O'Reilly, 2003). According to this framework, the hippocampus is involved in the recall of details from specific events (i.e., recollection) whereas the medial-temporal cortices – including the PrC – learned the statistical regularities of the environment (i.e., familiarity). Intriguingly, it seems difficult to implement these two functions in a single system (McClelland, McNaughton, & O'Reilly, 1995). Therefore, Norman & O'Reilly (2003) developed two separate networks for recognition: the hippocampal model for recollection and the neocortical model for familiarity. Basically, the neocortical model (Norman, 2010; Norman & O'Reilly, 2003) encodes regularities in the input layer (i.e. a stimulus) with Hebbian learning and assigns similar representations to similar stimuli. When the same stimulus is presented repeatedly to the neocortical model, the internal representation of this stimulus will sharpen

gradually and fewer neurons will respond to the stimulus. However, these neurons will be strongly activated. Here, familiar stimuli will strongly activate a small number of neurons whereas novel stimuli will weakly activate many neurons (see Sohal & Hasselmo, 2000 for another model using Hebbian plasticity and competition). Paradoxically, the idea behind familiarity-based recognition is the ability to recognize events that only have occurred once (Yonelinas et al., 2022). This assumption seems therefore incompatible with the gradual learning postulated by the neocortical model.

Another major limitation of the neocortical model – as well as for other architectures – is that they used formal binary patterns (i.e., sequences of 0s and 1s) as direct inputs for memorizing (Bogacz & Brown, 2003b; Norman & O’Reilly, 2003; Sohal & Hasselmo, 2000). One could reasonably assume that this kind of inputs are not congruent with the processing occurring in human brain. As a matter of fact, our judgments of familiarity arise from events involving real stimuli instead of artificial patterns. Eichenbaum et al. (2007) proposed a functional organization for visual processing in the median temporal lobes including the PrC. In this organization, most of the neocortical input to the PrC comes from association areas called the ventral pathway (Eichenbaum et al., 2007). The ventral pathway process unimodal sensory information about qualities of objects: the so-called “*what*” information (Humphreys & Riddoch, 2006). The representation of a stimulus formed by the ventral pathway allows subsequent judgment of familiarity.

Trying to fulfill the gap between modelling and human brain processing, some models used convolutional neural networks (CNN) to mimic the ventral pathway processing to the PrC (Ji-An, Stefanini, Benna, & Fusi, 2022; Kazanovich & Borisyuk, 2021; Tyulmankov, Yang, & Abbott, 2022). CNN consist of the succession of several layers of artificial neurons programmed to reproduce the detection of features in a stimulus, mimicking the cells in the visual cortex (Bengio & Lecun, 1997). In this type of networks, a gradient backpropagation

algorithm adjusts weights so that the neurons of each layer manage to detect particular patterns such as vertical and horizontal bars. For example, for the first convolutional layer which could be compared to visual layer V1 in the brain, neurons will detect a particular orientation line such as a vertical bar regardless of where that line appears in the image. In principle, every pre-trained CNN followed by a simple neural network could successfully model familiarity on natural images with an adequate synaptic plasticity rule (Kazanovich & Borisyuk, 2021).

Accordingly, two synaptic plasticity rules seem very promising to model familiarity-based recognition: the Hebbian and the anti-Hebbian learning rules (Bogacz & Brown, 2003b; Tyulmankov et al., 2022). The functioning of these learning rules are based on Hebb's works (Hebb, 1949):

*“When an axon of cell A is near enough to excite cell B and repeatedly or persistently takes part in firing it, some growth process or metabolic change takes place in one or both cells such that A's efficiency, as one of the cells firing B, is increased.”*

Computational models using Hebb's theory to model familiarity are essentially designed as two-layers feedforward networks (Androulidakis, Lulham, Bogacz, & Brown, 2008; Bogacz & Brown, 2003a; Bogacz, Brown, & Giraud-Carrier, 2001; Sohal & Hasselmo, 2000). In these networks, weight modification is implemented such as the connection strengths are either strengthened or weakened in response to co-activated neurons. The direction of this modification (i.e., strengthening or weakening) depends on the chosen synaptic plasticity rule. Respectively, the Hebbian plasticity potentiates connection strengths while the anti-Hebbian plasticity depresses them in response to a stimulus.

The advantage of these learning rules is that they are built to reproduce patterns of activity observed in the PrC during familiarity, which correspond to physiological evidence (Brown & Aggleton, 2001; Brown & Xiang, 1998). In that way, models of that kind provide a biologically plausible implementation for familiarity recognition (Bogacz & Brown, 2003a;

Bogacz et al., 2001; Tyulmankov et al., 2022). Nevertheless, Hebbian and anti-Hebbian trainings seem to have distinct properties and thus operate differently from each other. For example, Bogacz & Brown (2003b) observed differences in performance whether inputs are correlated or not. More importantly, some authors argue that Hebbian learning is more biologically plausible than its anti-Hebbian counterpart. According to these authors, this is due to the fact that anti-Hebbian learning tries to reproduce synaptic mechanisms that they declare were not observed in the brain yet (Bogacz & Brown, 2003b). However, this lack of biological plausibility is still debated. In fact, recent works with meta-learning algorithms seems to be in favor of anti-Hebbian models. In the model proposed by Tyulmankov et al. (2022), the network learns from itself (i.e., meta-learns) which one the two learning rules, that is the Hebbian or the anti-Hebbian, should be preferred during training. They showed that a network with meta-learning optimization is more likely to converge to the anti-Hebbian solution. Moreover, anti-Hebbian plasticity seems to generalize better and has a larger memory capacity than Hebbian plasticity (Tyulmankov et al., 2022). So, the question remains: which one of these learning rules should be preferred when one is trying to model familiarity using artificial neural networks?

The goal of this article is to understand, by means of computational models, the inner functioning of the Hebbian and anti-Hebbian training. By comparing how they operate, we want to explain differences in models' abilities on natural images. Therefore, we built two models respectively with Hebbian and anti-Hebbian type of learning rules. The models are preceded by a pre-trained CNN to extract features of images. In this article, we compared the two models' performance under two criteria. First, we implemented and administered Standing's behavioral experiment to the Hebbian and anti-Hebbian models as this experiment is frequently used to test recognition models' performance. Standing's apparatus showed that familiarity has an almost unlimited capacity during a forced-choice recognition (FCR) task (Standing, 1973); results that we were able to replicate only with our anti-Hebbian model, as in Kazanovich &

Read et al. (2023). *Hippocampus*

Borisyuk (2021). Secondly, we compared our models' performance during a FCR task by varying the degree of homogeneity inside the dataset. As expected, we showed poor performance of our Hebbian model for highly similar stimuli. In the discussion, we propose an interpretation of these results based on the idea both models exhibit similarities with two distinct familiarity mechanisms regarding their respective pattern of performance.

## METHODS

Model architectures and recognition task simulations were implemented with the Python 3.9.11 software (<https://www.python.org/>, RRID:SCR\_008394). The code is available in open access on GitHub (<https://github.com/JRead98/master.git>, RRID:SCR\_002630). Note that for our modelling, we used basic model of artificial neurons and not spiking neurons.

### 2.1. Model's architecture

As our model was inspired by previous works, it therefore functions in a similar way (Ji-An et al., 2022; Kazanovich & Borisyuk, 2021). That is, it was designed as a two-step network combining deep learning and simple feedforward neural networks (see **Figure 1**). The goal of this architecture is to reproduce patterns of activity observed in the PrC leading to a familiarity decision during a recognition task.

**[Insert Figure 1]**

Training operates in two stages. First, an image is presented to a pre-trained CNN – in this case ResNet50 – for feature extraction. This mimics the processing in the ventral pathway from visual associative areas to the PrC (Eichenbaum et al., 2007; Le Cun, 2019). This is the feature extraction module. Second, the output of the second-last layer of the CNN is used as an input for a memory module. The memory module is a simple two-layers feedforward network which will learn the features of an image thanks to an Hebbian or an anti-Hebbian learning rule

(similar two-layers networks were also used in Androulidakis et al., 2008). The output of the memory module is used for familiarity discrimination during the testing phase.

## 2.2. Feature extraction module

We used a CNN called ResNet50 as our feature extraction module (He, Zhang, Ren, & Sun, 2015, 2016). More precisely, we used the version ResNet50 v1.5 which was previously trained on PyTorch with 1.2 million high-resolution photographs of natural images from ImageNet (Deng et al., 2009). ResNet50 was initialized as described in He et al. (2015). Originally, ResNet50 allows the classification of images in 1000 different categories with a high rate of accuracy. ResNet50 is built with 48 convolutional layers and 2 pooling layers to identify an image and define its characteristics according to different degrees of complexity. The penultimate layer of the model is a fully connected layer of 2048 features. We use this layer which corresponds to the embedding of the many successive convolutional layers to represent the characteristics of an image. Note that in the complete architecture of ResNet50, the fully connected layer projects onto a SoftMax layer. This SoftMax allows the network to classify images. We do not use this last layer in our architecture.

Before going into the extraction module, the RGB representation of each image was normalized to the size  $3 \times 224 \times 224$  (as in Kazanovich & Borisyuk, 2021); 3 being the number of channels corresponding to the RGB colors and  $224 \times 224$  the size of the images. For a given image, we retrieved a vector of 2048 features obtained at the penultimate fully connected layer of ResNet50. We consider this vector to represent the characteristics of this image. This vector is further used for image learning in the memory module. After passing through the CNN, the vector of size 2048 for a given image is collected then normalized. That is, the distribution of vector values has a mean of zero and a standard deviation of 1. We used this vector of real numbers as inputs for the memory module.

### 2.3. Memory modules

To reproduce familiarity decision, we implemented versions of the memory module that are similar to the version designed by Kazanovich & Borisyuk (2021). In contrast to Bogacz & Brown (2003b), we used simple neural networks instead of a rate-based network.

Both Hebbian and anti-Hebbian modules are two-layers fully connected feedforward networks. Input layers consist of  $n = 2048$  neurons and output layers consist of  $m = 2048$  novelty neurons. Connection strengths (i.e., weights between inputs and outputs) are denoted  $w_{ij}$  and were initialized randomly between -1 and 1. The two learning rules differ in terms of weight modifications (**Figure 2**). Nevertheless, the formula to compute the activity in the output layer is the same for the Hebbian and anti-Hebbian model. That is, we used a forward propagation to compute the activity  $h_j$  of novelty neurons  $j$  according to the following formula:

$$h_j = \sum_{i=1}^n w_{ij} x_i, \quad j = 1, \dots, n \quad (1)$$

where  $x_i$  is the vector of neurons activity for an image  $X$  after normalization in the feature extraction module and  $w_{ij}$  denotes the connection strengths.

Authors originally introduced the notion of active neurons as neurons whose number in the output layer must be limited (Bogacz et al., 2001; Bogacz & Brown, 2003a). We decided to reproduce this distinction between active neurons and neurons at rest using competition and inhibition as previously done in Androulidakis et al. (2008). More precisely, we used an  $m/2$ -winners rule, meaning that half of the novelty neurons  $m$  with the highest activity  $h_j$  are selected to be active (red circles in Figure 2A). The other half are considered inactive and should not participate in the weight modification during the training phase (blue circles). We used the median of the overall activity to determine which neuron is active ( $>$  median) or inactive ( $\leq$  median). Active neurons took the value  $y_j = 1$  and inactive neurons took the value  $y_j = 0$  (see **Figure 2A**).

**[Insert Figure 2]**

### 2.3.1. Hebbian learning rule

In the Hebbian learning rule, we assumed that the novelty neuron  $j$  is active only if the corresponding input neuron  $i$  is also active (Bogacz & Brown, 2003b). Consequently, at the first presentation of an image  $X$ , the activity pattern of novelty neurons  $j$  ( $y$ ) is equal to the activity pattern of input neurons  $i$  ( $x$ ). Thus, in vector form, we consider that the initial response of the networks would be:

$$x^X = y^X$$

where  $x^X$  is the vector of neurons activity for an image  $X$  after normalization in the feature extraction module and where  $y^X$  is the vector of novelty neurons activity for an image  $X$ . In the Hebbian model, we didn't use the activity of novelty neurons during training given this assumption that the initial response of the network is equal to the activity of input neurons. Instead, we started by applying the *m/2-winners* rule on the vector  $x^X$  to obtain the vector  $y^X$  constituted of 0s and 1s. We then applied the following weight modification formula (one should note that weights are not bounded and could thus be subject to saturation):

$$w_{ij} = w_{ij} + \eta y_j x_i \tag{2}$$

where  $\eta = 0.01$  is the learning rate (this value has been found as the global minimum in Kazanovich & Borisyuk, 2021),  $x_i$  corresponds to the input neurons after normalization and  $y_j$  corresponds to the input neurons after inhibition and competition. If  $y_j$  and  $x_i$  represent features of the input, then the learning rule will amplify the  $w_{ij}$  link between features that appear together. Here, learning occurs through the increase in connections strengths between co-occurring features as if by Long-Term Potentiation (Bliss & Collingridge, 1993; Bogacz et al., 2001; Sohal & Hasselmo, 2000). This weight modification is implemented once for each image

of the training set. It will lead to an overall higher activity in the output layer when a familiar stimulus is presented again.

However, to correctly mimic the pattern of neuronal firing in the PrC during the presentation of a familiar stimulus, the activity of novelty neurons should be lower for familiar stimuli than novel ones (Brown & Aggleton, 2001; Brown & Xiang, 1998). To overcome this problem, the Hebbian model originally described by Bogacz et al. (2001) used an inhibitory interneuron to model the familiarity decision in the PrC. This inhibitory interneuron is computed from the activity of novelty neurons. It will represent the level of inhibition that should be used to reduce the activity of novelty neurons when a familiar stimulus is presented again. They argued that familiarity decision in their model could be implemented with two methods (Bogacz et al., 2001; Bogacz et Brown, 2003b). First, with the level of inhibition itself, which should be higher for familiar stimuli than for novel ones. Second, with the reduced activity of the novelty neurons after inhibition by the inhibitory interneuron. If the second method allows to reproduce the pattern of neural firing observed in the PrC, it however requires an additional step. Hence, we decided to implement the first option during the testing phase as the authors stated that both methods are viable (Bogacz & Brown, 2003b). We used the activity of the output layer to compute an inhibition level  $d(X)$  as:

$$d(X) = \sum_{j=1}^m x_j h_j \quad (3)$$

where  $h_j$  are the components from the vector of the novelty neurons computed with formula (1) and  $x_j$  are the components from the vector of inputs neurons after normalization. Familiar images should present a higher level of inhibition compared to novel images. Thus, during a recognition task where a pair of images ( $X, Z$ ) is presented to the model, where  $X$  is an old and  $Z$  is a novel image, a correct familiarity decision is made by the model if  $d(X) > d(Z)$ . This can be easily seen by presenting the same image several times to the model during training.

This will amplify the  $w_{ij}$  links between active features, increasing  $h_j$  and consequently increasing  $d(X)$  compared to a novel image  $d(Z)$ .

### 2.3.2. *Anti-Hebbian learning rule*

In the anti-Hebbian learning rule, on the opposite of the Hebbian learning rule, we started by computing the activity  $h_j$  of the novelty neurons with formula (1) before applying the weight modification formula. Thus, there was a diffusion of activity before the weight were modified. Once the output layer is computed, we applied the *m/2-winners* rule on the components  $h_j$  to obtain  $y_j$ .

Here, learning occurs through the decrease in connections strengths between input neurons and active novelty neurons as if by Long-Term Depression (Androulidakis et al., 2008; Bogacz & Brown, 2003a; Ito, 1989). Therefore, weights are modified during training with the following formula:

$$w_{ij} = w_{ij} - \eta x_i y_j \tag{4}$$

where  $\eta = 0.01$  denotes the learning rate,  $x_i$  corresponds to the input neurons after normalization and  $y_j$  corresponds to the components of the vector of novelty neurons after the *m/2-winners* rule. This weight modification will slightly reduce the variance inside the vector of novelty neurons  $h_j$  when computed again with formula (1). As in the Hebbian solution, this weight modification is only implemented a single time for each image of the training set. Nevertheless, the variance reduction can be easily objectified if we repeatedly present a sole stimulus to the anti-Hebbian model. Indeed, after several presentations, the differences between values of novelty neurons for a given image will gradually decrease.

After each image has been studied by the model, we fix the connection strengths before the testing phase. Overall, we should observe an average activity in the output layer that is lower when a familiar stimulus is presented compared to a novel one. During the testing phase, we

computed the average output activity to model the familiarity decision with the following formula (Kazanovich & Borisyuk, 2021):

$$d(X) = \frac{1}{m} (\sum_{j \in M1} h_j - \sum_{j \in M2} h_j) \quad (5)$$

where  $h_j$  are the components from the vector of the novelty neurons computed with formula (1) and  $M1$  and  $M2$  are respectively the sets of  $m/2$ -winners (active neurons) and -losers (inactive neurons) in the output layer. Familiar images should produce lower activity than novel images (Bogacz & Brown, 2003a). Indeed, if we present several times the same image to the network, the  $w_{ij}$  links will decrease, reducing  $h_j$  and consequently decreasing  $d(X)$ . Thus, during a recognition task where a pair of images ( $X$ ,  $Z$ ) is presented to the model, where  $X$  is an old and  $Z$  is a novel image, a correct familiarity decision is made if  $d(X) < d(Z)$ .

#### 2.4. Simulation methodology

The simulation methodology is depicted in **Figure 3** and was similar to that of Kazanovich & Borisyuk (2021). The methodological pipeline is identical for every simulation with a training phase followed by a testing phase. During the training phase, a model was trained on a subsample constituted of  $N$  images randomly taken from the corresponding dataset. Images were learned one-by-one with the weight modification specific to the selected memory module. Each image was presented once to the model for learning. In the testing phase, we implemented a forced-choice recognition (FCR) task. During the FCR task,  $N$  pairs of images were presented simultaneously to the network: a new image as well as an image previously learned during training. The model had thus to decide which image is familiar depending on the memory module. If the model has chosen the new image as familiar, a recognition error was logged.

**[Insert Figure 3]**

## RESULTS

All the plots were generated by using ggplot2 package (Wickham, 2009, <https://cran.r-project.org/web/packages/ggplot2/index.html>, RRID:SCR\_014601). The data obtained during the different simulations and the script used to visualize them are openly available on the OSF platform (<https://osf.io/vpgdm/>, RRID:SCR\_003238).

As a first simulation, we reproduced Standing’s experiment to evaluate the memory capacity of the models with the methodological pipeline depicted in **Figure 3** (Standing, 1973). The dataset consisted of a database of about 30 000 natural images divided into 256 object categories (Caltech 256 Image Dataset; Griffin, Holub, & Perona, 2007). All categories contained on average 119 images and a minimum of 80 images.

As part of the simulation, we estimated the error probability ( $P_{err}$ ) for the entire task then averaged it on 100 runs of the models. Each run was realized with a different training and testing set. We also computed the number of images retained in memory, similarly to Kazanovich & Borisyuk (2021):

$$N_{ret} = N(1 - 2P_{err}) \quad (6)$$

where  $N$  is the number of images presented during training and  $P_{err}$  is the error probability for the entire task. Results from the first simulation are shown in **Figure 4**.

**[Insert Figure 4]**

As expected, we observed for both our models that performance decreases gradually as the dataset size increases (**Figure 4B**). That is, the error probability is on average worse when the models are tested with large datasets than with small datasets (**Table 1**). In the medium dataset size condition ( $N = 100$ ), both models still have good accuracy. However, when this threshold is crossed, the performance of the Hebbian model started to decrease more drastically than its anti-Hebbian counterpart.

In comparison with human data, we can see that anti-Hebbian model outperformed human performance until 1000 images are presented. As a matter of fact, it is only for the two biggest datasets ( $N = 4000$  and  $N = 10000$ ) that the anti-Hebbian model performs worse than human. One should note that the performance still reaches more than 65% accuracy with the highest dataset size, suggesting that the model didn't perform at chance level even in this condition. Regarding the Hebbian model, the probability of error is similar to human behavioral performance up to 40 images. Passed this dataset size, performance of the Hebbian model gradually decreased to reach random choices between familiar and novel images for the highest dataset size ( $P_{err} = 0.5$ ). This random choice pattern of answers tends to come up when more than 1000 images were presented during the training phase.

Moreover, we observed that the memory capacity for the anti-Hebbian solution is strikingly similar to human performance with on average  $\mu = 3-171.760$  ( $\sigma = 101.402$ ) images retained in memory for  $N = 10\ 000$ . In fact, it managed to have near perfect memory for most of the dataset sizes (**Table 1**). Overall, it tends to fit the power law observed in Standing's original experiment (**Figure 4A**). In comparison, the Hebbian solution seems to have a poor memory capacity which didn't exceed 376 images when 10 000 pictures are learned during training ( $\mu = 150.760$ ;  $\sigma = 90.912$ ). On average, the number of images retained in memory by the Hebbian model seems to be constant for every dataset size that exceeds a hundred pictures.

**[Insert Table 1]**

Next, we wanted to check whether the models could display a recency effect. To highlight such an effect, we estimated the probability that a network will make an error for a given pair of images during the testing phase and averaged it over 100 runs of the models. We performed the simulations at the threshold where models' performance started to diverge while they both kept more than 80% accuracy ( $N = 100$ ). Graphically, a recency effect should be marked by a gradual decrease in the average error probability as a function of the image position

in the training phase. Results were then smoothed with a Loess Regression function and plotted in **Figure 5**.

**[Insert Figure 5]**

Interestingly, it seems that the anti-Hebbian model exhibits a recency effect that is not observed with the Hebbian model. The former has indeed lower probabilities of error for images learned at the end of the training (i.e., recent images) compared to images learned at the beginning of the training. This is not the case for the Hebbian model which showed no tendency to make less mistakes for recent images.

For the second simulation, we decided to test whether the models are sensitive to homogeneity between the inputs. We tested the performance of the models in three conditions of homogeneity: heterogeneity, mild homogeneity, high homogeneity. The heterogeneity condition consisted of random pictures selected from the Caltech 256 database (Griffin et al., 2007). The two homogeneous conditions consisted of two datasets, each constituted with only one semantic category of images, respectively dogs and cats. The mild homogeneity condition thus corresponded exclusively to dogs' pictures randomly selected for the dog's category folder from the Caltech 256 database (Griffin et al., 2007). Regarding the high homogeneity condition, we used exclusively cats' pictures randomly selected from the so-called "Cat Dataset", which consists of nearly 10 000 pictures of cats divided in 7 sub-folders (W. Zhang, Sun, & Tang, 2008). We justify our choices on the fact the dogs have a wider variety of perceptual features than cats (i.e. dogs are more heterogeneous than cats, French, Quinn, & Mareschal, 2001; Mareschal, French, & Quinn, 2000).

Simulations took place similarly as in the first simulation (see **Figure 3**). The only difference is that for the mild and high homogeneity conditions, models were trained exclusively with dogs or cats' pictures, respectively. For example, in the high homogeneity condition, models had to learn 40 images of cats ( $N = 40$ ). During the testing phase,  $N$  pairs of

cats' pictures were presented to the model: a new and an old cat. The models had to decide which one was familiar.

As previously done, the results were average over 100 runs of the models. Each run was realized with a different training and testing set. The average  $P_{err}$  and standard deviations for the three homogeneity conditions are plotted in **Figure 6**.

**[Insert Figure 6]**

Foremost, the anti-Hebbian model has a better accuracy than the Hebbian model in every homogeneity condition. With the anti-Hebbian learning rule, model performance still reaches high accuracy in the high homogeneity condition. Performance is furthermore stable for the heterogeneity to the mild homogeneity, and we observed no decrease in accuracy between the two conditions. In fact, the anti-Hebbian model has a near perfect accuracy when trained and tested with low and no homogeneity between the inputs. With the Hebbian learning rule, we observed a gradual decrease as the homogeneity between the pictures increases during the learning phase. Moreover, we can see that when the Hebbian model is trained with cat pictures only (i.e., high homogeneity), the model responds randomly during the FCR task.

**[Insert Table 2]**

A summary of our key results is detailed in **Table 2**. For each simulation, we estimated the error probability ( $P_{err}$ ) for the entire task then averaged it on 100 runs of the models. Each run was realized with a different training and testing set. We can observe that both the Hebbian and anti-Hebbian model have more than 80% accuracy on small, medium, and mildly homogeneous datasets. Besides, the accuracy is numerically higher in the anti-Hebbian model in every conditions. Regarding the performance of the Hebbian model on large and highly homogeneous dataset, it seems that the model failed 1 out of 2 times to correctly choose the familiar image. We interpreted these results as random answers.

## DISCUSSION

The goal of the paper is to compare two learning rules which can be used to model familiarity by reproducing the pattern of neural firing observed in the PrC. Here, by differentiating Hebbian and anti-Hebbian learning on natural images, we want to provide insight into the operations at hands when a stimulus becomes familiar. We replicated previous results showing that the anti-Hebbian solution has on average a higher memory capacity than the Hebbian solution (Bogacz & Brown, 2003b; Kazanovich & Borisyuk, 2021). Besides, the former fits relatively well Standing's behavioral data (Standing, 1973) whereas the later only fits the data when the training set doesn't exceed 40 images. Regarding their ability to manage homogeneity between the inputs, we showed that the anti-Hebbian model once again has better accuracy than its Hebbian counterpart as previously suggested by the work of Androulidakis et al. (2008). In fact, the anti-Hebbian model still reaches high accuracy even with highly homogeneous stimuli (i.e., cats). The Hebbian model reaches more than 80% accuracy in the mild homogeneity condition (i.e., dogs). Nevertheless, it fails to perform above chance in the high homogeneity condition suggesting high vulnerability to homogeneity.

On one hand, our results with the anti-Hebbian model are in line with previous networks using anti-Hebbian learning to model familiarity (Androulidakis et al., 2008; Kazanovich & Borisyuk, 2021). Interestingly, in the model proposed by Kazanovich & Borisyuk (2021), they did not implement inhibition and competition *per se*. Rather, they only applied formula (5) to withdraw the activity from the sets of losers (i.e., half of the neurons in the output layer with the lowest activity) for the pair of images presented during the FRC task. Besides, they used AlexNet for features extraction instead of ResNet50 as in our modeling (Krizhevsky, Sutskever, & Hinton, 2012). Despite these slightly different implementations of the anti-Hebbian model, we still managed to reproduce their results on Standing's experiment.

In addition, our results showed that the anti-Hebbian model can react to the more recent (i.e., familiar) stimuli with greater accuracy. More importantly, by reducing the overall activity in the output layer, it successfully reproduces the repetition suppression mechanisms observed in the brain when a stimulus becomes familiar (Grill-Spector et al., 2006; Meyer & Rust, 2018). According to Tanaka, Saito, Fukada, & Moriya (1991), repetition suppression is thought to be very selective for complex visual stimuli. In fact, it provides the specific information that would permit recognizing a recent stimulus. Taken together, this suggests that it is the anti-Hebbian learning rule ability to reduce the variance of the vector of novelty neurons that allows it to accurately model familiarity recognition (Bogacz & Brown, 2003a). If the target has lower variance in its output layer than the lure, it should mean that the target has more recency – or familiarity – than the lure. Our simulations showed that this ability is impaired neither by the number of presented stimuli nor by the similarity between targets lures.

On the other hand, to our knowledge, this is the first time that a Hebbian learning rule was implemented on natural images instead of artificial inputs. This makes the comparison with other networks difficult. Nevertheless, Bogacz & Brown (2003b) have previously shown that its performance should be lower than an anti-Hebbian model when there were dependances between the stimuli features. To address this issue, Kazanovich & Borisyuk (2021) have computed this dependances for the images from the Caltech 256 database. As expected, they showed that the co-occurrence between pairs of features could be high for some pictures. It is then plausible that differences in models' performance to reproduce behavioral data could be explained to some extent by co-occurrence between the features of an image. This is also in line with our results showing high sensitivity to homogeneity between inputs in the Hebbian model only. However, this raises the question: can the Hebbian solution provide an accurate modeling framework for familiarity recognition in human?

Based on the results of our simulations, we can reasonably admit that the Hebbian model can successfully discriminate between old and new pictures under certain conditions (small to medium dataset set, mild homogeneity in the training data). We also showed that our version of the Hebbian learning rules operates by encoding co-occurrence between features that appeared together in an image. This means that the learning rule will increase the connection strength between two active features of an input. For example, consider a picture of an old man as the stimulus. He has glasses, a beard, and a baldness that are considered as active features. The Hebbian model will increase the link between the glasses and the beard, between the baldness and the beard, and so on. In other words, the Hebbian model will create a global representation of a stimulus. This means, regarding recognition, that stimuli where glasses appear with beard and where baldness appears with glasses will be more familiar to the system.

Interestingly, this description of our Hebbian model is consistent with the global matching models (GMM) of recognition (Clark & Gronlund, 1996; Osth & Dennis, 2020). The assumption behind GMM is that an item is constituted of several memory representations (i.e., several features). During a recognition task, a cued item will activate these representations. The activation of these components of an item will be combined to produce global match. If the match signal is high enough, it will lead to a familiarity judgment. More importantly, GMM predicts that high number of stimuli and similarity between stimuli (i.e., homogeneity) will both lead to impaired recognition judgment (Brandt, Zaiser, & Schnuerch, 2019; Cary, 2003). Along with the results from our simulations, this suggest that the Hebbian model could indeed correspond to a mechanism for familiarity recognition.

It has long been thought that familiarity could involve different co-existing mechanisms (Bastin et al., 2019; Mandler, 1980; Mecklinger & Bader, 2020). Therefore, the anti-Hebbian and Hebbian model should not be mutually exclusive. Instead, we believe that our models are quite complementary and can provide insight into answering questions of that sort. In a review

article, Mecklinger & Bader (2020) highlight the distinction between a relative familiarity and an absolute or lifetime familiarity. The former would be associated with stimuli that have been recently encountered whereas the latter would be linked to stimuli that have been frequently encountered in our lifetime. These two familiarity mechanisms have been dissociated in ERP studies (which showed two distinct early familiarity signals, the FN400 vs. N400 components, see Bridger, Bader, & Mecklinger (2014) for an example), in neuropsychological studies in patients with PrC lesions (Anderson, Baena, Yang, & Köhler, 2021) and in functional MRI studies (Duke, Martin, Bowles, McRae, & Köhler, 2017; Yang, McRae, & Köhler, 2023). Specifically, taken together, results from fMRI studies show that both types of familiarity are automatically and jointly elicited, regardless of the task at-hand, in association with the activity of the PrC (Yang et al., 2023). Besides, the PrC exhibits different patterns of activity in association with both relative and absolute familiarity (Duke et al., 2017), with a decrease in the BOLD signal strength that was associated with relative familiarity, consistent with the overall decrease observed in the activity of the novelty neurons layer in our anti-Hebbian model, and an increase in the signal strength associated with absolute familiarity which would be in turn consistent with the overall increase in the activity of novelty neurons in our Hebbian model.

In the same vein, Xiang & Brown (1998) identified three types of neurons that responded differentially during a recognition task in different regions of the anterior temporal lobe – including the PrC – in macaque monkeys, namely the familiarity neurons, recency neurons and novelty neurons. Familiarity neurons signaled “the lifetime familiarity of the stimulus but not whether it had been seen recently (i.e., relative familiarity)”, recency neurons signaled “the recency of presentation but not the lifetime/absolute familiarity of stimuli”, whereas novelty neurons signaled “both the absolute familiarity of a stimulus and whether it had been seen recently”. Although it is currently not known if two different plasticity mechanisms could underlie these familiarity mechanisms or if only one plasticity can account

for both familiarities, it seems plausible that the response pattern observed in novelty neurons of the PrC might be generated through a combination of two plastic mechanisms corresponding to what is observed in our models, consistent with a wide range of studies (Bridger et al., 2014; Duke et al., 2017; Tyulmankov et al., 2022; Xiang & Brown, 1998; Yang et al., 2023). Respectively, one similar to that of recency neurons which would track relative familiarity with anti-Hebbian plasticity and the other similar to that of familiarity neurons which would track absolute familiarity with Hebbian plasticity.

For now, we don't know precisely how these two familiarity processes are articulated together. Coane, Balota, Dolan, & Jacoby (2011) tried to answer this question by clarifying the time course of the familiarity signal. Previous works showed that items already have a baseline familiarity whose level depends on how often an item has been encountered during the lifespan (Joordens & Hockley, 2000; Reder et al., 2000). Coane et al. (2011) hypothesized that when an item is studied, it acquires a temporary increase in its familiarity signal in addition to a permanent increase in its absolute level of familiarity (**Figure 7A**). This temporary familiarity boost corresponds to the relative level of familiarity. Unfortunately, this framework does not tell us about the conditions for a mechanism to take precedence over another.

Yet, a limitation of this framework is that, from an evolutionary point of view, it seems unlikely that the PrC employs both Hebbian and anti-Hebbian models to produce familiarity on different time scales. In fact, a model operating with different time-constants and forgetting rates for different neurons might produce the same outcomes within a single learning rule instead. Since the anti-Hebbian model is more accurate than the Hebbian model, a network employing only an anti-Hebbian learning rule on spiking neurons could be more accurate than a hybrid Hebbian/anti-Hebbian network. Still, it is worth mentioning that Tyulmankov et al. (2022) already designed a hybrid model that takes into account both plastic mechanisms. Despite their discussion mostly focus on the anti-Hebbian solution, they did not exclude the

Hebbian solution *per se* as an alternative plasticity mechanism. Thus, more research is needed to disentangle these two working hypotheses. The advantage of the current framework is that it indeed allows to make new testable predictions by separating the contribution of absolute and relative familiarity to the phenomenological feeling of familiarity (**Figure 7**).

Here, we assume that the anti-Hebbian learning rule models exclusively the relative familiarity through the recency of a stimulus (blue curve on **Figure 7**) and that the Hebbian learning rule models exclusively absolute familiarity through the overall structure of the stimulus, like in GMM (green curve on **Figure 7**). At first, both mechanisms could participate in familiarity decisions when the distinctiveness between stimuli is high (**Figure 7B**). However, when the number of stimuli learned increases, absolute familiarity alone would not be efficient anymore as shown in **Figure 7C** (Brandt et al., 2019; Cary, 2003), because stimuli become too homogeneous. In turn, more homogeneity – meaning less distinctiveness between the stimuli – increases the response criterion necessary to make familiarity judgments. In these conditions, we could only rely on relative familiarity mechanism to maintain recognition accuracy, as shown by the blue line.

According to Mandler (1980), an important feature regarding the two familiarity mechanisms is that relative familiarity should always be greater for items with low lifetime/absolute familiarity, allowing the word frequency mirror effect observed in recognition memory. The word frequency mirror effect suggests that low-frequency words produce higher hit rates and lower false alarm rates than high-frequency words due to their higher level of absolute familiarity (Joordens & Hockley, 2000; Reder et al., 2000). Therefore, when studied, the latter would receive greater familiarity boost than the former, which are insensitive to a single additional exposure (Coane et al., 2011, see **Figure 7**). Although this frequency mirror effect has mostly been studied with words, few studies showed it holds for other variables such as pictures (see Glanzer & Adams (1985) for a review). The frequency mirror effect on pictures

could thus be another testable hypothesis of our modeling framework. Along with the number of items and the similarity between items, the frequency of items could also be a factor of saturation for our Hebbian model so that, when confronted to highly frequent items, the Hebbian model would not be able to perform correctly. In order to test this hypothesis, it would be interesting to vary the baseline/absolute familiarity of pictures from the training set. To do so, one could perform simulation implementing more than one weight modification for a random subset of pictures by presenting them several times during the training phase. The idea is that the more a picture is seen during training, the more its absolute level of familiarity will be similar to high frequency items in the frequency mirror effect. According to **Figure 7**, we can speculate that performance of the anti-Hebbian would not be impacted by the frequency of items whereas the Hebbian model would performed poorly on items that have been seen several times during the training and therefore have higher degree of absolute familiarity.

Interestingly, our hypotheses could also be tested against neuropsychological evidence from the literature on patients with impaired relative familiarity (Anderson et al., 2021; Köhler & Martin, 2020). In particular, one can think about patient NB who exhibits impaired relative familiarity following left PrC damage (Bowles et al., 2010), especially for discriminating between concepts with high semantic features overlap (i.e., high homogeneity, which would saturate the Hebbian model) compared to limited features overlap (Bowles, Duke, Rosenbaum, McRae, & Köhler, 2016). Similarly, patients at risk of Alzheimer's disease, presenting with PrC lesions due to the disease's neuropathology (Braak, Thal, Ghebremedhin, & Del Tredici, 2011), show a selective relative familiarity impairment, with preserved absolute familiarity (Anderson et al., 2021). Therefore, our framework allows us to make additional testable predictions with the aim of disentangling the respective contributions of relative and absolute familiarity in Alzheimer's patients. One could easily imagine a recognition task administered in different homogeneity conditions such as in our second simulation (see Delhaye, Folville, &

Bastin, 2019, for an example of paradigm). The prediction would be that Alzheimer's patients should exhibit impaired relative familiarity responses regardless of the homogeneity condition (see Frick, Besson, Salmon, & Delhaye, 2023). In contrast, they should exhibit preserved absolute familiarity responses only in low homogeneity conditions.

As expected, our study has several limitations that should be acknowledged. To begin with, we wanted to highlight the influence of the CNN on our results. Indeed, it is plausible that the reason why our Hebbian model is highly sensitive to homogeneity is due to the way features are extracted by ResNet50. ResNet50 – as most of other CNN – was trained to categorize pictures (He et al., 2016; Krizhevsky et al., 2012); i.e., this picture is a dog, this picture is a cat. That is what a CNN is trained to do but its inner mechanism is still a black box. Thus, in our models, the vector of 2048 features extracted from the second-last layer of the network has been designed during the training to represent the concept of cat. If we present a picture of another cat to the CNN, the new vector of features could be highly similar to the last picture categorized as a cat by ResNet50. Put in other words, it is plausible that our CNN does not extract the feature of the image *per se* but rather the features of the concept of “cat” itself. This would explain why it is more difficult to choose correctly between two similar images as in our second simulation. However, the lack of similarity effect on the ability of the anti-Hebbian model lets us think that our CNN does not impact that much the results from our simulations.

Another limitation is linked to our reproduction of the initial Hebbian model of Bogacz et al. (2001b). Indeed, in their original paper, they used an inhibitory interneuron to reduce the activity in the output layer after a stimulus is presented for the first time and thus to adequately reproduce the functioning of the PrC. By ease of computation, we decided not to implement the additional step of downsizing the activity of novelty neurons (Bogacz & Brown, 2003b). Rather, we directly used the so-called level of inhibition - which should be used to reproduce repetition suppression – in our decision function. One could therefore say that our model is

incomplete in comparison to the model of Bogacz et al. (2001b). It would therefore be interesting to enhance our Hebbian model to see if our arbitrary simplification could have a profound impact on its performance. Furthermore, it seems apparent that both models are too simplistic to account for the whole properties of real neurons. For example, artificial neural networks as used in our works don't even consider the temporal dimension of synaptic plasticity (L. I. Zhang, Tao, Holt, Harris, & Poo, 1998). One way to overcome this problem would be to implement more realistic neural networks taking into account electro-physiological properties of biological neurons. Yet, it is worth mentioning that this has been done in other computational models of recognition without changing the general pattern of results obtained (Bogacz & Brown, 2003b; Ji-An et al., 2022; Norman & O'Reilly, 2003).

Finally, it is well known that other cerebral structures are involved in familiarity (Bastin et al., 2019). Unlike previous works, we did not take into account the contribution of other brain regions to model recognition (see for example Hasselmo & Wyble (1997) who implemented the contribution of the hippocampal formation to recognition). In fact, our modeling framework considers exclusively the ventral pathway to simulate familiarity given that previous computational models of vision suggested that the PrC - along with the ventral pathway - enables complex visual discrimination behaviors (Bonnen, Yamins, & Wagner, 2021). However, it is known that the PrC is a polymodal association area that receives inputs from many brain regions (Suzuki, 1996), thus enabling familiarity recognition in other sensory domains. We thus cannot account for the entire familiarity phenomenon with ventral pathway inputs only. In addition, the PrC may also act as a conceptual representation site (Clarke & Tyler, 2015; Wright, Randall, Clarke, & Tyler, 2015). More precisely, defining conceptual features of stimuli – such as their category – are represented in the anterior temporal area and the integration of meaning to stimuli representations will occur in the PrC thanks to projections with these regions (see Bastin et al., 2019). Closely related, it has been shown that the

anterolateral entorhinal cortex is also associated with familiarity recognition on images with overlapping features (Besson, Simon, Salmon, & Bastin, 2020). Further work is needed to enhance the biological plausibility of our models. The integration in our modelling framework of other brain regions such as parts of the transentorhinal cortex or even other modalities such as words instead of images could be a promising way to capture more adequately the functioning of familiarity mechanisms (Bastin & Delhaye, 2023; Besson et al., 2020).

## CONCLUSION

In conclusion, we designed two computational models of familiarity in the PrC, the anti-Hebbian and the Hebbian models. We argued that these models should be viewed as complementary as they account for two distinct familiarity mechanisms, respectively relative and absolute familiarity. On one hand, the anti-Hebbian model reduces the variance inside the output layer to compute the recency of an item, which would be a suitable mechanism for relative familiarity. On the other hand, the Hebbian model increases the link between co-occurring features to produce a global match between features activation and a cued item, which would in turn be related to absolute familiarity. We also hypothesized that the contributions of these familiarity processes to recognition can be dissociated when there is not enough distinctiveness between items. To extent this framework, we could challenge predictions made by the models with experimental studies on real subjects.

## REFERENCES

- Aggleton, J. P., & Brown, M. W. (2006). Interleaving brain systems for episodic and recognition memory. *Trends in Cognitive Sciences*, *10*(10), 455–463.  
<https://doi.org/10.1016/j.tics.2006.08.003>
- Anderson, N. D., Baena, E., Yang, H., & Köhler, S. (2021). Deficits in recent but not lifetime familiarity in amnesic mild cognitive impairment. *Neuropsychologia*, *151*, 107735.  
<https://doi.org/10.1016/j.neuropsychologia.2020.107735>
- Androulidakis, Z., Lulham, A., Bogacz, R., & Brown, M. W. (2008). Computational models can replicate the capacity of human recognition memory. *Network: Computation in Neural Systems*, *19*(3), 161–182. <https://doi.org/10.1080/09548980802412638>
- Bastin, C., Besson, G., Simon, J., Delhaye, E., Geurten, M., Willems, S., & Salmon, E. (2019). An Integrative Memory model of recollection and familiarity to understand memory deficits. *Behavioral and Brain Sciences*. <https://doi.org/10.1017/S0140525X19000621>
- Bastin, C., & Delhaye, E. (2023). Targeting the function of the transentorhinal cortex to identify early cognitive markers of Alzheimer’s disease. *Cognitive, Affective, & Behavioral Neuroscience*.  
<https://doi.org/10.3758/s13415-023-01093-5>
- Bengio, Y., & Lecun, Y. (1997). *Convolutional Networks for Images, Speech, and Time-Series*.
- Besson, G., Ceccaldi, M., & Barbeau, E. J. (2012). L’évaluation des processus de la mémoire de reconnaissance. *Rev Neuropsychol*, *4*(4), 242–254. <https://doi.org/10.1684/nrp.2012.0238>
- Besson, G., Simon, J., Salmon, E., & Bastin, C. (2020). Familiarity for entities as a sensitive marker of antero-lateral entorhinal atrophy in amnesic mild cognitive impairment. *Cortex*, *128*, 61–72.  
<https://doi.org/10.1016/j.cortex.2020.02.022>
- Bliss, T. V. P., & Collingridge, G. L. (1993). A synaptic model of memory: Long-term potentiation in the hippocampus. *Nature*, *361*(6407), 31–39. <https://doi.org/10.1038/361031a0>
- Bogacz, R., & Brown, M. W. (2003a). An anti-Hebbian model of familiarity discrimination in the perirhinal cortex. *Neurocomputing*, *52*, 1–6. [https://doi.org/10.1016/s0925-2312\(02\)00738-5](https://doi.org/10.1016/s0925-2312(02)00738-5)

Read et al. (2023). *Hippocampus*

Bogacz, R., & Brown, M. W. (2003b). Comparison of computational models of familiarity discrimination in the perirhinal cortex. *Hippocampus*, *13*(4), 494–524.

<https://doi.org/10.1002/hipo.10093>

Bogacz, R., Brown, M. W., & Giraud-Carrier, C. (1999). High Capacity Neural Networks for Familiarity Discrimination. *Proceedings of the 9th International Conference on Artificial Neural Networks*, 773–778. Edinburgh, UK. <https://doi.org/10.1023/A:1008925909305>

Bogacz, R., Brown, M. W., & Giraud-Carrier, C. (2001). Model of Familiarity Discrimination in the Perirhinal Cortex. *Journal of Computational Neuroscience*, *10*, 5–23.

Bonnen, T., Yamins, D. L. K., & Wagner, A. D. (2021). When the ventral visual stream is not enough: A deep learning account of medial temporal lobe involvement in perception. *Neuron*, *109*(17), 2755–2766.e6. <https://doi.org/10.1016/j.neuron.2021.06.018>

Bowles, B., Crupi, C., Pigott, S., Parrent, A., Wiebe, S., Janzen, L., & Köhler, S. (2010). Double dissociation of selective recollection and familiarity impairments following two different surgical treatments for temporal-lobe epilepsy. *Neuropsychologia*, *48*(9), 2640–2647. <https://doi.org/10.1016/J.NEUROPSYCHOLOGIA.2010.05.010>

Bowles, B., Duke, D., Rosenbaum, R. S., McRae, K., & Köhler, S. (2016). Impaired assessment of cumulative lifetime familiarity for object concepts after left anterior temporal-lobe resection that includes perirhinal cortex but spares the hippocampus. *Neuropsychologia*, *90*, 170–179. <https://doi.org/10.1016/j.neuropsychologia.2016.06.035>

Braak, H., Thal, D. R., Ghebremedhin, E., & Del Tredici, K. (2011). Stages of the Pathologic Process in Alzheimer Disease: Age Categories From 1 to 100 Years. *Journal of Neuropathology & Experimental Neurology*, *70*(11), 960–969. <https://doi.org/10.1097/NEN.0b013e318232a379>

Brandt, K. R., Eysenck, M. W., Nielsen, M. K., & von Oertzen, T. J. (2016). Selective lesion to the entorhinal cortex leads to an impairment in familiarity but not recollection. *Brain and Cognition*, *104*, 82–92. <https://doi.org/10.1016/J.BANDC.2016.02.005>

Read et al. (2023). *Hippocampus*

Brandt, M., Zaiser, A.-K., & Schnuerch, M. (2019). Homogeneity of item material boosts the list length effect in recognition memory: A global matching perspective. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 45(5), 834–850.

<https://doi.org/10.1037/xlm0000594>

Bridger, E. K., Bader, R., & Mecklinger, A. (2014). More ways than one: ERPs reveal multiple familiarity signals in the word frequency mirror effect. *Neuropsychologia*, 57, 179–190.

<https://doi.org/10.1016/j.neuropsychologia.2014.03.007>

Brown, M. W., & Aggleton, J. P. (2001). Recognition memory: What are the roles of the perirhinal cortex and hippocampus? *Nature Reviews Neuroscience*, 2(1), 51–61.

<https://doi.org/10.1038/35049064>

Brown, M. W., & Xiang, J. Z. (1998). Recognition memory: Neuronal substrates of the judgement of prior occurrence. *Progress in Neurobiology*, 55(2), 149–189. [https://doi.org/10.1016/S0301-0082\(98\)00002-1](https://doi.org/10.1016/S0301-0082(98)00002-1)

Cary, M. (2003). A dual-process account of the list-length and strength-based mirror effects in recognition. *Journal of Memory and Language*, 49(2), 231–248.

[https://doi.org/10.1016/S0749-596X\(03\)00061-5](https://doi.org/10.1016/S0749-596X(03)00061-5)

Clark, S. E., & Gronlund, S. D. (1996). Global matching models of recognition memory: How the models match the data. *Psychonomic Bulletin & Review*, 3(1), 37–60.

<https://doi.org/10.3758/BF03210740>

Clarke, A., & Tyler, L. K. (2015). Understanding What We See: How We Derive Meaning From Vision.

*Trends in Cognitive Sciences*, 19(11), 677–687. <https://doi.org/10.1016/j.tics.2015.08.008>

Coane, J. H., Balota, D. A., Dolan, P. O., & Jacoby, L. L. (2011). Not all sources of familiarity are created equal: The case of word frequency and repetition in episodic recognition. *Memory & Cognition*, 39(5), 791–805.

<https://doi.org/10.3758/s13421-010-0069-5>

Cowell, R. A. (2012). Computational models of perirhinal cortex function. *Hippocampus*, 22(10), 1952–1964. <https://doi.org/10.1002/hipo.22064>

Read et al. (2023). *Hippocampus*

Delhaye, E., Folville, A., & Bastin, C. (2019). How to induce an age-related benefit of semantic

relatedness in associative memory: It's all in the design. *Psychology and Aging, 34*(4), 572–

586. <https://doi.org/10.1037/pag0000360>

Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., & Fei-Fei, L. (2009). ImageNet: A Large-Scale Hierarchical

Image Database. *CVPR09*.

Diana, R. A., Reder, L. M., Arndt, J., & Park, H. (2006). Models of recognition: A review of arguments

in favor of a dual-process account. *Psychonomic Bulletin & Review, 13*(1), 1–21.

<https://doi.org/10.3758/BF03193807>

Duke, D., Martin, C. B., Bowles, B., McRae, K., & Köhler, S. (2017). Perirhinal cortex tracks degree of

recent as well as cumulative lifetime experience with object concepts. *Cortex, 89*, 61–70.

<https://doi.org/10.1016/j.cortex.2017.01.015>

Eichenbaum, H., Yonelinas, A. P., & Ranganath, C. (2007). The Medial Temporal Lobe and Recognition

Memory. *Annual Review of Neuroscience, 30*(1), 123–152.

<https://doi.org/10.1146/annurev.neuro.30.051606.094328>

French, R. M., Quinn, P. C., & Mareschal, D. (2001). Reversing Category Exclusivities in Infant

Perceptual Categorization: Simulations and Data. *Proceedings of the Annual Meeting of the*

*Cognitive Science Society, (23)*, 23.

Frick, A., Besson, G., Salmon, E., & Delhaye, E. (2023). Perirhinal cortex is associated with fine-grained

discrimination of conceptually confusable objects in Alzheimer's disease. *Neurobiology of*

*Aging, 130*, 1–11. <https://doi.org/10.1016/j.neurobiolaging.2023.06.003>

Glanzer, M., & Adams, J. K. (1985). The mirror effect in recognition memory. *Memory & Cognition,*

*13*(1), 8–20. <https://doi.org/10.3758/BF03198438>

Griffin, G., Holub, A. D., & Perona, P. (2007). *Caltech 256. Object category dataset: Caltech Technical*

*Report*. <https://doi.org/www.kaggle.com/datasets/jessicali9530/caltech256>

Read et al. (2023). *Hippocampus*

Grill-Spector, K., Henson, R., & Martin, A. (2006). Repetition and the brain: Neural models of stimulus-specific effects. *Trends in Cognitive Sciences*, 10(1), 14–23.

<https://doi.org/10.1016/J.TICS.2005.11.006>

Hasselmo, M. E., & Wyble, B. P. (1997). Free recall and recognition in a network model of the hippocampus: Simulating effects of scopolamine on human memory function. *Behavioural Brain Research*, 89(1–2), 1–34. [https://doi.org/10.1016/S0166-4328\(97\)00048-X](https://doi.org/10.1016/S0166-4328(97)00048-X)

He, K., Zhang, X., Ren, S., & Sun, J. (2015). Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. *Proceedings of the IEEE International Conference on Computer Vision*, 1026–1034.

He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 770–778.

Hebb, D. O. (1949). *The Organization of Behavior*. New York: Wiley.

Humphreys, G. W., & Riddoch, M. J. (2006). Features, objects, action: The cognitive neuropsychology of visual object processing, 1984–2004. *Cognitive Neuropsychology*, 23(1), 156–183.

<https://doi.org/10.1080/02643290542000030>

Ito, M. (1989). Long-Term Depression. *Annual Review of Neuroscience*, 12(1), 85–102.

<https://doi.org/10.1146/annurev.ne.12.030189.000505>

Jacoby, L. L. (1991). A process dissociation framework: Separating automatic from intentional uses of memory. *Journal of Memory and Language*, 30(5), 513–541. [https://doi.org/10.1016/0749-596X\(91\)90025-F](https://doi.org/10.1016/0749-596X(91)90025-F)

Ji-An, L., Stefanini, F., Benna, M. K., & Fusi, S. (2022). Face familiarity detection with complex synapses. *bioRxiv*, 854059. <https://doi.org/10.1101/854059>

Joordens, S., & Hockley, W. E. (2000). Recollection and familiarity through the looking glass: When old does not mirror new. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 26(6), 1534–1555. <https://doi.org/10.1037/0278-7393.26.6.1534>

Read et al. (2023). *Hippocampus*

Kazanovich, Y., & Borisyuk, R. (2021). A computational model of familiarity detection for natural pictures, abstract images, and random patterns: Combination of deep learning and anti-Hebbian training. *Neural Networks, 143*, 628–637.

<https://doi.org/10.1016/j.neunet.2021.07.022>

Köhler, S., & Martin, C. B. (2020). Familiarity impairments after anterior temporal-lobe resection with hippocampal sparing: Lessons learned from case NB. *Neuropsychologia, 138*, 107339.

<https://doi.org/10.1016/j.neuropsychologia.2020.107339>

Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). ImageNet Classification with Deep Convolutional Neural Networks. In F. Pereira, C. J. Burges, L. Bottou, & K. Q. Weinberger (Eds.), *Advances in Neural Information Processing Systems* (Vol. 25, pp. 1097–1105). Curran Associates, Inc.

Le Cun, Y. (2019). *Quand la machine apprend: La révolution des neurones artificiels et de l'apprentissage profond*. Odile Jacob.

Mandler, G. (1980). Recognizing: The Judgment of Previous Occurrence. *Psychological Review, 87*(3), 252–271. <https://doi.org/10.1037/0033-295X.87.3.252>

Mareschal, D., French, R. M., & Quinn, P. C. (2000). A connectionist account of asymmetric category learning in early infancy. *Developmental Psychology, 36*(5), 635–645.

<https://doi.org/10.1037/0012-1649.36.5.635>

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*(3), 419–457.

<https://doi.org/10.1037/0033-295X.102.3.419>

Mecklinger, A., & Bader, R. (2020). From fluency to recognition decisions: A broader view of familiarity-based remembering. *Neuropsychologia, 146*, 107527.

<https://doi.org/10.1016/j.neuropsychologia.2020.107527>

Read et al. (2023). *Hippocampus*

Merkow, M. N., Burke, J. F., & Kahana, M. J. (2015). The human hippocampus contributes to both the recollection and familiarity components of recognition memory. *Proceedings of the National Academy of Sciences*, *112*(46), 14378–14383. <https://doi.org/10.1073/pnas.1513145112>

Meyer, T., & Rust, N. C. (2018). *Single-exposure visual memory judgments are reflected in inferotemporal cortex*. <https://doi.org/10.7554/eLife.32259.001>

Montaldi, D., & Mayes, A. R. (2010). The role of recollection and familiarity in the functional differentiation of the medial temporal lobes. *Hippocampus*, *20*(11), 1291–1314. <https://doi.org/10.1002/hipo.20853>

Norman, K. A. (2010). How hippocampus and cortex contribute to recognition memory: Revisiting the complementary learning systems model. *Hippocampus*, *20*(11), 1217–1227. <https://doi.org/10.1002/hipo.20855>

Norman, K. A., & O'Reilly, R. C. (2003). Modeling Hippocampal and Neocortical Contributions to Recognition Memory: A Complementary-Learning-Systems Approach. *Psychological Review*, *110*(4), 611–646. <https://doi.org/10.1037/0033-295X.110.4.611>

Osth, A. F., & Dennis, S. (2020). *Global matching models of recognition memory* [Preprint]. PsyArXiv. <https://doi.org/10.31234/osf.io/mja6c>

Reder, L. M., Nhouyvanisvong, A., Schunn, C. D., Ayers, M. S., Angstadt, P., & Hiraki, K. (2000). A mechanistic account of the mirror effect for word frequency: A computational model of remember–know judgments in a continuous recognition paradigm. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, *26*(2), 294–320. <https://doi.org/10.1037/0278-7393.26.2.294>

Ritchey, M., Libby, L. A., & Ranganath, C. (2015). Cortico-hippocampal systems involved in memory and cognition. In *Progress in Brain Research* (Vol. 219, pp. 45–64). Elsevier. <https://doi.org/10.1016/bs.pbr.2015.04.001>

Read et al. (2023). *Hippocampus*

Scalici, F., Caltagirone, C., & Carlesimo, G. A. (2017). The contribution of different prefrontal cortex regions to recollection and familiarity: A review of fMRI data. *Neuroscience & Biobehavioral Reviews*, *83*, 240–251. <https://doi.org/10.1016/J.NEUBIOREV.2017.10.017>

Sohal, V. S., & Hasselmo, M. E. (2000). A model for experience-dependent changes in the responses of inferotemporal neurons. *Network: Computation in Neural Systems*, *11*(3), 169–190. [https://doi.org/10.1088/0954-898X\\_11\\_3\\_301](https://doi.org/10.1088/0954-898X_11_3_301)

Standing, L. (1973). Learning 10,000 pictures. *Quarterly Journal of Experimental Psychology*, *25*, 207–222. <https://doi.org/10.1080/14640747308400340>

Suzuki, W. A. (1996). The anatomy, physiology and functions of the perirhinal cortex. *Current Opinion in Neurobiology*, *6*(2), 179–186. [https://doi.org/10.1016/S0959-4388\(96\)80071-7](https://doi.org/10.1016/S0959-4388(96)80071-7)

Tanaka, K., Saito, H., Fukada, Y., & Moriya, M. (1991). Coding visual images of objects in the inferotemporal cortex of the macaque monkey. *Journal of Neurophysiology*, *66*(1), 170–189. <https://doi.org/10.1152/jn.1991.66.1.170>

Tulving, E. (1985). Memory and consciousness. *Canadian Psychology/Psychologie Canadienne*, *26*(1), 1. <https://doi.org/10.1037/h0080017>

Tyulmankov, D., Yang, G. R., & Abbott, L. F. (2022). Meta-learning synaptic plasticity and memory addressing for continual familiarity detection. *Neuron*, *110*(3), 544-557.e8. <https://doi.org/10.1016/j.neuron.2021.11.009>

Wickham, H. (2009). *ggplot2: Elegant Graphics for Data Analysis*. New York, NY: Springer New York. <https://doi.org/10.1007/978-0-387-98141-3>

Wolk, D. A., Dunfee, K. L., Dickerson, B. C., Aizenstein, H. J., & Dekosky, S. T. (2011). A Medial Temporal Lobe Division of Labor: Insights From Memory in Aging and Early Alzheimer Disease. *Hippocampus*, *21*(5), 461–466. <https://doi.org/10.1002/hipo.20779>

Wright, P., Randall, B., Clarke, A., & Tyler, L. K. (2015). The perirhinal cortex and conceptual processing: Effects of feature-based statistics following damage to the anterior temporal

Read et al. (2023). *Hippocampus*

lobes. *Neuropsychologia*, 76, 192–207.

<https://doi.org/10.1016/j.neuropsychologia.2015.01.041>

Xiang, J. Z., & Brown, M. W. (1998). Differential neuronal encoding of novelty, familiarity and recency in regions of the anterior temporal lobe. *Neuropharmacology*, 37(4–5), 657–676.

[https://doi.org/10.1016/S0028-3908\(98\)00030-6](https://doi.org/10.1016/S0028-3908(98)00030-6)

Yang, H., McRae, K., & Köhler, S. (2023). Perirhinal cortex automatically tracks multiple types of familiarity regardless of task-relevance. *Neuropsychologia*, 187, 108600.

<https://doi.org/10.1016/j.neuropsychologia.2023.108600>

Yonelinas, A. P. (2002). The nature of recollection and familiarity: A review of 30 years of research.

*Journal of Memory and Language*, 46(3), 441–517. <https://doi.org/10.1006/jmla.2002.2864>

Yonelinas, A. P., Aly, M., Wang, W. C., & Koen, J. D. (2010). Recollection and familiarity: Examining controversial assumptions and new directions. *Hippocampus*, 20(11), 1178–1194.

<https://doi.org/10.1002/hipo.20864>

Yonelinas, A. P., Otten, L. J., Shaw, K. N., & Rugg, M. D. (2005). Separating the Brain Regions Involved in Recollection and Familiarity in Recognition Memory. *The Journal of Neuroscience*, 25(11),

3002–3008. <https://doi.org/10.1523/JNEUROSCI.5295-04.2005>

Yonelinas, A. P., Ramey, M. M., & Riddell, C. (2022). *Recognition Memory: The Role of Recollection and Familiarity*. Handbook of Human Memory: Foundations and Applications. Oxford University Press.

Zhang, L. I., Tao, H. W., Holt, C. E., Harris, W. A., & Poo, M. (1998). A critical window for cooperation and competition among developing retinotectal synapses. *Nature*, 395(6697), 37–44.

<https://doi.org/10.1038/25665>

Zhang, W., Sun, J., & Tang, X. (2008). Cat head detection-how to effectively exploit shape and texture features. *European Conference on Computer Vision*, 802–816. Springer.

**TABLES****Table 1**

Number of images retained in memory ( $N_{ret}$ ) by the Hebbian and anti-Hebbian models for different dataset sizes.

| <b>Anti-Hebbian</b> | <b>Dataset Size</b> |           |            |            |            |             |             |              |
|---------------------|---------------------|-----------|------------|------------|------------|-------------|-------------|--------------|
|                     | <b>20</b>           | <b>40</b> | <b>100</b> | <b>200</b> | <b>400</b> | <b>1000</b> | <b>4000</b> | <b>10000</b> |
| Mean                | 20.00               | 39.860    | 97.90      | 189.680    | 357.70     | 744.720     | 1786.180    | 3171.760     |
| Std. Deviation      | 0.00                | 0.586     | 2.418      | 4.720      | 10.661     | 22.256      | 63.583      | 101.402      |
| Minimum             | 20.00               | 36.00     | 90.00      | 180.00     | 332.00     | 678.00      | 1632.00     | 2932.00      |
| Maximum             | 20.00               | 40.00     | 100.00     | 200.00     | 378.00     | 798.00      | 1928.00     | 3564.00      |
| <b>Hebbian</b>      |                     |           |            |            |            |             |             |              |
| Mean                | 19.920              | 38.380    | 73.440     | 91.760     | 100.140    | 107.820     | 128.560     | 150.760      |
| Std. Deviation      | 0.394               | 1.879     | 6.781      | 12.112     | 17.294     | 29.666      | 61.671      | 90.912       |
| Minimum             | 18.00               | 32.00     | 52.00      | 60.00      | 52.00      | 50.00       | 12.00       | 0.00         |
| Maximum             | 20.00               | 40.00     | 90.00      | 124.00     | 144.00     | 188.00      | 282.00      | 376.00       |

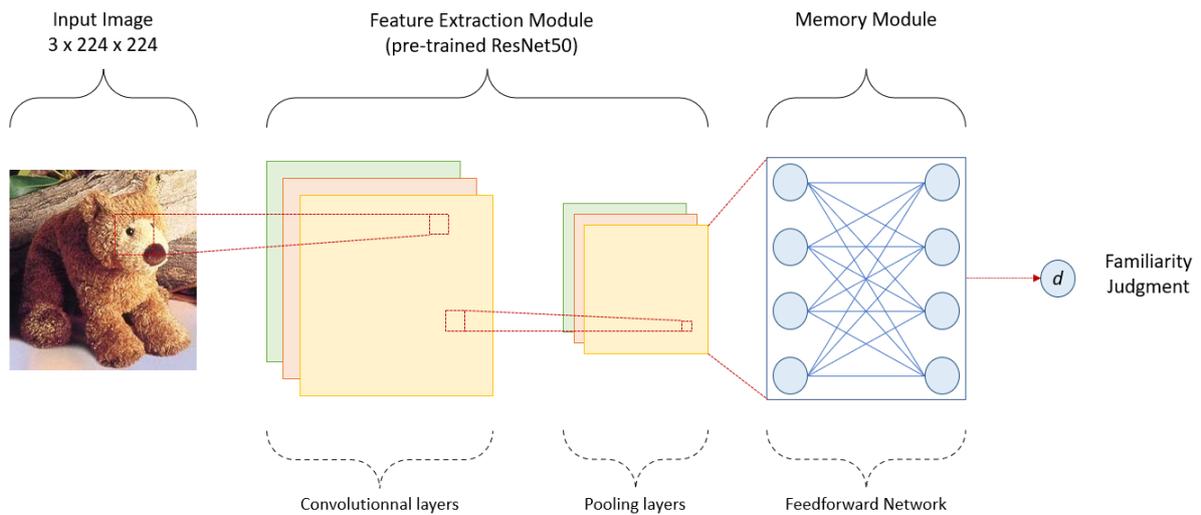
**Table 2**

Accuracy for the Hebbian and anti-Hebbian model computed during the testing phase in every dataset condition.

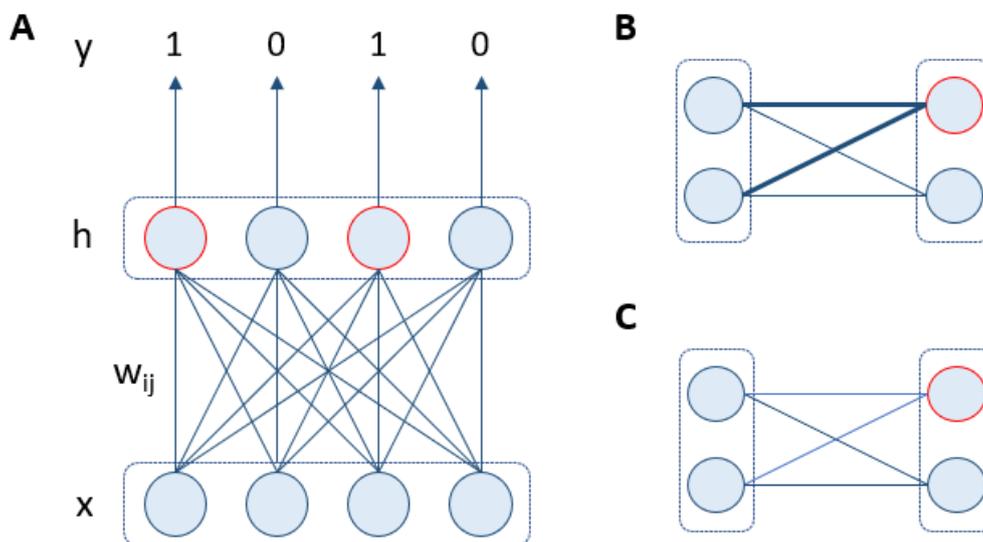
| <b>Dataset Size</b>          | <b>Model</b>        |                |
|------------------------------|---------------------|----------------|
|                              | <b>Anti-Hebbian</b> | <b>Hebbian</b> |
| Small dataset (N = 20)       | 100.00 (0.00)       | 99.80 (1.00)   |
| Medium dataset (N = 100)     | 98.90 (1.20)        | 86.70 (3.40)   |
| Large dataset (N = 1000)     | 87.20 (1.10)        | 55.40 (01.50)  |
| <b>Dataset Type (N = 40)</b> |                     |                |
| Heterogeneous (random)       | 99.80 (0.7)         | 98.00 (2.30)   |
| Mild homogeneity (dogs)      | 99.90 (0.5)         | 82.90 (5.10)   |
| High homogeneity (cats)      | 91.30 (4.50)        | 56.80 (7.50)   |

*Note.* Mean % over 100 runs (standard deviation).

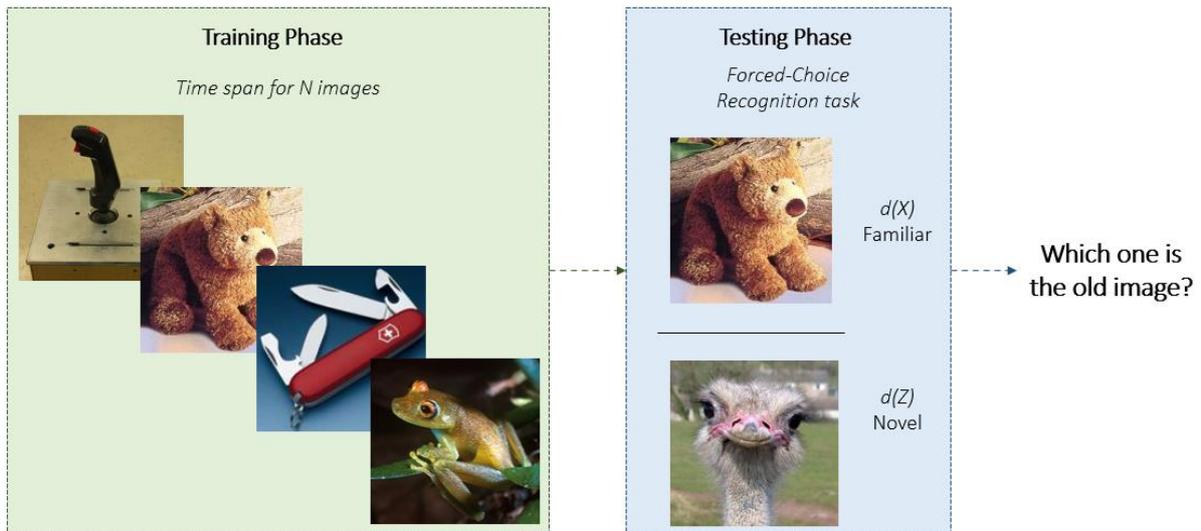
**FIGURES**



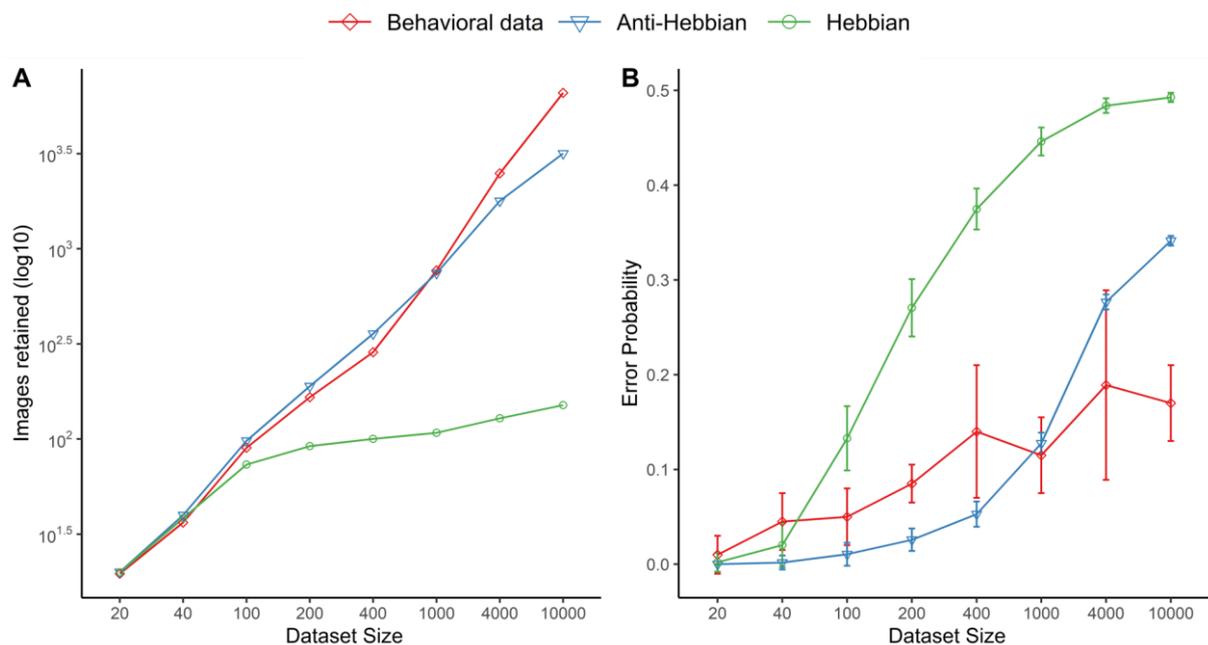
**Figure 1.** Global architecture of the models. An image goes through ResNet50 for features extraction then inside a memory module for learning. During the testing phase, a familiarity score  $d$  is computed for decision making.



**Figure 2.** Learning rules inside the memory modules. (A) General idea behind the functioning of the memory module. (B) Weight potentiation for active neurons in the Hebbian model. (C) Weight depression for active neurons in the anti-Hebbian model.

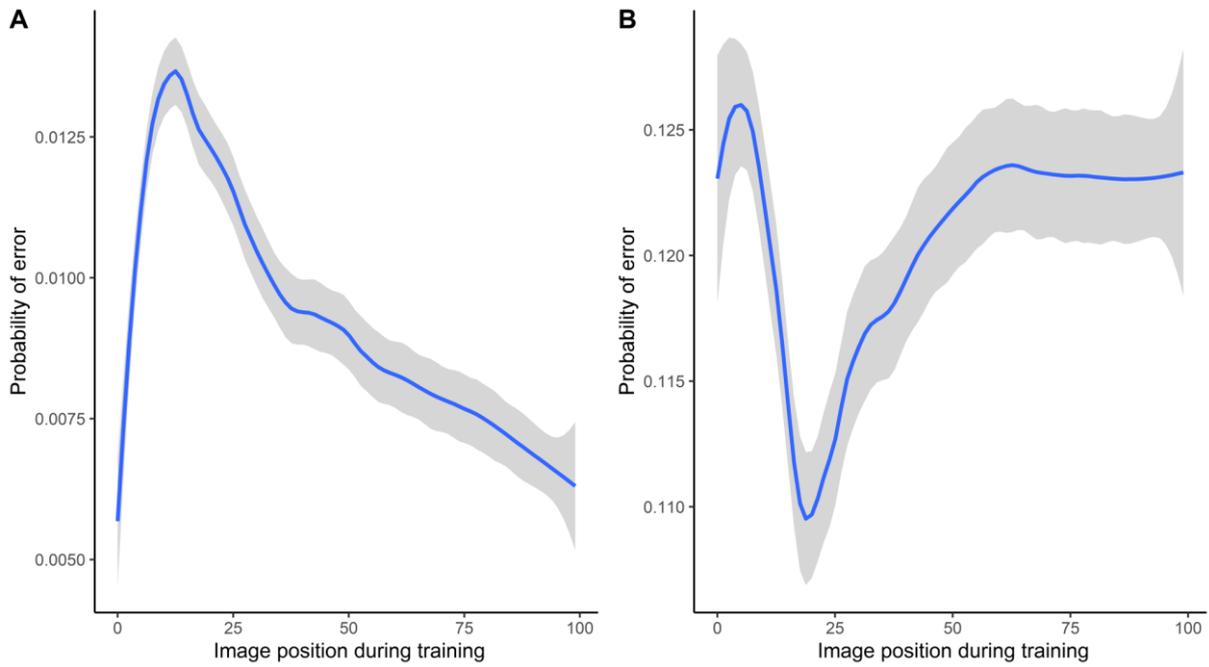


**Figure 3.** Simulation methodology. During the training phase,  $N$  images are learned one-by-one by the model. During the testing phase, pairs of images are presented to the model which has to decide which image is familiar.

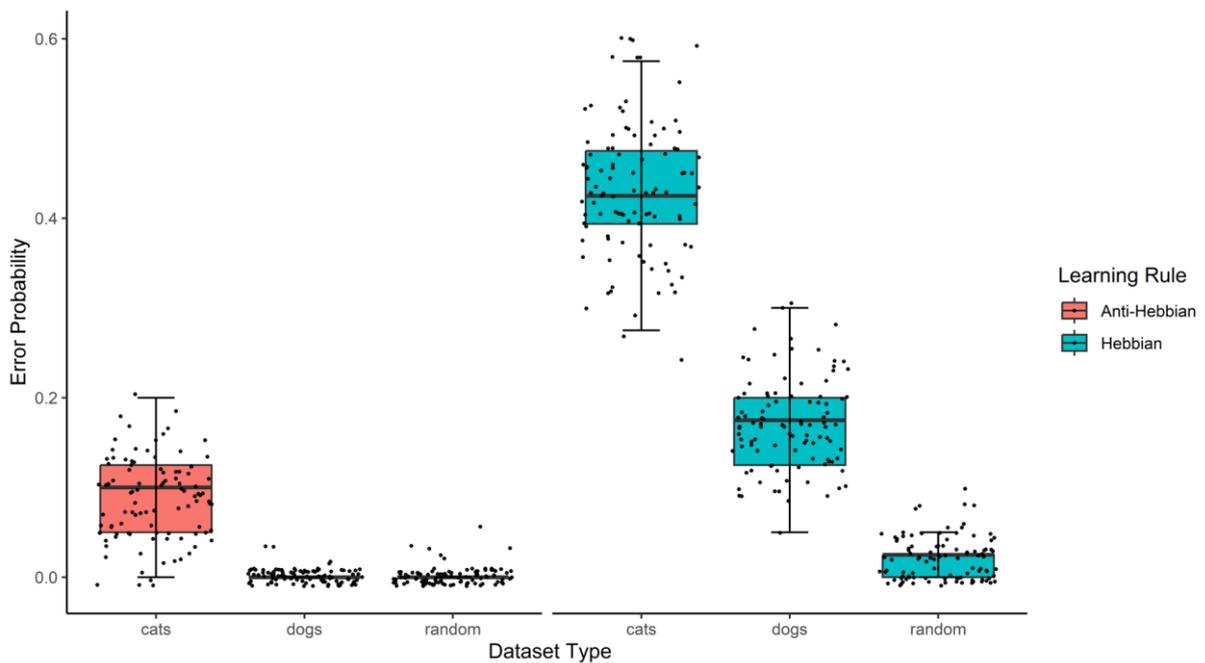


**Figure 4.** Results from the reproduction of Standing's experiment. (A) Number of images retained by the model as a function of the number of the dataset size during training (log<sub>10</sub> scale) (B) The probability of error as a function of the number of images learned during training. Red

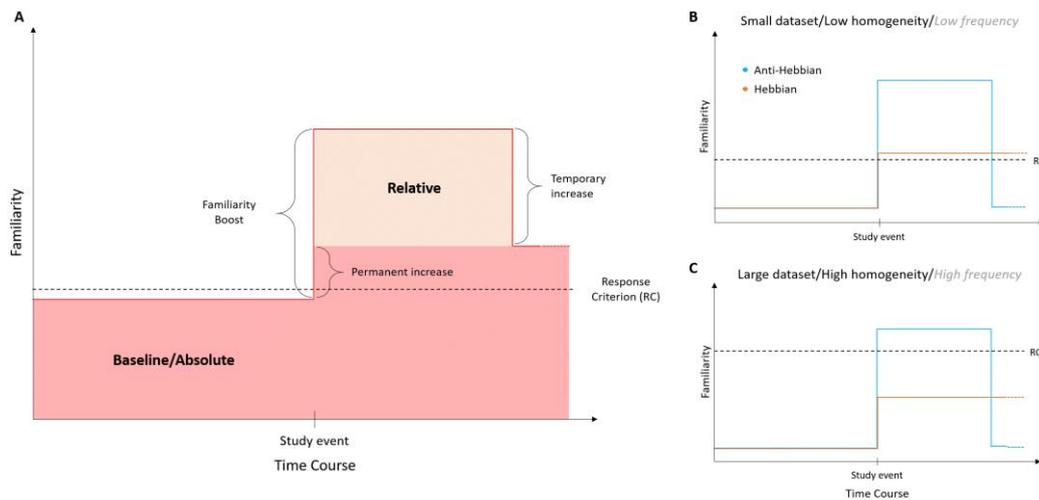
curves: Standing's behavioral data (Standing, 1973). Blue curves: performance for the anti-Hebbian solution. Green curves: performance for the Hebbian solution.



**Figure 5.** Mean probability of error for a given image and standard deviation (grey areas) as a function of the position of this image in the training phase ( $N = 100$ ). (A) Performance tested with the anti-Hebbian model. (B) Performance tested with the Hebbian model.



**Figure 6.** Mean probability of error and standard deviation when the two models are tested on dataset with different homogeneity levels ( $N = 40$ ).



**Figure 7.** Time course of the familiarity signal. (A) Collective contribution of both absolute and relative familiarity to the familiarity decision as described in Coane et al. (2011). (B) Separate contributions of the absolute familiarity as modeled by the Hebbian learning rule (orange curve) and the relative familiarity as modeled by the anti-Hebbian learning rule (blue curve) under high level of distinctiveness. (C) Separate contributions of the absolute familiarity as modeled by the Hebbian learning rule (orange curve) and the relative familiarity as modeled by the anti-Hebbian learning rule (blue curve) under low level of distinctiveness.