

1 **Title:** Meta-learning local synaptic plasticity for continual familiarity detection

2 **Authors:** Danil Tyulmankov, Guangyu Robert Yang, LF Abbott

3

4 **Abstract**

5 Over the course of a lifetime, a continual stream of information is encoded and retrieved from
6 memory. To explore the synaptic mechanisms that enable this ongoing process, we consider a continual
7 familiarity detection task in which a subject must report whether an image has been previously encountered.
8 We design a class of feedforward neural network models endowed with biologically plausible synaptic
9 plasticity dynamics, the parameters of which are meta-learned to optimize familiarity detection over long
10 delay intervals. After training, we find that anti-Hebbian plasticity leads to better performance than Hebbian
11 and replicates experimental results from the inferotemporal cortex, including repetition suppression. Unlike
12 previous models, this network both operates continuously without requiring any synaptic resets and
13 generalizes to intervals it has not been trained on. We demonstrate this not only for uncorrelated random
14 stimuli but also for images of real-world objects. Our work suggests a biologically plausible mechanism for
15 continual learning, and demonstrates an effective application of machine learning for neuroscience
16 discovery.

17

18 **Introduction**

19 Every day, a continual stream of sensory information and internal cognitive processing causes
20 lasting synaptic changes in our brains that alter our responses to future stimuli. It remains a mystery how
21 neural activity and local synaptic updates coordinate to support distributed storage and readout of
22 information and, in particular, how ongoing synaptic changes due to either new memories or homeostatic
23 mechanisms do not interfere with previously stored information.

24 Memory research in theoretical neuroscience and machine learning has addressed these questions
25 through modeling studies, but important features remain to be clarified. First, memories of an individual's
26 history are encoded in a one-shot manner – this is different from typical neural network models which use
27 a prolonged incremental training process to learn a complex task. Such training algorithms use a global
28 error signal and perform per-synapse credit assignment through knowledge of the entire network
29 (Rumelhart et al., 1986), whereas biological synapses typically only have access to local pre- and
30 postsynaptic activity (Hebb, 1949) and various modulatory signals (Frémaux and Gerstner, 2016; Gerstner
31 et al., 2018). Second, biological synapses change continually in response to ongoing activity, whereas
32 models commonly assume that synapses are fixed after training ends, such as the classical Hopfield
33 network (Hopfield, 1982) and most deep neural networks (LeCun et al., 2015). Unregulated continual
34 updating of synapses can cause catastrophic forgetting in which a network either erases previous memories
35 (Kirkpatrick et al., 2017; Zenke et al., 2017) or renders stored information unreadable (Parisi, 1986).
36 Recurrent neural networks are commonly used to perform tasks that involve memories sustained by neural
37 activity (Elman, 1991; Hochreiter and Schmidhuber, 1997; Mante et al., 2013), however most memories

38 are likely stored through synaptic potentiation and depression (Abbott and Nelson, 2000). Synaptic memory
39 has sometimes been studied through an ideal observer approach (Benna and Fusi, 2016; Fusi et al., 2005)
40 in which synaptic weights are directly accessible for readout, but biological organisms must read out
41 synaptic storage through neuronal activations. In fact, the readout may not be a dedicated circuit as in
42 attractor network models of memory (Hopfield, 1982), but rather manifested as a change in ongoing neural
43 processing (Hasson et al., 2015). Finally, machine learning research often eschews biological plausibility
44 and mechanistic understanding in favor of performance on benchmark tasks (Ba et al., 2016; Graves et al.,
45 2014, 2016; Miconi et al., 2018, 2019).

46 Familiarity detection – identifying whether a stimulus has been previously encountered – is a simple
47 and ubiquitous form of memory that serves as a useful testbed for addressing these issues. Classical
48 studies have demonstrated that human recognition memory capacity for images is "almost limitless" in a
49 two-alternative-forced-choice task with separate encoding and testing phases: retention follows a power
50 law as a function of the number of items viewed (Standing, 1973). Pioneering theoretical work has shown
51 that the number of memories stored by a familiarity detection network depends on the synaptic plasticity
52 rule and, in the case of uncorrelated inputs, capacity can scale proportionally to the number of synapses
53 (Bogacz and Brown, 2003). More recent behavioral work further demonstrates an impressive capacity in a
54 continual setting, the error rate as a function of the number of intervening items exhibiting a "power law of
55 forgetting" (Brady et al., 2008). Theoretical studies have shown that power-law forgetting is achievable by
56 synapses with metaplasticity, using both uncorrelated inputs (Fusi et al., 2005) and face images (Ji-An et
57 al., 2019). Neural signals of visual familiarity have been observed as reductions in responses to repeated
58 presentations of a stimulus, a phenomenon known as repetition suppression (Grill-Spector et al., 2006;
59 Meyer and Rust, 2018; Miller et al., 1991; Xiang and Brown, 1998). At the timescales relevant for this task
60 – one-shot memorization on the order of seconds and long-term forgetting on the order of days – this is
61 plausibly caused by depression of excitatory synapses or potentiation of inhibitory ones (Lim et al., 2015).

62 Previous modeling work on recognition memory used a predesigned architecture and plasticity rule
63 and both empirical and analytic evaluation of performance (Androulidakis et al., 2008; Bogacz and Brown,
64 2003; Norman and O'Reilly, 2003; Sohal and Hasselmo, 2000). An emerging approach uses a machine
65 learning technique known as "meta-learning," or "learning how to learn" (Confavreux et al., 2020; Thrun
66 and Pratt, 2012), which uses optimization tools to rapidly search for mechanisms that artificial neural
67 networks can use to solve a learning/memory task. In contrast to hand-designed models, meta-learning
68 enables unbiased exploration of a large family of architectures and plasticity rules. Importantly, it is possible
69 to impose constraints that ensure biological plausibility (Bengio et al., 1991). For example, given a network
70 architecture and a family of biologically plausible plasticity rules with tunable parameters, the meta-learning
71 algorithm can search for the optimal parameters for memorizing a set of inputs.

72 In this work, we consider a family of models that recognize previously experienced stimuli and,
73 importantly, learn and operate continuously without separate "learning" and "testing" phases. We
74 investigate a feedforward network architecture with continual Hebbian plasticity in its synaptic weights.

75 Parameters governing plasticity and other network parameters are meta-learned using gradient descent to
76 optimize the continual familiarity detection process. To isolate synaptic plasticity as the unique memory
77 mechanism, we avoid recurrent connectivity that could store memory through maintained neuronal
78 activations. This architecture, unlike recurrent networks, generalizes naturally over a range of repeat
79 intervals even if trained on a single interval. We show that an anti-Hebbian plasticity rule (co-activated
80 neurons cause synaptic depression) enables repeat detection over longer intervals than a Hebbian rule,
81 and this is the solution most frequently found by meta-learning. This rule leads to experimentally observed
82 features such as repetition suppression in the hidden layer neurons. Critically, the capacity of these
83 networks remains constant over time, so they can be continually fed new inputs with no reduction in steady-
84 state memory performance.

85

86 Results

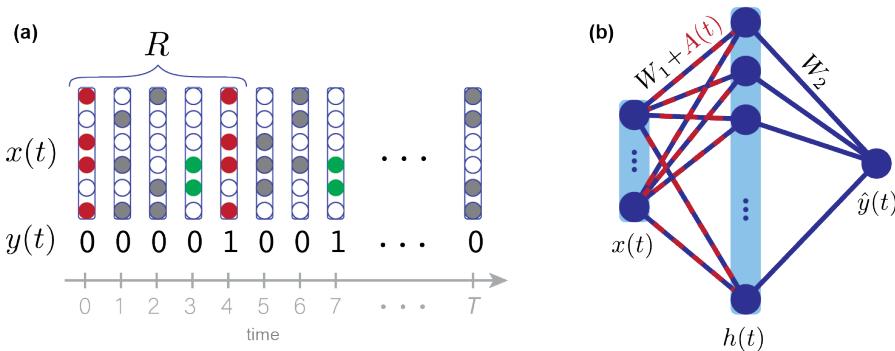
87 Continual familiarity detection task

88 We consider a continual familiarity detection task (Fig 1a) in which a stream of stimuli is presented
89 to a network. With probability $1 - p$, the stimulus at time t is chosen as a randomly generated d -dimensional
90 binary vector $x(t)$, where each component is either +1 or -1 (note that for sufficiently large d , spurious
91 chance repeats are extremely unlikely). With probability p , the stimulus is a copy of the stimulus presented
92 R time steps ago, so that $x(t) = x(t - R)$. However, we ensure that a stimulus is repeated at most once so,
93 if $x(t - R)$ is already a repeat, i.e. $x(t - R) = x(t - 2R)$, a new $x(t)$ is generated. As a result, the fraction
94 of novel stimuli, which we call f , is not equal to $1 - p$, but rather $f = \frac{1}{1+p}$. We use $p = \frac{1}{2}$, so $f = \frac{2}{3}$. The output
95 of the network should be $y(t) = 0$ if $x(t)$ is novel and $y(t) = 1$ if it is familiar, i.e. has appeared previously.

96 The accuracy of the network (P_{correct} , the probability of correctly responding to a stimulus) depends
97 on two factors: the true positive rate (P_{TP} , the probability of correctly reporting a repeated stimulus as
98 "familiar"), and the false positive rate (P_{FP} , the probability of incorrectly reporting a novel stimulus as
99 "familiar"). These two factors are weighted by the fraction of novel stimuli f , so that $P_{\text{correct}} = (1 - f)P_{TP} +$
100 $f(1 - P_{FP})$. Through our choice of loss function (*Online methods*), we are effectively training the networks
101 to maximize accuracy, so the "chance" level performance is f (for $f > \frac{1}{2}$), which a network can achieve by
102 reporting all stimuli as novel ($P_{TP} = P_{FP} = 0$).

103 In our paradigm, a given dataset has a single repeat interval R , which differs slightly from previously
104 studied experimental paradigms (Brady et al., 2008; Meyer and Rust, 2018). However, we evaluate
105 performance on multiple datasets with various values of R . For testing, this is analogous to evaluating a
106 single dataset with multiple repeat intervals and computing accuracy for each interval separately. We use
107 this approach because it allows us to test generalization by training on one value of R and testing on others.
108 It also allows us to train the network to its maximal capacity by gradually increasing R during "curriculum
109 training", and simplifies analytic calculations.

110 We begin by considering familiarity detection for uncorrelated stimuli, but, in later sections, we
 111 generalize to a task that requires simultaneous familiarity detection and binary classification, and to a
 112 dataset of real-world object images.



113
 114 **Figure 1. Continual familiarity detection task and HebbFF model.** (a) The continual familiarity detection task. Given a continual
 115 stream of stimuli $x(t)$, the desired output is $y(t) = 1$ if the stimulus has appeared previously and $y(t) = 0$ otherwise. For a given
 116 dataset, repeat stimuli always appear at an interval R after their first presentation. Although the task is continual, for the purposes of
 117 network training we use a finite-duration trial of length $T \gg R$. (b) The HebbFF network architecture. A feedforward layer is endowed
 118 with ongoing Hebbian plasticity, the parameters of which are optimized using stochastic gradient descent. The hidden units are linearly
 119 read out to produce the network's estimate of familiarity $\hat{y}(t)$.

120

121 HebbFF network architecture

122 To investigate the effectiveness of synaptic plasticity for solving this task, we use a feedforward
 123 neural network with a single hidden layer and, to implement the memory function, activity-dependent
 124 ongoing Hebbian plasticity (HebbFF) (Fig 1b). We purposefully do not include any recurrent connections to
 125 ensure that memory cannot be stored through persistent neuronal activity, thus isolating synaptic plasticity
 126 as the only possible memory mechanism.

127 In the HebbFF network, a group of hidden layer neurons with firing rates given by an N -dimensional
 128 vector $\mathbf{h}(t)$, receives a d -dimensional input $\mathbf{x}(t)$. The variable t indexes stimulus presentations which occur
 129 sequentially, so we refer to it as "time". The input to each hidden-layer neuron is weighted by its
 130 corresponding synaptic strength and then transformed into a hidden-layer firing rate through a nonlinear
 131 activation function $\sigma(\cdot)$. The synaptic strength between the postsynaptic neuron with rate $h_i(t)$ and the
 132 presynaptic neuron carrying the input $x_j(t)$ is the (i,j) component of an N -by- d matrix that is the sum of a
 133 fixed matrix \mathbf{W}_1 and a plastic matrix $\mathbf{A}(t)$. Thus, the firing rate of the hidden layer is given by

$$139 \quad \mathbf{h}(t) = \sigma((\mathbf{W}_1 + \mathbf{A}(t))\mathbf{x}(t) + \mathbf{b}_1)$$

140 where $\sigma(\cdot)$ is the logistic function applied element-wise and \mathbf{b}_1 is a vector representing baseline currents
 141 into the hidden layer. The plastic matrix $\mathbf{A}(t)$ is updated at every time step: its (i,j) component decays by
 142 a factor $0 < \lambda < 1$ and is incremented by a Hebbian product of the pre- and postsynaptic activities,
 143 $h_i(t)x_j(t)$. A plasticity rate parameter $-\infty < \eta < \infty$ controls the sign and magnitude of this increment. In
 144 matrix form, the synaptic update rule is then

$$145 \quad \mathbf{A}(t+1) = \lambda\mathbf{A}(t) + \eta\mathbf{h}(t)\mathbf{x}(t)^T$$

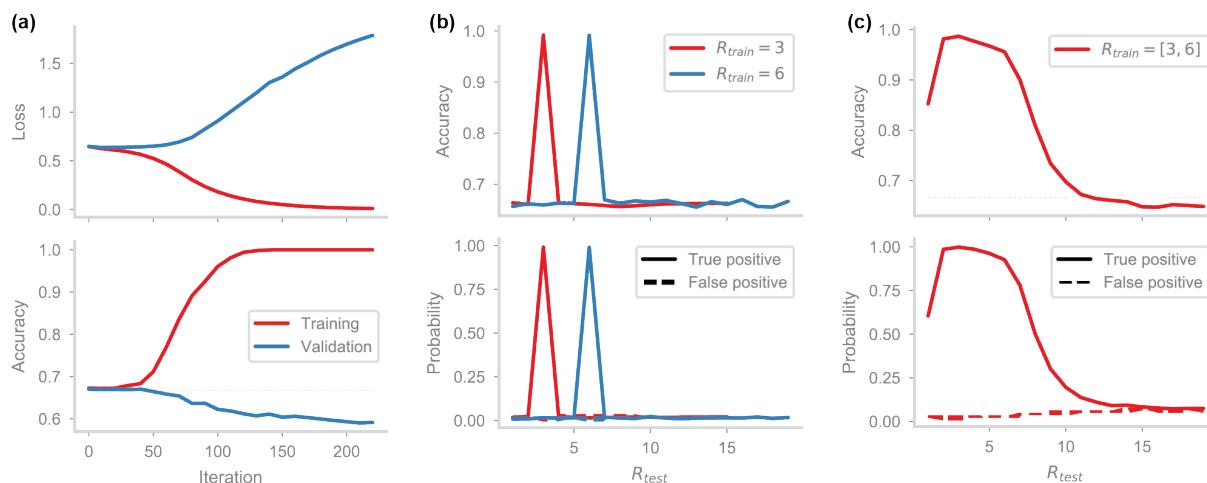
141 Finally, the output of the network $\hat{y}(t)$ is a linear readout of the hidden layer and, since the target $y(t)$ is
 142 binary, we bound the readout with the logistic function,

$$143 \quad \hat{y}(t) = \sigma(\mathbf{W}_2 \mathbf{h}(t) + \mathbf{b}_2)$$

144 The response of the network is taken to be familiar if $\hat{y}(t) > 1/2$ and novel otherwise.

145 To construct the network, we use backpropagation through time (BPTT) to "meta-learn" the
 146 parameters $\mathbf{W}_1, \mathbf{b}_1, \mathbf{W}_2, \mathbf{b}_2, \lambda, \eta$, which are fixed once training is completed (*Online methods*). The continual
 147 familiarity detection task – the "learning" task – is then performed exclusively by the ongoing synaptic
 148 dynamics of $A(t)$, determined by the fixed parameters. These dynamics are a biologically plausible
 149 mechanism for solving the continual memory task, but BPTT is simply used as an optimization tool to find
 150 suitable parameters of the network.

151



152
 153 **Figure 2. RNN performance on continual familiarity detection.** (a) Training an LSTM ($d = 100$ input dimension, $N = 100$ recurrent
 154 units) on a single dataset of a familiarity detection task ($T = 500$ stimulus presentations, repeat interval $R = 3$). Although the loss (top)
 155 approaches zero and accuracy (bottom) approaches 1 on the training dataset (red curves), performance on a validation dataset (blue)
 156 with the same parameters fails to generalize even when tested in-distribution with the same R . (b) Training the RNN using "infinite
 157 data." New datasets, each with $R = 3$, are generated at every epoch of training (red). Accuracy (top), as well as true positive and false
 158 positive probabilities (bottom) is shown as a function of the repeat interval on validation datasets. The same is repeated with another
 159 RNN using $R = 6$ (blue). The RNNs perform well in-distribution on datasets with the same repeat interval as used during training, but
 160 fail to generalize out-of-distribution to other repeat intervals. (c) Training the RNN with "infinite data," using datasets with both repeat
 161 intervals $R = 3$ and $R = 6$. The RNN interpolates between the intervals – performance is high when tested on repeat intervals $3 \leq R \leq$
 162 6 – but fails to extrapolate. Performance quickly drops for longer repeat intervals, and even for shorter ones.
 163

164 HebbFF generalizes across datasets and repeat intervals

165 As a benchmark for comparing HebbFF performance, we first train a long short-term memory
 166 (LSTM) network (Hochreiter and Schmidhuber, 1997) – a recurrent neural network (RNN) architecture well-
 167 suited for memory performance – on the continual familiarity detection task. Unlike HebbFF, which stores
 168 its input history in the plastic synaptic matrix $A(t)$, an RNN uses ongoing neuronal activity.

169 If we train the RNN using a single dataset with $T = 500$ image presentations (*Online methods*) and
 170 a repeat interval of $R = 3$, it successfully learns the training set, but entirely fails to generalize to new test
 171 sets with the same R (Fig. 2a). This is because the RNN learns to perform classification, not familiarity
 172 detection. To fix this, we use an "infinite data" approach in which we generate a new dataset for every

173 iteration of BPTT, each with the same value of $R = 3$. Trained in this way, the RNN now generalizes "in-
174 distribution" across datasets with $R = 3$ (i.e. to datasets drawn from the same distribution as the training
175 data, which is parameterized by R), but fails to generalize "out-of-distribution" to data with any other value
176 of R (i.e. to datasets from a different distribution) (Fig 2b). The same result holds if we train another RNN
177 with $R = 6$ (Fig 2b). We can further train the RNN with items spaced at intervals of both $R = 3$ and $R = 6$
178 (i.e. the value of R is chosen randomly for each familiar stimulus rather than being fixed). While the network
179 can interpolate between the trained values, it does not extrapolate well to larger or smaller ones (Fig 2c).
180 Although it is likely possible to train the RNN to perform well for multiple values of R by using more complex
181 training schedules, we believe that poor out-of-distribution generalization is a bottleneck of the RNN
182 approach.

183 In contrast, the HebbFF network exhibits both in-distribution and out-of-distribution generalization.
184 Even when trained on a single dataset with a fixed repeat interval R , the network generalizes not only to
185 new test sets with the same R , (Fig 3a) but even to those with different R values. Trained with "infinite data"
186 (the scheme we use in general), HebbFF generalizes to datasets with smaller and even larger R values
187 (Fig 3b). If we match the number of dynamic variables rather than the number of hidden neurons, HebbFF
188 still shows superior generalization compared to the RNN (Fig S1). This qualitative difference in performance
189 suggests that Hebbian plasticity provides a powerful inductive bias in the form of a more "natural"
190 mechanism of memory for the purpose of familiarity detection.

191 The generalization performance of HebbFF is due to the fact that the memory representation of an
192 item does not change over time, other than being scaled by a factor. A stimulus $\mathbf{x}(t)$ is initially stored as
193 the outer product of $\mathbf{h}(t)$ and $\mathbf{x}(t)$, multiplied by the plasticity rate η . The plastic component of the
194 connectivity matrix also contains terms arising from previously stored memories which, for the purposes of
195 this particular stimulus, act as additive noise ε :

$$\mathbf{A}(t+1) = \eta \mathbf{h}(t) \mathbf{x}(t)^T + \varepsilon$$

196 As subsequent stimuli are presented, the representation of $\mathbf{x}(t)$ maintains the same form, so that k time
197 steps later it is still stored as the outer product of $\mathbf{h}(t)$ and $\mathbf{x}(t)$, scaled by a factor λ^k :

$$\mathbf{A}(t+k) = \lambda^k \eta \mathbf{h}(t) \mathbf{x}(t)^T + \lambda^k \varepsilon + \varepsilon'$$

199 where further additive noise ε' arises from stimuli presented after $\mathbf{x}(t)$.

201 Unlike HebbFF, RNNs are poor at generalizing across intervals R because the dynamics of their
202 units allow the memory representation of a stimulus to change arbitrarily over time. The RNN only generates
203 the appropriate representation at the time when a query is expected, namely after a delay equal to the value
204 of R used during training. This makes it difficult to generalize across intervals.

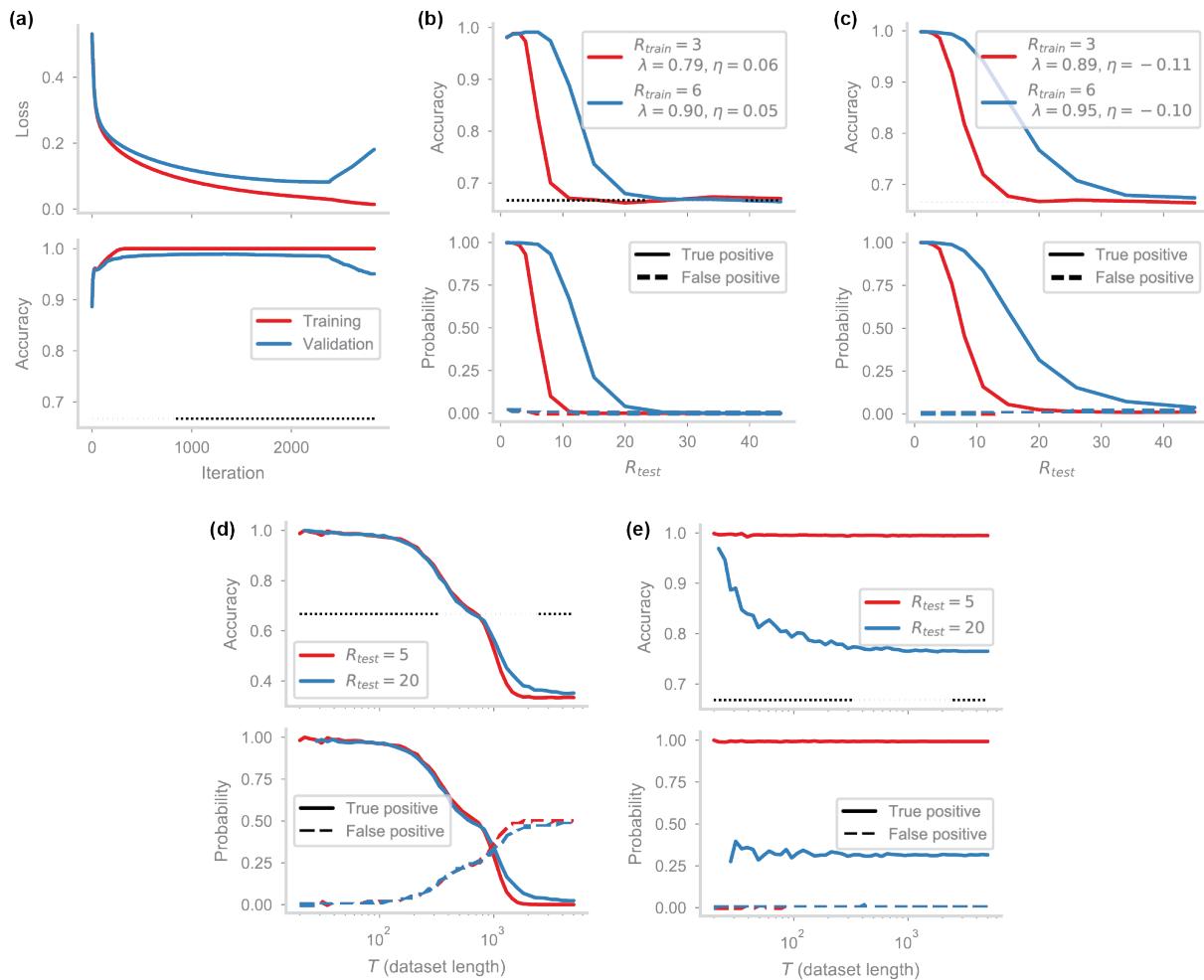


Figure 3. Hebbian vs. anti-Hebbian plasticity and continual operation. (a) Training the HebbFF network ($d = N = 100$), as in Fig 2a. Both training and validation loss decrease, and accuracy increases, for a single instance of the dataset with $R = 3$, indicating in-distribution generalization. Over many iterations, however, overtraining occurs due to the use of a single dataset, increasing the final validation loss. (b) HebbFF network trained with "infinite data" as in Fig 2b ($R = 3$, red; $R = 6$, blue) shows not only in-distribution generalization to any dataset with $R = 3$ ($R = 6$, resp.), but also out-of-distribution to datasets with any smaller R and some larger R 's. (c) HebbFF with a different initialization converges to a qualitatively different solution with a negative learning rate η , an anti-Hebbian learning rule, in contrast to the Hebbian solution in (b). The anti-Hebbian solution shows generalization performance over a larger range of R values than Hebbian. (d) Model from (Bogacz and Brown, 2003), evaluated on the continual familiarity detection task, varying the length T of the trial. Accuracy (top) is near-perfect regardless of the repeat interval R (blue vs. red curve) until the model reaches its capacity ($P^* \approx 100$ for network size $d = N = 100$) because the model reliably stores the first P^* patterns. Accuracy rapidly drops below chance for $T > P^*$ as the model begins to report familiar stimuli as novel (see Fig S2b). (e) HebbFF network operates continuously, as its accuracy (top) is consistent with the generalization curve from (c), with near-perfect performance for $R_{test} = 5$ and above 80% for $R_{test} = 20$ for any trial length. True and false probabilities (bottom) are better representations as accuracy (top) is artificially higher for small T due to the low proportion of familiar stimuli.

Anti-Hebbian outperforms Hebbian plasticity

The plasticity rate η in HebbFF can be positive or negative, resulting in either Hebbian or anti-Hebbian plasticity. For the Hebbian solution with $\eta > 0$, synapses are potentiated in response to a stimulus. When it is repeated, the hidden layer activity is higher than for a novel stimulus due to the increased strength of the synapses storing the memory. For anti-Hebbian plasticity, $\eta < 0$, synapses are depressed when a memory is stored. In this case, the hidden layer activity is lower for a familiar stimulus than a novel stimulus, which is consistent with experimental results of repetition suppression (Grill-Spector et al., 2006; Meyer and

228 Rust, 2018; Xiang and Brown, 1998). Furthermore, the meta-learning algorithm is more likely to converge
229 to the anti-Hebbian solution, especially when trained with a relatively large repeat interval, even if the initial
230 value of η is positive, and almost always does so when the initial value is negative.

231 Interestingly, anti-Hebbian plasticity enables successful familiarity detection over considerably
232 longer intervals than a Hebbian rule (Fig 3c). To understand this, note that the memory of a stimulus is
233 degraded in two ways: additional plasticity events obscure existing memories, and the plastic weights decay
234 over time. With an anti-Hebbian plasticity rule, the hidden layer activation $\mathbf{h}(t)$ is close to zero for a familiar
235 stimulus due to repetition suppression. As a result, the plasticity update $\eta \mathbf{h}(t) \mathbf{x}(t)^T$ when the stimulus is
236 repeated is negligible – as if a stimulus was not presented at that time step. This effectively reduces the
237 number of plasticity events, reducing the disruption of existing memories. As a secondary effect, the smaller
238 number of plasticity events allows a larger λ (smaller decay rate) to be used while still controlling the
239 amplitude of plastic weights. This slower decay rate further extends the lifetime of the memory. Due to their
240 superior performance and consistency with experimental results, we only consider anti-Hebbian solutions
241 throughout the following sections.

242
243 *HebbFF learns continually without catastrophic forgetting*
244 Previous modeling work using anti-Hebbian plasticity mechanisms for familiarity detection (Bogacz
245 and Brown, 2003) focused on a paradigm used in classic studies of recognition memory (Standing, 1973)
246 in which subjects are serially presented an entire dataset and later asked to identify which stimulus is
247 familiar in a **two-alternative-forced-choice** (2AFC) test. Analogously, this previous modeling work used
248 explicit "learning" and "testing" phases and demonstrated an impressive capacity for recognition memory
249 (Bogacz and Brown, 2003) (Fig S2a). When evaluated on the continual memory task that we use, the
250 Bogacz-Brown model has near-perfect performance if the number of stimuli T in the dataset is smaller than
251 the model's capacity P^* , independent of the value of the repeat interval R (Fig 3d). That is, the model
252 successfully stores all $T < P^*$ stimuli. As the dataset size T increases, however, the model performance
253 declines due to catastrophic interference (Fig 3d, S2b; *Online methods*). To store additional memories, the
254 old memories must be removed by resetting the synaptic weights.

255 In contrast, the HebbFF model operates continually rather than using separate learning and
256 evaluation phases. Its performance is independent of the length of the dataset, and it can operate
257 continuously without any need to reset the synaptic weights. For example, a HebbFF network trained with
258 $R = 5$ operates at near-perfect performance irrespective of the duration of the trial T when tested with $R =$
259 5 (Fig 3e). Similarly, when tested with $R = 20$, it operates continually at near 80% accuracy (Fig 3e), as
260 expected from the generalization curve in Fig 3c (note that for small T the accuracy (Fig 3e, top, blue) is
261 transiently elevated because the fraction of novel stimuli is more than $\frac{2}{3}$). In other words, the model has a
262 moving window in time within which it can successfully detect a familiar stimulus, and it forgets old stimuli
263 gracefully without suffering from catastrophic interference.

264

265 *A uniform readout*

266 Up to now, the readout of the hidden layer, $\hat{y}(t) = \sigma(\mathbf{W}_2\mathbf{h}(t) + \mathbf{b}_2)$, has involved a trained matrix
267 \mathbf{W}_2 . Because the inputs to the network are uniformly random, there is no reason to believe that, as far as
268 the output is concerned, one hidden unit would be statistically different than any other (although they may
269 still be statistically dependent). This observation leads us to restrict \mathbf{W}_2 to be a scaled 1-by- N matrix of
270 ones, $\mathbf{W}_2 = \alpha_2[1, \dots, 1]$, where α_2 is a trained scalar. Similarly, we restrict $\mathbf{b}_1 = \beta_1[1, \dots, 1]^T$. We verified that
271 performance is not affected by this choice of output weights (Fig S3a), the distribution of $\hat{y}(t)$ for familiar
272 and novel stimuli is the same (Fig S3b), the readout vector \mathbf{W}_2 and the bias term \mathbf{b}_1 have similar features
273 (Fig S3c), and anti-Hebbian plasticity is still the preferred form of plasticity. We use this uniform readout
274 throughout the remainder of our studies unless stated otherwise.

275

276 *Storage and readout mechanisms*

277 In the HebbFF network, the hidden layer plays a dual role. On the one hand, it must produce a
278 reliable familiarity signal for the readout to decode. On the other, it must create a robust representation of
279 the input stimulus during the Hebbian plasticity update. The hidden activity is controlled by the fixed
280 parameters \mathbf{W}_1 and \mathbf{b}_1 , as well as the plastic matrix $\mathbf{A}(t)$. Here, we investigate how \mathbf{W}_1 , \mathbf{b}_1 and $\mathbf{A}(t)$, impact
281 these two aspects of the familiarity detection task.

282 Networks trained with larger R have sparser hidden unit activity (Fig 4a-c): the sparser the activity,
283 the less plasticity is evoked, and thus the longer memories can be retained without overwriting. In the
284 limiting case we might expect that exactly one neuron is active for a novel stimulus and none are active for
285 familiar stimuli. Associated with this increased sparsity in activity, \mathbf{W}_1 is also sparser for larger R (Fig 4d-f).

286 To isolate the effect of \mathbf{W}_1 on hidden unit activity, we compute a histogram of the input current into
287 the hidden layer due to the non-plastic synapses, $\mathbf{W}_1\mathbf{x}(t) + \mathbf{b}_1$, across units and across time (Fig 4g-i). As
288 R increases, the distribution becomes multi-modal. In general, the number of peaks in this distribution
289 depends on the number of large-magnitude values of \mathbf{W}_1 per row. Critically, due to the logistic function
290 nonlinearity, only the rightmost peak in Fig 4i is large enough to elicit appreciable activity in the hidden
291 layer. This peak drives the small number of hidden units that are significantly activated by a novel stimulus.
292 In other words the \mathbf{W}_1 matrix acts like a hash function to select a small subset of hidden units to store the
293 memory of a given stimulus.

294 We next consider the effect of $\mathbf{A}(t)$, focusing on the network trained to maximum capacity (Fig
295 4c,f,i) (see next section). For a novel stimulus, the distribution of the input current due to the plastic
296 synapses $\mathbf{A}(t)\mathbf{x}(t)$ is unimodal and symmetric about zero (Fig 4j-l). For a familiar stimulus, however, there
297 is an additional peak at approximately $\lambda^{R-1}\eta d$. This peak is due to the dot product of the input vector $\mathbf{x}(t -$
298 $R)$ (stored in the matrix $\mathbf{A}(t)$ as $\lambda^{R-1}\eta\mathbf{h}(t-R)\mathbf{x}(t-R)^T$), and the familiar input vector $\mathbf{x}(t) = \mathbf{x}(t - R)$.
299 Importantly, the neurons that exhibit this behavior are the same neurons that were active due to \mathbf{W}_1 when
300 the stimulus was novel. This implies that, in addition to selecting a subset of neurons for storage of a

301 memory, \mathbf{W}_1 also allows the system to probe those same neurons during recall. Thus, \mathbf{W}_1 serves as an
302 addressing function during memory recall.

303 Finally, the total hidden layer input current is the sum of these two components, $(\mathbf{W}_1 + \mathbf{A}(t))\mathbf{x}(t) +$
304 \mathbf{b}_1 (Fig 4m-o). Comparing Fig4i and Fig4o, we see that the large central symmetric mode of the $\mathbf{A}(t)\mathbf{x}(t)$
305 distribution does not significantly impact the total hidden layer input current. Rather, the familiarity signal
306 arises because the smaller peak of the $\mathbf{A}(t)\mathbf{x}(t)$ distribution pushes the rightmost peak of the $\mathbf{W}_1\mathbf{x}(t) + \mathbf{b}_1$
307 distribution below zero (Fig 4m). Anti-correlation between the two input currents for familiar stimuli (Fig 4p-
308 r) indicates that this shift is caused by the input current from the plastic component of the synapse cancelling
309 the input current from the fixed component, resulting in lower activation, i.e. repetition suppression.

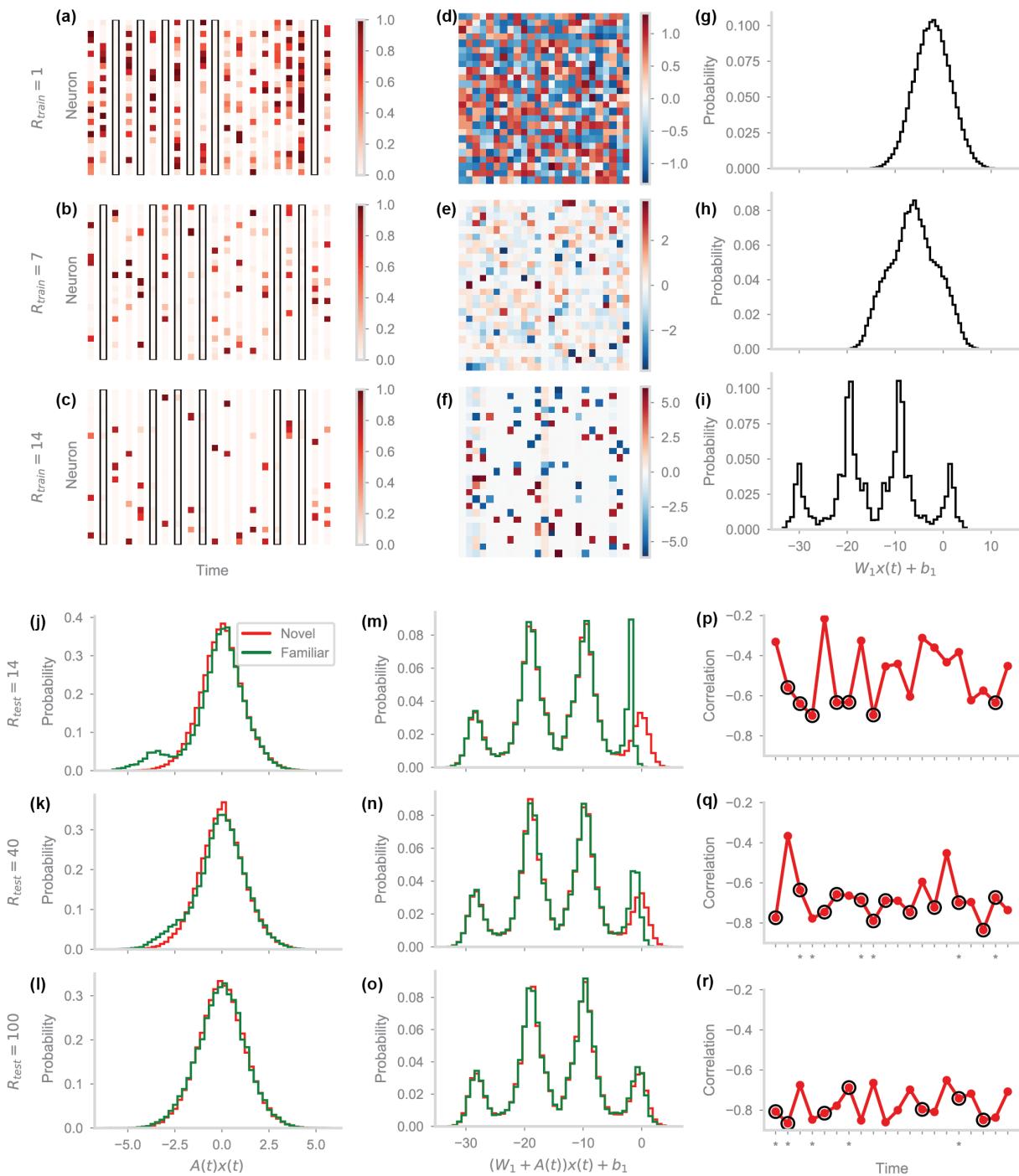


Figure 4. Storage and readout mechanism. (a-c) Hidden layer activity $h(t)$ over 20 consecutive timesteps for networks with input dimension $d = 25$ and $N = 25$ hidden units, trained on datasets with $R = 1, 7$, or 14 , respectively. Familiar stimuli (black rectangles) cause silencing i.e. repetition suppression of hidden layer activity. Activity for novel stimuli becomes sparser for networks trained with larger R . (d-f) Static weight matrix W_1 of the networks from (a-c). The weight matrix becomes sparser, and individual weight magnitudes increase for networks trained with larger R , enabling more sparse activity in the hidden layer for novel stimuli. (g-i) Distributions of input current into the hidden layer due to the static component of the synapses, i.e. the matrix W_1 and bias b_1 , for the networks from (a-c). We do not distinguish familiar and novel stimuli since the current due to the static component is the same regardless of novelty. For networks trained with larger R , the distribution becomes multi-modal, with the number of modes equal (approximately) to the number of high-magnitude values per row of W_1 , plus one. Due to the bias, only the rightmost mode has the potential to produce firing rates that are significantly above zero. (j-l) Distributions of input current into the hidden layer due to the plastic component of the synapses, i.e. the matrix $A(t)$, for novel (red) and familiar (green) stimuli. We only consider the trained network from (c,f,i) and evaluate its behavior on test sets with $R = 14, 40$, or 100 , corresponding to perfect, intermediate, and chance performance.

310

311

312

313

314

315

316

317

318

319

320

321

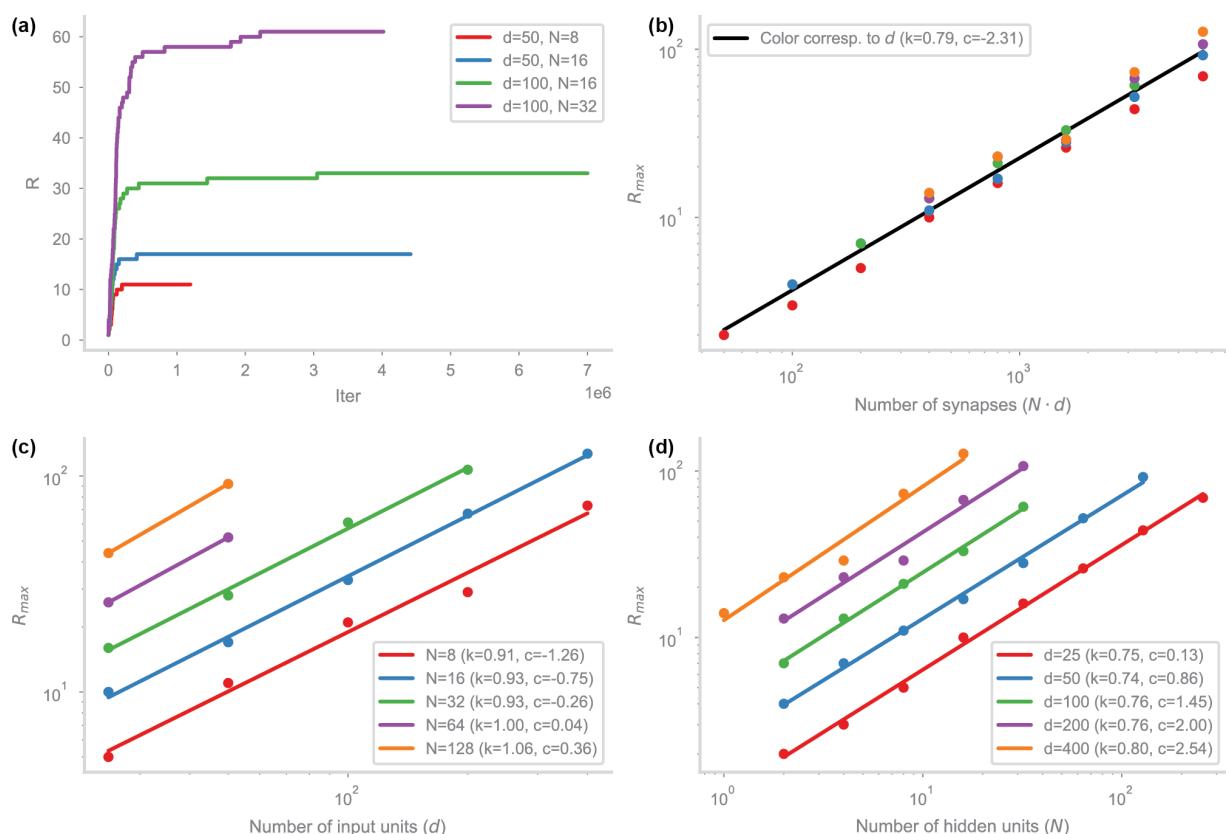
322

323 accuracy. The large central mode occurs due to stored stimuli uncorrelated with the input stimulus $x(t)$. In the novel case, the input
 324 is uncorrelated with all the stored stimuli by definition, and thus there is only one mode. Similarly, in the familiar case with a long delay
 325 interval $R = 100$ the stored stimulus has decayed sufficiently that its signal is lost. In the case of familiar stimuli presented at shorter
 326 delay intervals, $R = 14$ or 40, there is an additional mode due to the correlation between the input $x(t)$ and its copy $x(t - R)$ previously
 327 stored in the plastic matrix $A(t)$. (m-o) Distributions of the total input current into the hidden layer on test sets with $R = 14, 40$, or 100.
 328 Only the values above zero cause high firing rates after applying the logistic sigmoid nonlinearity. Since all the input currents are low
 329 for familiar stimuli (green) for small values of R , there is repetition suppression. (p-r) Correlation between the input current into the
 330 hidden layer from static and plastic synapse components at each of 20 consecutive timepoints. Asterisks indicate output response
 331 errors. For sufficiently small R , the input currents are more anti-correlated for familiar stimuli (black circles) than for novel. Combined
 332 with the distributions of input currents, this indicates that the units receiving positive input current from the static synapses receive
 333 negative input current from the plastic synapses.

334

335 Curriculum training and memory capacity

336 A randomly-initialized HebbFF network is typically unable to find a solution if trained with a large
 337 value of R (e.g. a network with $d = 50, N = 100$ rarely converges if trained on $R > 10$). Instead, we use a
 338 curriculum training procedure to bootstrap the optimized solution. First, the network is trained on data with
 339 $R = 1$, using the "infinite data" regime. Once the accuracy is above 99%, R is incremented by one and
 340 training continues on data with $R = 2$. This process continues until R becomes large enough that the
 341 network cannot find a solution with accuracy above 99%, i.e. if R is not incremented for at least 2 million
 342 iterations (Fig 5a). We thus define the memory capacity R_{\max} as the largest value of R for which the
 343 familiarity detection accuracy is above 99%.



344 **Figure 5. Curriculum training and empirical capacity.** (a) The value of R used over the course of curriculum training for four different
 345 network sizes. R is incremented once the network achieves an accuracy > 0.99 . Training is considered converged when the value of
 346 R is not incremented for at least 1 million iterations. (b) The final value of R after curriculum training (i.e. network capacity) as a function
 347 of the number of plastic synapses in the network, plotted on a log-log scale. The color of the points corresponds to the number of input

348 units, colors from panel (d). The least-squares fit (black line, slope k , bias c) indicates that the empirical network capacity scales sub-
349 linearly with the number of synapses. (c) Capacity as a function of the input dimension d for various hidden layer sizes N . (d) Capacity
350 as a function of the hidden layer size for various input dimensions d . Capacity primarily depends only on the number of synapses,
351 rather than on the hidden or input layer sizes.
352

353 We curriculum-train networks of different sizes and plot the capacity R_{\max} for each one (Fig 5b).
354 For consistency and ease of training, we restrict the networks to the anti-Hebbian solution and use the
355 uniform readout. We find that the capacity depends primarily on the number of synapses, rather than on
356 the number of pre- or postsynaptic neurons (Fig 5c,d), consistent with previous familiarity detection results
357 (Bogacz and Brown, 2003). To estimate the scaling, we compute a linear least-squares fit of $\log(R_{\max})$ as
358 a function of $\log(Nd)$. Empirically, we find that the capacity of the network scales as

360
$$R_{\max} \approx 0.10(Nd)^{0.79}$$

359 which is sublinear in the number of plastic input synapses to the hidden layer, Nd .

361 The model of Bogacz and Brown (Bogacz and Brown, 2003) for the non-continual task has a
362 capacity that is linear in the number of synapses. To determine whether the difference between the
363 empirical performance of HebbFF and the Bogacz-Brown model reflects a fundamental limitation in the
364 feedforward architecture, we developed an idealized version of the model (Fig 6a) that we could study
365 analytically (*Online methods*).

366 We noted above that the limiting behavior of the network at maximum capacity appears to have
367 W_1 activate just a single unit for memory storage. We build this limiting behavior into the idealized model
368 through a specific choice of W_1 and b_1 , set by design rather than through a training procedure. Specifically,
369 we use the first $n \ll d$ components of $x(t)$ as an identifier by choosing the first n columns of W_1 so that a
370 unique hidden unit is activated by each possible n -bit combination of these components, and set the
371 remaining columns of W_1 to zero. To simplify the model, we do not allow plasticity to operate on the inputs
372 from these bits and set the first n columns of $A(t)$ to zero (Fig 6a). This isolates the "hashing" function of
373 the fixed matrix from the memory storage. Furthermore, instead of a sigmoid nonlinearity for the hidden
374 units, we use a Heaviside step function $\Theta(\cdot)$. Thus, the hidden layer in the idealized model is governed by

375
$$h(t) = \Theta((W_1 + A(t))x(t) + b_1)$$

376 For the nonzero entries of $A(t)$, plasticity is the same as in the trained model. However, because
377 the Heaviside function does not depend on the scale of the input, we can set the plasticity rate to $\eta = -1$
378 without loss of generality. The optimal synaptic decay rate λ can be computed analytically. Finally, a
379 stimulus is considered familiar if all hidden unit activities are identically zero, and novel otherwise (*Online*
380 *methods*).

381 This idealized model exhibits qualitatively similar behavior to HebbFF. We can
382 fit the analytic functional form of the true and false positive probabilities computed from the idealized model
383 (Fig 6b) to the corresponding probabilities of HebbFF (Fig 6c). Furthermore, the histograms of inputs to the
384 hidden layer are qualitatively similar: $W_1x(t) + b_1$ has the same multimodal distribution with more
385 prominent peaks in the middle (Fig 4i, 6d), a bimodal distribution of $A(t)x(t)$ with a large symmetric central
386 peak and a smaller one corresponding to the familiarity signal (Fig 4j, 6e), and a similar distribution of the

387 total input current ($\mathbf{W}_1 + \mathbf{A}(t)\mathbf{x}(t) + \mathbf{b}_1$) (Fig 4o, 6f). From this, we conclude that the memory storage and
388 readout mechanisms are analogous in the meta-learned HebbFF network and the idealized model.

389 Along with the true and false positive probabilities, the memory capacity of the idealized model can
390 be computed analytically (*Online methods*). As in the Bogacz-Brown model (Bogacz and Brown, 2003), the
391 capacity, as characterized by 99% accuracy, is proportional to the number of synapses Nd . There are
392 several possible reasons for the discrepancy between this analytic capacity, as well as that of the Bogacz-
393 Brown model, relative to the empirical capacity for HebbFF.

394 First, the idealized HebbFF model uses a dedicated set of synapses through the fixed \mathbf{W}_1 matrix,
395 and the Bogacz and Brown model selects the units that have the highest input current implicitly through
396 inhibitory competition. Both of these are dedicated addressing functions for the hidden layer, but meta-
397 learned HebbFF must multiplex this functionality with memory storage, leading to correlations between the
398 hidden layer input currents from the plastic and fixed synapse components (Fig S4a).

399 In addition, replacing the logistic function with a Heaviside function means that familiar stimuli truly
400 generate no plasticity in the idealized model, reducing overwriting at the cost of not reinforcing partially-
401 decayed memories (Fig S4b-c). For the same reason, in contrast to HebbFF, the idealized model achieves
402 maximal plasticity for any suprathreshold level of input to a hidden layer unit.

403 Finally, training the HebbFF model may lead to specialized solutions for small d and N that have
404 better performance than that predicted by the asymptotic analysis. Similarly, training may not converge to
405 the optimal solution for large d and N because it requires the use of very long repeat intervals R . This
406 means the dataset size T must be very large to include a sufficient number of familiar examples, which may
407 lead to practical issues such as vanishing gradients. Thus, the empirical capacity may scale sublinearly
408 with the number of synapses because of over-performance at low R , under-performance at high R , or both.

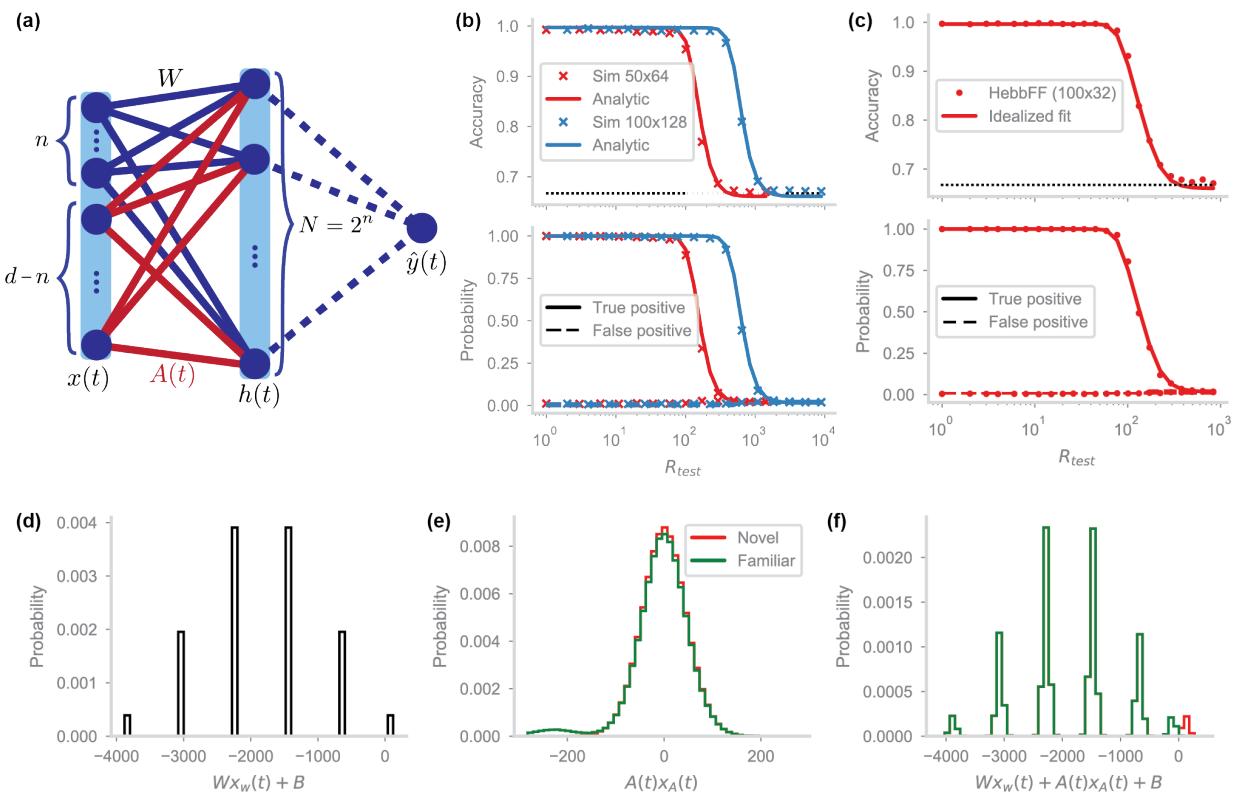


Figure 6. Idealized model. (a) The idealized HebbFF network architecture. In contrast to the original HebbFF network with a single effective matrix $W_1 + A(t)$, the input $x(t)$ is effectively split into two sections of size n and $D = d - n$ that serve as inputs into separate static and plastic synaptic matrices W_1 and $A(t)$, respectively (*Online methods*). The hidden layer size is $N = 2^n$. The readout unit outputs $\hat{y}(t) = 1$ whenever any of the hidden units is active. (b) The analytic calculation of network performance (solid line) matches simulation results for the idealized network (x's), shown for two different network sizes (red, blue). (c) A least-squares fit of the analytic performance curve of the idealized network to a trained HebbFF network of the same size for two network sizes. The idealized network has similar performance to the HebbFF model if its decay rate and bias are scaled appropriately: $\lambda \approx 0.986$, $b_1 \approx -4.771$ (for all units) for $d = 200$, $N = 32$, and $\lambda \approx 0.993$, $b_1 \approx -4.771$ for $d = 200$, $N = 32$. (d-f) Same as Fig 4(l,j,m), but for the idealized network ($D = 400$, $N = 32$, $R = 300$).

HebbFF recapitulates neural data from inferotemporal cortex

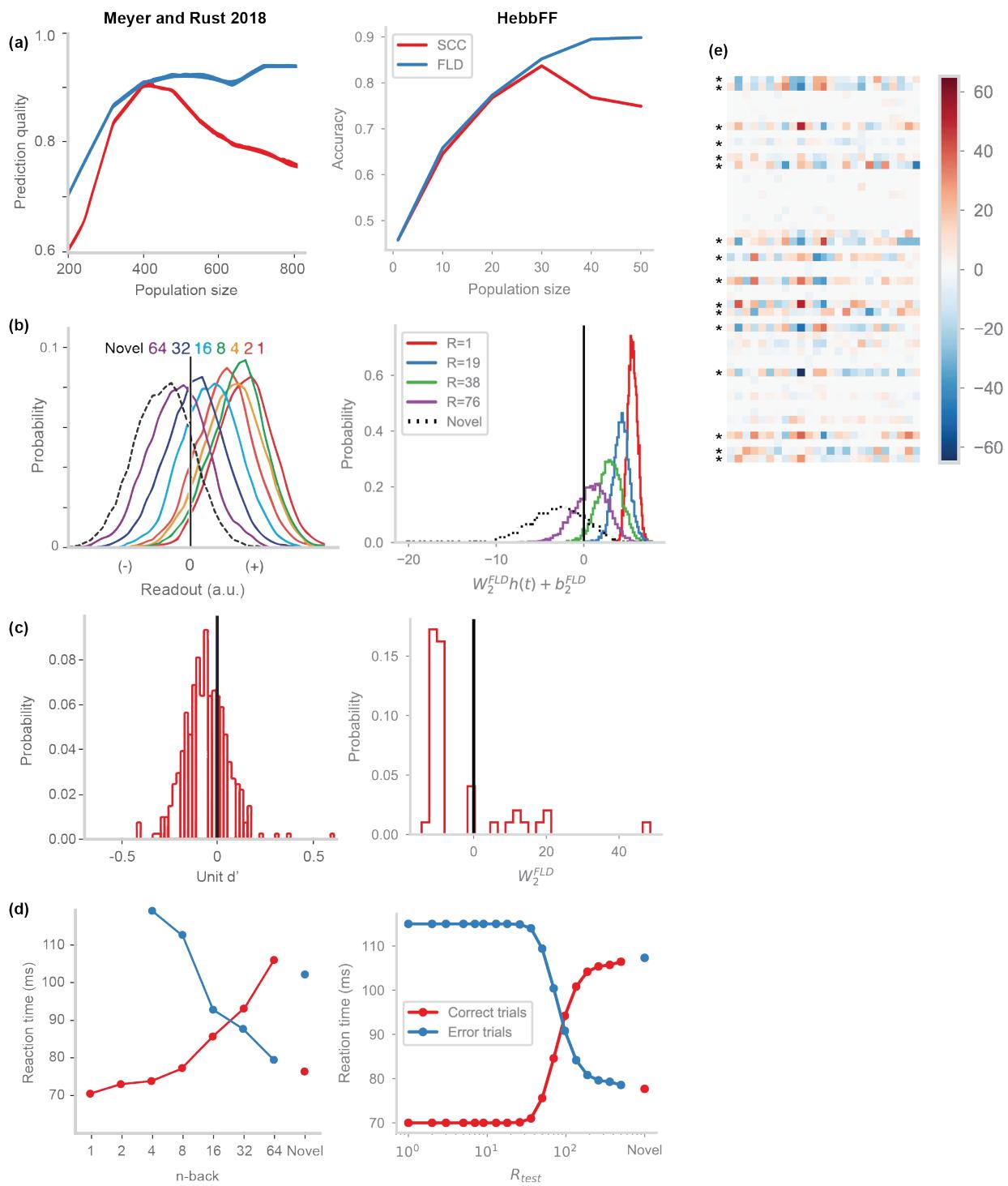
We next compare the optimized HebbFF model with experimental results. Meyer and Rust (Meyer and Rust, 2018) recorded neurons from the inferotemporal (IT) cortex of monkeys performing familiarity detection and compared the quality of two decoders in predicting behavior from neural data as a function of neural population size. The authors considered a "spike count classifier" (SCC) decoder, which amounts to comparing a simple average of neuronal firing rates to a threshold, as well as a Fisher linear discriminant (FLD), which instead considers a weighted average, with weights computed from the data (Meyer and Rust, 2018).

We perform a similar analysis. We first construct an FLD decoder of the hidden unit firing rates and rank the units in reverse order of their FLD readout weights (i.e. units with the most negative weights are top-ranked; *Online methods*). We then consider decoders that use increasingly larger subsets of hidden units, adding them according to their ranking (Meyer and Rust, 2018). As in the experimental data, performance saturates for the FLD and declines for the SCC readout beyond a certain number of decoded units (Fig 7a). This can be explained by the fact that some units do not provide a reliable signals of familiarity

434 – in the network this is due to suboptimal training, and in the IT cortex possibly due to those neurons
435 performing an unrelated task (see next section). Including them hurts performance of the SCC decoder, but
436 since the FLD readout weight for these units is close to zero, they do not alter its familiarity detection
437 performance.

438 Comparing the experimental and model distributions of readout activity shows a qualitatively similar
439 pattern for outputs to novel and familiar stimuli (Fig 7b). Both distributions shift toward smaller values as R
440 increases, as the outputs for familiar stimuli approach those for novel. The fact that the distribution of
441 outputs becomes narrower for HebbFF as R decreases, unlike in the data, may be due to repetition
442 suppression causing hidden units to have near-zero responses for highly familiar (low R) stimuli, thus
443 causing the readout distribution to cluster around its minimal value. On the other hand, biological neurons
444 that exhibit repetition suppression may never be fully silenced – for example if it takes multiple repetitions
445 to achieve maximal familiarity or if neurons are multiplexed with another task that requires a baseline level
446 of activity. Furthermore, as in the data, the distribution of readout weights is biased towards negative values
447 (Fig 7c).

448 Finally, using a similar "strength theory" analysis as in the experimental results (Meyer and Rust,
449 2018; Murdock, 1985), which suggests that reaction times are inversely proportional to the distance of the
450 readout from the threshold, we can qualitatively reproduce the x-shaped reaction time curves seen in the
451 data (Meyer and Rust, 2018). We used the same proportionality constant determined experimentally to
452 compute network "reaction times" (Fig 7d). Overall, we find that the HebbFF model captures a number of
453 features seen in the experimental results.



454
 455 **Figure 7. Comparison to IT cortex data.** (a) Left: neurons from the IT cortex used to predict the behavioral outputs of a monkey
 456 performing continual familiarity detection, decoded using the Fisher linear discriminant (FLD, blue) or spike count classifier (SCC, red).
 457 Right: units from the hidden layer of a trained HebbFF network (trained with an unconstrained readout W_2 rather than uniform) used
 458 to detect familiarity obtained from SCC and FLD decoders. In both cases, the number of neurons/units available to the decoder was
 459 varied, added in order of increasing weight according to the FLD decoder. While the FLD decoder accuracy saturates, the SCC
 460 decoder accuracy peaks and begins to decline as more neurons/units are included in the decoder. (b) Distribution the FLD decoder
 461 output for IT cortex neurons (left) and HebbFF hidden units (right) for novel stimuli, and familiar stimuli at varying delay intervals. In
 462 both cases, the distribution shifts towards lower values as delay interval increases. For HebbFF, the distribution gets narrower for
 463 shorter delay intervals due to saturation in the hidden layer units. (c) Distribution of the FLD decoder weights for decoding IT cortex
 464 data (left) or HebbFF hidden unit activity (right). In both cases, the majority of output weights are negative, with some positive values.

464 (d) Left: measured reaction time as a function of delay interval for correct and error trials (red, blue curves) in monkeys performing the
465 continual familiarity detection task. Black lines indicate reaction times predicted using strength theory analysis. Right: HebbFF
466 predicted reaction times using analogous strength theory analysis, using constants of proportionality from (Meyer and Rust, 2018)
467 (*Online methods*). Both result in a qualitatively similar x-shaped pattern. Plots on the left side of (a-d) adapted from (Meyer and Rust,
468 2018). (e) The weight matrix W_1 of the HebbFF network trained on the augmented task, requiring simultaneous classification and
469 familiarity detection. Hidden units split into "classification" and "familiarity" units, with classification units (marked with asterisks) having
470 very strong input weights to overcome the noise from the plastic $A(t)$ matrix.
471

472 *Two subpopulations emerge in a classification-augmented task*

473 IT cortex encodes object identity as well as familiarity (Lehky and Tanaka, 2016; Lueschow et al.,
474 1994). To match this dual functionality, we augment familiarity detection with object classification. We first
475 create a large pool of random vectors and randomly assign a binary label to each one. We then generate
476 a familiarity detection dataset as before, except that each novel input is drawn from this pool (without
477 replacement) rather than being generated anew. In addition to the scalar readout of familiarity, the network
478 must now report the class of the stimulus through a second binary output. Critically, both outputs are read
479 out from the same hidden layer activity (*Online methods*).

480 The augmented task could be solved by having all the neurons multiplexed to encode both
481 familiarity and object identity. Alternatively, the neurons could split into two subpopulations, one of which
482 detects familiarity and the other classifies objects (Rutishauser et al., 2015). We find that the HebbFF model
483 converges to this second solution, an even split between familiarity and classifier units, as evident from
484 inspecting the W_1 matrix (Fig 7e). Consistent with this, the capacity of the classifier-augmented HebbFF
485 with 50 hidden units ($R_{\max} \approx 13$) is approximately the same as the original network with 25 units ($R_{\max} \approx$
486 14). In accord with this split, SCC decoder performance peaks in the split-task network when half of the
487 top-ranked units are included (Fig S5d) because including units responsible for object identity but not
488 familiarity degrades the familiarity readout. The other similarities to experimental results discussed in the
489 previous section also hold for the task-augmented network (Fig S5).

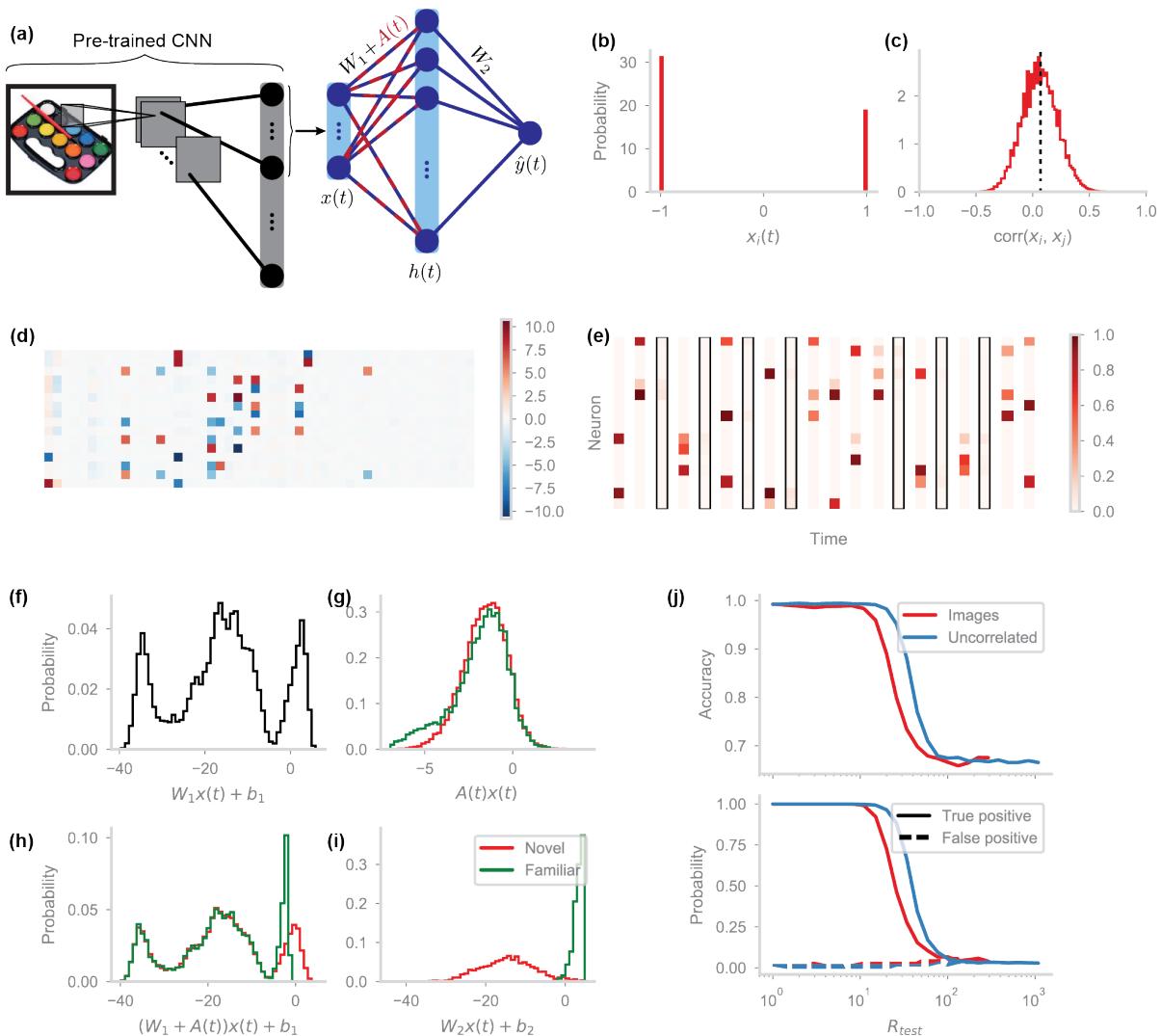


Figure 8. HebbFF performance on real-world images. (a) Network architecture for familiarity detection of real-world images. The activity of the penultimate layer of a convolutional neural network (ResNet18, pre-trained on ImageNet) is downsampled and passed to the HebbFF network ($d = 50, N = 16$) for familiarity detection. Only the HebbFF portion of this network is trained, via curriculum training. (b) Distribution of inputs $x(t)$ to HebbFF. After down-sampling by extracting the first 50 units of the CNN, the activity is centered at zero and binarized. (c) Histogram of the correlations between all pairs of input stimuli $x(t)$. On average (vertical dashed line) the correlation is slightly positive. (d-h) Same plots as Fig 4(f,c,i,j,m), respectively ($R_{\text{train}} = R_{\text{test}} = 12$). (j) Generalization performance, compared to a network of the same size trained on uncorrelated binary random vectors, is lower due to correlations in the input images.

501 Familiarity detection of real images

502 To validate the HebbFF model in a more realistic scenario, we evaluate its performance on real-
 503 world object images. We consider the dataset used by Brady et al. to study familiarity detection in humans
 504 (Brady et al., 2008). As a stand-in for the processing done by the visual stream before the inferotemporal
 505 or perirhinal cortices, we use a pre-trained convolutional neural network (CNN), and extract the activity in
 506 its penultimate layer (before the final classification step). We use the ResNet18 network (He et al., 2015),
 507 although any CNN could, in principle, be used. This activity is a 512-dimensional vector, which, if used as
 508 the HebbFF input dimension d , would lead to the capacity R_{max} being prohibitively large for training

509 purposes. To keep the performance in a reasonable range, we downsample to $d = 50$, either by partial
510 sampling (Fig 8a) or by introducing an intermediate layer (Fig S6a).

511 As the first method of downsampling, we truncate the output of the CNN (Fig 8a). To keep the same
512 input datatype as in previous sections, we also shift the inputs to zero mean and binarize them by taking
513 the sign of each input component (Fig 8b). Unlike in previous sections, however, the inputs to HebbFF now
514 have correlations that tend to be positive (Fig 8c). Nevertheless, this network has qualitatively similar
515 features as the networks trained on uncorrelated vectors. The W_1 matrix has a similar structure (Fig 4f, 8d),
516 the hidden layer activity is sparse (Fig 4c, 8e), and the hidden unit input current distributions have similar
517 shapes (Fig 4i,j,m, 7b, 8f-i). Due to the added correlations, however, there is a decline in performance
518 compared to a network of the same size trained and evaluated on uncorrelated binary random vectors (Fig
519 8j).

520 As another way to downsample, we add a trainable linear layer that transforms the CNN output to
521 a 50-dimensional real-valued vector (Fig S6a). After training, the resulting inputs to HebbFF are no longer
522 binary, but they are zero-mean (Fig S6b) and have zero-mean correlations Fig S6c). Interestingly, the
523 network learns to generate this representation automatically to optimize familiarity detection over long
524 intervals, which further supports storing uncorrelated stimuli. Although the W_1 matrix (Fig S6d) and the
525 distribution of input currents from the fixed component of the synapses (Fig S6f) have a different structure
526 compared to the original network, the operating principle remains the same: the W_1 matrix acts as a hash
527 function to induce sparse activity in the hidden layer (Fig S6e) that is then suppressed for a familiar stimulus
528 through the $A(t)$ matrix (Fig S6g-h). The network maintains its generalization performance across repeat
529 intervals R , and across permutations of the sequence of images (Fig S6j). However, it does not generalize
530 well to images it has not been trained on. It is possible that this difficulty is due to the relatively small number
531 of images used during training and may be addressed by using a much larger dataset such as ImageNet
532 (Deng et al., 2009).

533

534 **Discussion**

535 Continual familiarity detection is a memory task that we perform every day, typically without being
536 aware that we are doing it. We have used meta-learning to generate networks that solve this task using
537 synaptic plasticity. This is distinct from memory storage in RNNs that maintain memory traces through
538 persistent activity. Given the extraordinary capacity and robustness of recognition memory, the idea that
539 biological networks use ongoing activity for this purpose appears untenable (Lundqvist et al., 2018; Masse
540 et al., 2020). If a neuronal network is storing a stimulus by maintaining a particular firing rate pattern across
541 its neurons, any other computation risks disrupting that memory trace. In contrast, storage through synaptic
542 updates leaves the neuronal activity free to perform other computations unrelated to memory storage (Ba
543 et al., 2016). In addition, we found that synaptic plasticity provides a better inductive bias than recurrence
544 for familiarity detection. After optimization through meta-learning, the HebbFF network not only outperforms
545 RNNs on the task, but also easily generalizes both in-distribution and out-of-distribution of the training data.

546 Although RNNs are a common approach to tasks that require storage of the input history (Elman, 1991;
547 Hochreiter and Schmidhuber, 1997; Mante et al., 2013), this result highlights the importance of considering
548 alternative architectures and storage mechanisms, such as those that rely on synaptic, rather than
549 neuronal, dynamics.

550 We found that anti-Hebbian plasticity, in which neuronal co-activation causes synaptic depression,
551 is a better storage mechanism for familiarity detection than Hebbian plasticity. An anti-Hebbian rule
552 generalizes better, has a larger capacity, and is discovered by meta-learning more frequently and reliably.
553 Although this result is consistent with previous work (Bogacz and Brown, 2003), the underlying reasons are
554 different. Bogacz and Brown showed that in a non-continual version of the familiarity detection task, an anti-
555 Hebbian plasticity rule leads to a larger storage capacity, although this advantage only held in the case of
556 correlated inputs. In their case, the anti-Hebbian rule automatically suppresses common input features,
557 effectively storing only the uncorrelated components, leading to an increased capacity. In contrast, anti-
558 Hebbian HebbFF shows an advantage even for uncorrelated inputs in the continual task. This is due to an
559 effective decrease in the number of plasticity events – a synaptic update is weak for a familiar stimulus
560 because the postsynaptic activity is low, leading to smaller updates that are less disruptive to stored
561 memories.

562 In addition to reproducing prior experimental results such as repetition suppression, the HebbFF
563 model allows us to make a novel experimental prediction. Although it is obvious that the true positive rate
564 (probability of correctly identifying a repeated stimulus as "familiar") should decrease with longer delay
565 intervals R , we also observe that the false positive rate slightly increases. Because anti-Hebbian plasticity
566 causes repetition suppression, false negative responses arising from spurious activity in the hidden layer
567 cause additional plasticity, depressing a subset of synapses. Subsequently, a novel stimulus is more likely
568 to silence the hidden layer units, signaling familiarity and resulting in a false positive response. Prior
569 experimental paradigms have not measured this effect because each trial had familiar stimuli interleaved
570 at various delay intervals. As a result, novel stimuli could not be separated and scored depending on the difficulty
571 of the dataset, and false positive probability was reported in aggregate. If biological networks implement
572 familiarity detection through an anti-Hebbian plasticity mechanism, we expect an increase in the false
573 positive rate to coincide with a decrease in the true positive rate.

574 There are experimental results that the HebbFF model does not capture. For example, data from
575 human subjects shows a very slow decrease in performance as a function of R that begins at relatively
576 small value (Brady et al., 2008). In contrast, HebbFF has near-perfect performance for all $R < R_{\max}$, and
577 then performance drops off quickly. However, it is likely that errors in the experiments do not reflect
578 limitations on recognition memory but rather are due to factors such as fatigue and lack of attention that
579 were not included in the model.

580 Finally, our work demonstrates the utility of meta-learning as a tool for neuroscience discovery. We
581 used meta-learning to optimize a network architecture and plasticity rule that solves the continual familiarity
582 detection task, contrasted it with an alternative sub-optimal solution, and subsequently used analytic

583 methods to understand its mechanism. A similar approach can be used for other networks, plasticity rules,
584 datasets, and tasks.

585

586 **Acknowledgement**

587 We thank Stefano Fusi, Ken Miller, Dmitriy Aronov, James Murray, Marcus Benna, SueYeon Chung, Juri
588 Minxha, Taiga Abe, and Denis Turcu for helpful discussions. Research supported by NSF NeuroNex Award
589 DBI-1707398, the Gatsby Charitable Foundation, and the Simons Collaboration for the Global Brain. G.R.Y.
590 was additionally supported by the Simons Foundation. We acknowledge computing resources from
591 Columbia University's Shared Research Computing Facility project, which is supported by NIH Research
592 Facility Improvement Grant 1G20RR030893-01, and associated funds from the New York State Empire
593 State Development, Division of Science Technology and Innovation (NYSTAR) Contract C090171, both
594 awarded April 15, 2010.

595

596 **References**

- 597 Abbott, L.F., and Nelson, S.B. (2000). Synaptic plasticity: taming the beast. *Nature Neuroscience* **3**, 1178–
598 1183.
- 599 Androulidakis, Z., Lulham, A., Bogacz, R., and Brown, M.W. (2008). Computational models can replicate
600 the capacity of human recognition memory. *Network: Computation in Neural Systems* **19**, 161–182.
- 601 Ba, J., Hinton, G., Mnih, V., Leibo, J.Z., and Ionescu, C. (2016). Using Fast Weights to Attend to the Recent
602 Past. *ArXiv:1610.06258 [Cs, Stat]*.
- 603 Bengio, Y., Bengio, S., and Cloutier, J. (1991). Learning a synaptic learning rule. p.
- 604 Benna, M.K., and Fusi, S. (2016). Computational principles of synaptic memory consolidation. *Nature
605 Neuroscience* **19**, 1697–1706.
- 606 Bogacz, R., and Brown, M.W. (2003). Comparison of computational models of familiarity discrimination in
607 the perirhinal cortex. *Hippocampus* **13**, 494–524.
- 608 Brady, T.F., Konkle, T., Alvarez, G.A., and Oliva, A. (2008). Visual long-term memory has a massive storage
609 capacity for object details. *PNAS* **105**, 14325–14329.
- 610 Confavreux, B., Agnes, E.J., Zenke, F., Lillicrap, T., and Vogels, T.P. (2020). A meta-learning approach to
611 (re)discover plasticity rules that carve a desired function into a neural network. *BioRxiv* 2020.10.24.353409.
- 612 Deng, J., Dong, W., Socher, R., Li, L., Kai Li, and Li Fei-Fei (2009). ImageNet: A large-scale hierarchical
613 image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pp. 248–255.
- 614 Elman, J.L. (1991). Distributed representations, simple recurrent networks, and grammatical structure.
615 *Mach Learn* **7**, 195–225.
- 616 Frémaux, N., and Gerstner, W. (2016). Neuromodulated Spike-Timing-Dependent Plasticity, and Theory of
617 Three-Factor Learning Rules. *Front. Neural Circuits* **9**.
- 618 Fusi, S., Drew, P.J., and Abbott, L.F. (2005). Cascade Models of Synaptically Stored Memories. *Neuron*
619 **45**, 599–611.

- 620 Gerstner, W., Lehmann, M., Liakoni, V., Corneil, D., and Brea, J. (2018). Eligibility Traces and Plasticity on
621 Behavioral Time Scales: Experimental Support of NeoHebbian Three-Factor Learning Rules. *Front. Neural*
622 *Circuits* 12.
- 623 Graves, A., Wayne, G., and Danihelka, I. (2014). Neural Turing Machines. ArXiv:1410.5401 [Cs].
- 624 Graves, A., Wayne, G., Reynolds, M., Harley, T., Danihelka, I., Grabska-Barwińska, A., Colmenarejo, S.G.,
625 Grefenstette, E., Ramalho, T., Agapiou, J., et al. (2016). Hybrid computing using a neural network with
626 dynamic external memory. *Nature* 538, 471–476.
- 627 Grill-Spector, K., Henson, R., and Martin, A. (2006). Repetition and the brain: neural models of stimulus-
628 specific effects. *Trends in Cognitive Sciences* 10, 14–23.
- 629 Hasson, U., Chen, J., and Honey, C.J. (2015). Hierarchical process memory: memory as an integral
630 component of information processing. *Trends in Cognitive Sciences* 19, 304–313.
- 631 He, K., Zhang, X., Ren, S., and Sun, J. (2015). Deep Residual Learning for Image Recognition.
632 ArXiv:1512.03385 [Cs].
- 633 Hebb, D.O. (1949). *The organization of behavior* (New York: Wiley).
- 634 Hochreiter, S., and Schmidhuber, J. (1997). Long Short-Term Memory. *Neural Computation* 9, 1735–1780.
- 635 Hopfield, J.J. (1982). Neural networks and physical systems with emergent collective computational
636 abilities. *PNAS* 79, 2554–2558.
- 637 Ji-An, L., Stefanini, F., Benna, M.K., and Fusi, S. (2019). Face familiarity detection with complex synapses.
638 BioRxiv 854059.
- 639 Kingma, D.P., and Ba, J. (2017). Adam: A Method for Stochastic Optimization. ArXiv:1412.6980 [Cs].
- 640 Kirkpatrick, J., Pascanu, R., Rabinowitz, N., Veness, J., Desjardins, G., Rusu, A.A., Milan, K., Quan, J.,
641 Ramalho, T., Grabska-Barwinska, A., et al. (2017). Overcoming catastrophic forgetting in neural networks.
642 *PNAS* 114, 3521–3526.
- 643 LeCun, Y., Bengio, Y., and Hinton, G. (2015). Deep learning. *Nature* 521, 436–444.
- 644 Lehky, S.R., and Tanaka, K. (2016). Neural representation for object recognition in inferotemporal cortex.
645 *Current Opinion in Neurobiology* 37, 23–35.
- 646 Lim, S., McKee, J.L., Woloszyn, L., Amit, Y., Freedman, D.J., Sheinberg, D.L., and Brunel, N. (2015).
647 Inferring learning rules from distributions of firing rates in cortical neurons. *Nature Neuroscience* 18, 1804–
648 1810.
- 649 Lueschow, A., Miller, E.K., and Desimone, R. (1994). Inferior Temporal Mechanisms for Invariant Object
650 Recognition. *Cerebral Cortex* 4, 523–531.
- 651 Lundqvist, M., Herman, P., and Miller, E.K. (2018). Working Memory: Delay Activity, Yes! Persistent
652 Activity? Maybe Not. *J. Neurosci.* 38, 7013–7019.
- 653 Mante, V., Sussillo, D., Shenoy, K.V., and Newsome, W.T. (2013). Context-dependent computation by
654 recurrent dynamics in prefrontal cortex. *Nature* 503, 78–84.
- 655 Masse, N.Y., Rosen, M.C., and Freedman, D.J. (2020). Reevaluating the Role of Persistent Neural Activity
656 in Short-Term Memory. *Trends in Cognitive Sciences* 24, 242–258.

- 657 Meyer, T., and Rust, N.C. (2018). Single-exposure visual memory judgments are reflected in inferotemporal
658 cortex. *ELife* 7.
- 659 Miconi, T., Clune, J., and Stanley, K.O. (2018). Differentiable plasticity: training plastic neural networks with
660 backpropagation. ArXiv:1804.02464 [Cs, Stat].
- 661 Miconi, T., Rawal, A., Clune, J., and Stanley, K.O. (2019). Backpropamine: training self-modifying neural
662 networks with differentiable neuromodulated plasticity. 15.
- 663 Miller, E.K., Li, L., and Desimone, R. (1991). A Neural Mechanism for Working and Recognition Memory in
664 Inferior Temporal Cortex. *Science* 254, 1377–1379.
- 665 Murdock, B.B. (1985). An analysis of the strength-latency relationship. *Memory & Cognition* 13, 511–521.
- 666 Norman, K.A., and O'Reilly, R.C. (2003). Modeling hippocampal and neocortical contributions to recognition
667 memory: A complementary-learning-systems approach. *Psychological Review* 110, 611–646.
- 668 Parisi, G. (1986). A memory which forgets. *J. Phys. A: Math. Gen.* 19, L617–L620.
- 669 Rumelhart, D.E., Hinton, G.E., and Williams, R.J. (1986). Learning representations by back-propagating
670 errors. *Nature* 323, 533–536.
- 671 Rutishauser, U., Ye, S., Koroma, M., Tudusciuc, O., Ross, I.B., Chung, J.M., and Mamelak, A.N. (2015).
672 Representation of retrieval confidence by single neurons in the human medial temporal lobe. *Nature
673 Neuroscience* 18, 1041–1050.
- 674 Sohal, V.S., and Hasselmo, M.E. (2000). A model for experience-dependent changes in the responses of
675 inferotemporal neurons. *Network: Computation in Neural Systems* 11, 169–190.
- 676 Standing, L. (1973). Learning 10000 pictures. *Quarterly Journal of Experimental Psychology* 25, 207–222.
- 677 Thrun, S., and Pratt, L. (2012). *Learning to Learn* (Springer Science & Business Media).
- 678 Xiang, J.-Z., and Brown, M.W. (1998). Differential neuronal encoding of novelty, familiarity and recency in
679 regions of the anterior temporal lobe. *Neuropharmacology* 37, 657–676.
- 680 Zenke, F., Poole, B., and Ganguli, S. (2017). Continual Learning Through Synaptic Intelligence.
681 ArXiv:1703.04200 [Cs, q-Bio, Stat].
- 682

683 **Online methods**

684 *HebbFF and RNN training*

685 To set the fixed HebbFF parameters $\mathbf{W}_1, \mathbf{b}_1, \mathbf{W}_2, \mathbf{b}_2, \lambda, \eta$, as well as the RNN weight and bias
686 matrices, we use the PyTorch implementation of the Adam optimizer with the suggested default
687 hyperparameters (Kingma and Ba, 2017). For a single trial, we use a dataset containing T stimuli, with
688 familiar ones appearing at a repeat interval R . We present stimuli to the network sequentially, and compute
689 the binary cross-entropy loss

$$692 \quad \mathcal{L} = \frac{1}{T} \sum_{t=1}^T y(t) \log \hat{y}(t) + (1 - y(t)) \log(1 - \hat{y}(t))$$

690 Since this is a dynamic task (the state of the network at time $t+1$ depends on the state at time t),
691 backpropagation through time is used to compute the gradient of the loss with respect to the parameters.

693 For each trial, we either use the same pre-generated length- T dataset, or we generate a new length-
694 T dataset using the same repeat interval R . We refer to the latter case as the "infinite data" training regime
695 since the sample space is much larger than the network would explore during training. Note that in the
696 infinite data regime, we do not consider a validation dataset, since the training set is new every time and
697 the training accuracy is therefore the same as the validation accuracy. In both cases, one trial corresponds
698 to one step of gradient descent. To train the HebbFF network, the plastic matrix $\mathbf{A}(t)$ is reset to a matrix of
699 zeros at the start of each trial. Similarly, when training the RNN, hidden unit activity is reset to zero. In
700 practice, the plastic matrix of HebbFF reaches its steady state distribution quickly and the transient does
701 not contribute significantly to the gradient, so any reasonable initialization can be used.

702 To train the HebbFF network on the augmented familiarity detection/object classification task, we
703 simply sum the cross-entropy losses from the classifier and familiarity output units:

$$704 \quad \mathcal{L} = \frac{1}{2T} \sum_{t=1}^T \sum_{a=1}^2 y_a(t) \log \hat{y}_a(t) + (1 - y_a(t)) \log(1 - \hat{y}_a(t))$$

705 For every trial, we draw a new dataset from the pre-generated pool of stimuli. The class of each stimulus
706 remains the same across datasets, but the ordering and repeats are chosen randomly each time. Although
707 the network will have seen all of the stimuli during training in order to learn their classes, we can test
708 generalization performance on the familiarity subtask by varying R and generating previously unseen
709 permutations of the stimuli.

710 The PyTorch implementation of each of these can be found at <https://github.com/dtyulman/hebbff>.

711

712 *Bogacz-Brown (Bogacz and Brown, 2003) model implementation*

713 To validate it on the (non-continual) two-alternative-forced-choice (2AFC) familiarity detection task,
714 we implement the anti-Hebbian model as described by Bogacz and Brown (Bogacz and Brown, 2003), with
715 the exception that the distribution of weights in the plastic weight matrix must be normalized such that its
716 variance is equal to $\frac{1}{N}$, rather than unit variance as stated in the paper. In the encoding phase, the network

717 is presented a sequence of P random patterns. In the testing phase, it is shown the original P patterns, as
718 well as P novel ones. Critically, there are no plastic updates in the testing phase. A stimulus is reported as
719 "familiar" if the output unit activity is below the mean across all $2P$ test patterns and "novel" otherwise. We
720 see that this model performs well on the 2AFC task with a range of plasticity rates η (Fig S2a), so we
721 arbitrarily choose $\eta = 0.7$ to test its performance on the continual task.

722 The continual task, unlike the 2AFC task, does not have an equal proportion of novel and familiar
723 stimuli since we ensure that a stimulus is repeated at most once. So, we set the readout threshold such
724 that an item is considered novel if it is in the f^{th} quantile of output unit activity for that trial, where f is the
725 fraction of novel stimuli in the trial. This ensures that the fraction of stimuli reported as "novel" is equal to
726 the true fraction of novel stimuli. In the case of equal proportions of novel and familiar stimuli, this reduces
727 to the threshold being equal to the mean of the output unit activity for that trial.

728 Finally, note that unlike in the 2AFC task (Fig S2a), the performance of this model does not go to
729 chance levels for large dataset sizes T in the continual task (Fig 3d). Rather, the true positive rate goes to
730 zero and the false positive rate is ≈ 0.5 , so accuracy is ≈ 0.33 . The reason for this difference is that the
731 second presentation of a stimulus in the continual task causes an additional plasticity event, unlike the
732 2AFC task where the test phase is offline. As a result, for datasets much larger than the network capacity
733 $T \gg P^*$, the output unit activity for familiar stimuli becomes larger than the activity for novel stimuli (Fig S2b).

734

735 *Training FLD and SCC decoders*

736 To construct the Fisher linear discriminant (FLD) and spike count classifier (SCC) decoders, we
737 first generate a dataset of length $T = 1000$. To better match the experimental dataset (Meyer and Rust,
738 2018), we use multiple values of R in this single stream. For each familiar stimulus, the value of R is drawn
739 uniformly at random from 34 unique values, log-spaced from 1 to 100 (in practice, the results are
740 qualitatively the same regardless of the number of items, the range, or whether the spacing is linear or
741 logarithmic). We evaluate the trained network on this dataset and use the firing rates of the hidden layer to
742 perform analyses analogous to those reported in (Meyer and Rust, 2018).

743 We compute the readout weight and bias terms for the FLD decoder as

744

$$\mathbf{W}_2^{\text{FLD}} = \Sigma^{-1}(\bar{\mathbf{h}}_{\text{nov}} - \bar{\mathbf{h}}_{\text{fam}}), b_2^{\text{FLD}} = -\mathbf{W}_2^{\text{FLD}} \cdot \frac{1}{2}(\bar{\mathbf{h}}_{\text{nov}} + \bar{\mathbf{h}}_{\text{fam}})$$

745 where $\bar{\mathbf{h}}_{\text{nov}}$ and $\bar{\mathbf{h}}_{\text{fam}}$ are the average firing rates of the hidden layer for novel and familiar stimuli,
746 respectively, and the mean covariance matrix is calculated as

747

$$\Sigma = \frac{\Sigma_{\text{fam}} + \Sigma_{\text{nov}}}{2}$$

748 where Σ_{fam} and Σ_{nov} are the covariance matrices of the firing rates of the hidden layer for familiar and novel
749 stimuli, respectively. The SCC decoder is a simple weighted average

750

$$\mathbf{W}_2^{\text{SCC}} = \frac{1}{N}(\bar{\mathbf{h}}_{\text{nov}} - \bar{\mathbf{h}}_{\text{fam}}), b_2^{\text{SCC}} = -\mathbf{W}_2^{\text{SCC}} \cdot \frac{1}{2}(\bar{\mathbf{h}}_{\text{nov}} + \bar{\mathbf{h}}_{\text{fam}})$$

751 To get the ranking of the units for both decoders, we sort their readout weights and consider the
752 most negative weights as the highest ranked. Note that for both decoders, the sign of the weights is flipped
753 compared to (Meyer and Rust, 2018), and high-ranked units have the most negative weights rather than
754 positive. This is due to the fact that we ask the network to label familiar stimuli as $y(t) = 1$, whereas (Meyer
755 and Rust, 2018) readout a familiar stimulus as $y(t) = 0$. The two cases are symmetric and this does not
756 change the results.

757

758 *Idealized model analytic capacity derivation*

759 For notational simplicity, we only consider the nonzero submatrices of \mathbf{W}_1 and $\mathbf{A}(t)$, each of which
760 acts on its corresponding subset of the input vector $\mathbf{x}(t)$. Thus, equivalently, input layer of the idealized
761 network is a d -dimensional vector split into two parts $\mathbf{x}(t) = [\mathbf{x}_W(t), \mathbf{x}_A(t)]$, of dimension n and D
762 respectively ($d = n + D$). Thus, the firing rate of the hidden layer is given by

763
$$\mathbf{h}(t) = \Theta(\mathbf{W}_1 \mathbf{x}_W(t) + \mathbf{A}(t) \mathbf{x}_A(t) + \mathbf{b}_1)$$

764 for an $N \times n$ matrix \mathbf{W}_1 , an $N \times D$ matrix $\mathbf{A}(t)$, and an $N \times 1$ vector \mathbf{b}_1 . In other words, the firing rate of the i^{th}
765 hidden unit is

$$h_i(t) = \Theta\left(\sum_{j=1}^n W_{ij}x_j(t) + \sum_{k=1}^D A_{ik}(t)x_{n+k}(t) + b\right) \quad (1)$$

766 for $i = 1, \dots, N$, where $\Theta(\cdot)$ is the Heaviside step function, i.e. $\Theta(z) = 0$ for $z < 0$ and 1 for $z \geq 0$. We fix the
767 value of b to be the same for all i . As before, the elements of $\mathbf{x}(t)$ are $+1$ or -1 with equal probability. We
768 would like to specify the network parameters such that exactly one hidden neuron is active for a novel
769 stimulus and none for familiar, which will serve as the familiarity readout mechanism.

770 The $N \times n$ matrix \mathbf{W}_1 is designed such that the vector $\mathbf{W}_1 \mathbf{x}_W(t)$ has exactly one maximal entry given
771 any such $\mathbf{x}(t)$. Importantly, this matrix must act like a hash function such that different values of $\mathbf{x}_W(t)$ result
772 in different entries of $\mathbf{W}_1 \mathbf{x}_W(t)$ attaining the maximum value. One such \mathbf{W}_1 is one whose rows enumerate
773 all of the binary length- n strings consisting of entries $+1$ and -1 . This sets the number of rows N to be
774 equal to the total number of such strings, $N = 2^n$. To set the overall scale of the input current (the term
775 inside the nonlinearity), we scale this matrix by a factor K , to be determined later. For example, if $n = 3$,

776

$$\mathbf{W}_1 = K \begin{bmatrix} +1 & +1 & +1 \\ +1 & +1 & -1 \\ +1 & -1 & +1 \\ +1 & -1 & -1 \\ -1 & +1 & +1 \\ -1 & +1 & -1 \\ -1 & -1 & +1 \\ -1 & -1 & -1 \end{bmatrix}$$

777 Thus, we have $\sum_{j=1}^n W_{ij}x_j(t) = Kn$ for exactly one value of i , specifically the row where $W_{ij} = x_j(t)$

778 for all j . This is the unique maximal value and will correspond to a different row for each instance of $\mathbf{x}_W(t)$.

779 Subsequently, $\sum_{j=1}^n W_{ij}x_j(t) = K(n - 2)$ for n values of i , specifically those where $W_{ij} \neq x_j(t)$ for exactly

780 one j , and so on. Assuming that the vector $\mathbf{A}(t)\mathbf{x}_A(t)$ is zero-mean with sufficiently small variance (this will

781 be made rigorous shortly), we can now choose the scalar offset b such that exactly one element of $\mathbf{h}(t)$ is

782 equal to 1 and all others are zero.

783 The $N \times D$ matrix $\mathbf{A}(t)$ is updated at every timestep by $\mathbf{A}(t + 1) = \lambda\mathbf{A}(t) - \eta\mathbf{h}(t)\mathbf{x}_A(t)^T$, where the
784 plasticity rate η is now restricted to be positive, corresponding to an anti-Hebbian learning rule. Considering
785 one entry in this matrix and unrolling this recurrence, we find that

$$\begin{aligned} 788 \quad A_{ik}(t + 1) &= \lambda A_{ik}(t) - \eta h_i(t)x_{n+k}(t) \\ 789 \quad &= \lambda^{t+1}A_{ik}(0) - \eta \sum_{t'=0}^t \lambda^{t-t'}h_i(t')x_{n+k}(t') \\ 790 \quad &= -\eta \sum_{t'=0}^t \lambda^{t-t'}h_i(t')x_{n+k}(t') \end{aligned}$$

786 where the last equality holds if we assume that the network is in steady-state, so t is large, i.e. $t \rightarrow \infty$, and
787 therefore $\lambda^{t+1}A_{ik}(0) \rightarrow 0$.

788 We can now consider the middle term of eq. 1, which we denote by $\varepsilon_i(t)$. We consider it as a
789 random variable and compute its mean and variance. By definition, we have

$$\begin{aligned} 793 \quad \varepsilon_i(t) &= \sum_{k=1}^D A_{ik}(t)x_{n+k}(t) \\ 794 \quad &= \sum_{k=1}^D \left(-\eta \sum_{t'=0}^{t-1} \lambda^{t-1-t'}h_i(t')x_{n+k}(t') \right) x_{n+k}(t) \\ &= -\eta \sum_{t'=0}^{t-1} \lambda^{t-1-t'}h_i(t') \sum_{k=1}^D x_{n+k}(t') x_{n+k}(t) \end{aligned} \tag{2}$$

795 In the case where $\mathbf{x}(t)$ is novel, $x_{n+k}(t')$ and $x_{n+k}(t)$ are independent Bernoulli random variables
796 that take on values ± 1 with probability $1/2$. Thus, $X_k(t') = x_{n+k}(t')x_{n+k}(t)$ is also a Bernoulli random variable
797 with the same distribution, zero mean and unit variance, so

$$798 \quad \varepsilon_i(t) = -\eta \sum_{t'=0}^{t-1} \lambda^{t-1-t'}h_i(t') \sum_{k=1}^D X_k(t')$$

799 Since the entries of $\mathbf{x}(t)$ are independent by definition, the $X_k(t')$ are also independent across k , so
800 summing over these indices, the variances add. Therefore, $X(t') = \sum_{k=1}^D X_k(t')$ is a random variable with
801 mean 0 and variance D , and

$$802 \quad \varepsilon_i(t) = -\eta \sum_{t'=0}^{t-1} \lambda^{t-1-t'}h_i(t')X(t')$$

803 Next, we need the statistics of the term $h_i(t')$. Since it is a function of the random variable $x(t)$, we
 804 also consider it as a random variable. Let f_{eff} denote the fraction of stimuli reported as "novel" by the
 805 network. Note that there are two ways for a network to report a stimulus as "novel" – by correctly identifying
 806 a novel stimulus ("true negative"), or incorrectly identifying a familiar one ("false negative") – so if we let f
 807 denote the true fraction of novel input stimuli, we have

$$808 \quad f_{\text{eff}} = P_{TN}f + P_{FN}(1-f) = (1-P_{FP})f + (1-P_{TP})(1-f)$$

809 where P_{TN} , P_{FN} , P_{TP} and P_{FP} are the true negative, false negative, true positive, and false positive rates,
 810 respectively. Since by design there is exactly one hidden unit active for a novel stimulus, we have $h_i(t') =$
 811 1 with probability $\frac{f_{\text{eff}}}{N}$, and $h_i(t') = 0$ with probability $1 - \frac{f_{\text{eff}}}{N}$. So, $h_i(t')$ is a Bernoulli random variable with
 812 mean $\frac{f_{\text{eff}}}{N}$ and variance $\frac{f_{\text{eff}}}{N}\left(1 - \frac{f_{\text{eff}}}{N}\right)$. Now, we let $H_i(t') = h_i(t')X(t')$, so

$$813 \quad \varepsilon_i(t) = -\eta \sum_{t'=0}^{t-1} \lambda^{t-1-t'} H_i(t')$$

814 Although $h_i(t')$ is, in principle, a function of $x(t')$, we assume they are independent. Since $X(t')$ is
 815 zero-mean, the mean of $H_i(t')$ is also zero. Using the identity $\text{var}(XY) = \text{var}(X)\text{var}(Y) + \text{var}(X)\mathbb{E}^2[Y] +$
 816 $\text{var}(Y)\mathbb{E}^2[X]$, which holds for independent random variables X and Y , we have that the variance of $H_i(t')$
 817 is $\frac{f_{\text{eff}}D}{N}$. Finally, for convenience we can rewrite this as

$$818 \quad \varepsilon_i(t) = -\eta \sqrt{\frac{f_{\text{eff}}D}{N}} \sum_{t'=0}^{t-1} \lambda^{t-1-t'} \xi_i(t')$$

819 where $\xi_i(t')$ is a zero-mean, unit-variance random variable. Furthermore, we now see that by the Central
 820 Limit Theorem $\varepsilon_i(t)$ is a Gaussian random variable since we are considering the steady-state performance
 821 at large t , so we can take $t \rightarrow \infty$.

822 We can now compute the mean and variance of $\varepsilon_i(t)$. First, since $x_{n+k}(t)$ is zero-mean and
 823 independent of $A_{ik}(t)$,

$$824 \quad \mathbb{E}[\varepsilon_i(t)] = \mathbb{E}\left[\sum_{k=1}^D A_{ik}(t)x_{n+k}(t)\right] = 0$$

825 To compute the variance,

$$\begin{aligned} 826 \quad \text{var}(\varepsilon_i(t)) &= \mathbb{E}[\varepsilon_i^2(t)] - \mathbb{E}^2[\varepsilon_i(t)] = \mathbb{E}[\varepsilon_i^2(t)] \\ 827 &= \mathbb{E}\left[\left(-\eta \sqrt{\frac{f_{\text{eff}}D}{N}} \sum_{t'=0}^{t-1} \lambda^{t-1-t'} \xi_i(t')\right)^2\right] \\ 828 &= \mathbb{E}\left[\left(-\eta \sqrt{\frac{f_{\text{eff}}D}{N}} \sum_{t'=0}^{t-1} \lambda^{t-1-t'} \xi_i(t')\right)\left(-\eta \sqrt{\frac{f_{\text{eff}}D}{N}} \sum_{t''=0}^{t-1} \lambda^{t-1-t''} \xi_i(t'')\right)\right] \\ 829 &= \eta^2 \frac{f_{\text{eff}}D}{N} \sum_{t'=0}^{t-1} \sum_{t''=0}^{t-1} \lambda^{t-1-t'} \lambda^{t-1-t''} \mathbb{E}[\xi_i(t') \xi_i(t'')] \end{aligned}$$

830 In general, we have $\mathbb{E}[\xi_i(t')\xi_i(t'')] = 1$ for $t' = t''$ since $\xi_i(t')$ is a zero-mean, unit-variance random
 831 variable. For $t' \neq t''$, we again make a simplifying independence assumption. In principle, $\xi_i(t'')$ is not
 832 independent of $\xi_i(t')$ since $h_i(t'')$ depends on $h_i(t')$ for $t'' > t'$ through the memory stored in the $A(t)$
 833 matrix. This dependence, however, is sufficiently weak, so we let $\mathbb{E}[\xi_i(t')\xi_i(t'')] = 0$ for $t' \neq t''$. As a result,
 834 the double-sum collapses and we have

$$\begin{aligned} 835 \quad \text{var}(\varepsilon_i(t)) &= \eta^2 \frac{f_{\text{eff}}D}{N} \sum_{t'=0}^{t-1} \lambda^{2(t-1-t')} \\ 836 \quad &= \eta^2 \frac{f_{\text{eff}}D}{N} \frac{1 - \lambda^{2t}}{1 - \lambda^2} \end{aligned}$$

837 where the second equality comes from the standard geometric series. As before, since we are considering
 838 the steady-state with $t \rightarrow \infty$, we have $\gamma^{2t} \rightarrow 0$, so

$$839 \quad \text{var}(\varepsilon_i(t)) = \eta^2 \frac{f_{\text{eff}}D}{N} \frac{1}{1 - \lambda^2}$$

840 Thus, for a novel input $x(t)$ we can write

$$\varepsilon_i(t) = \xi \eta \sqrt{\frac{f_{\text{eff}}D}{N(1 - \lambda^2)}} \tag{3}$$

841 for all i , where ξ is a zero-mean, unit-variance Gaussian random variable, since $\varepsilon_i(t)$ is Gaussian.

842 For a familiar stimulus, where $x(t) = x(t - R)$, clearly $x_{n+k}(t')$ and $x_{n+k}(t)$ are no longer
 843 independent for $t' = t - R$. Thus, we consider this term separately, rewriting the sum in eq. 2 as

$$844 \quad \varepsilon_i(t) = -\eta \lambda^{t-1-(t-R)} h_i(t-R) \sum_{k=1}^D x_{n+k}(t-R) x_{n+k}(t) - \eta \sum_{\substack{t'=0 \\ t' \neq t-R}}^{t-1} \lambda^{t-1-t'} h_i(t') \sum_{k=1}^D x_{n+k}(t') x_{n+k}(t)$$

845 Assuming no errors, by design, $h_i(t-R) = 1$ for exactly one neuron i , since the stimulus at time $t - R$ was
 846 guaranteed to be novel (we enforce that a stimulus is repeated at most once in this task). We consider the
 847 statistics of $\varepsilon_i(t)$ for this particular neuron. In the first term, the sum $\sum_{k=1}^D x_{n+k}(t-R) x_{n+k}(t) = D$ since by
 848 assumption $x_{n+k}(t) = x_{n+k}(t - R)$ for all k . The second term has the same distribution as the one for a
 849 novel input since we have only removed one term from the sum and t is large. Thus, for a familiar stimulus
 850 we can write

$$851 \quad \varepsilon_i(t) = -\eta \lambda^{R-1} D + \xi \eta \sqrt{\frac{f_{\text{eff}}D}{N(1 - \lambda^2)}}$$

852 for exactly one value of i , where ξ is a zero-mean, unit-variance random variable as before. For all other
 853 values of i , eq. 3 holds.

854 Having established the statistics of the hidden layer input currents for a novel and a familiar
 855 stimulus, we can now write down the conditions for the model to work, use them to find the optimal values
 856 of the parameters and calculate the true positive and false positive probabilities, and compute the capacity
 857 – the largest value of R for which the error is below a predetermined threshold. First, to ensure that exactly

858 one unit is active for a novel stimulus (true negative), since we are using a step function nonlinearity, we
 859 must have the largest input current take on a positive value,

$$860 \quad Kn + \xi\eta \sqrt{\frac{f_{\text{eff}}D}{N(1-\lambda^2)}} + b > 0$$

861 and second-largest to be below zero,

$$862 \quad K(n-2) + \xi\eta \sqrt{\frac{f_{\text{eff}}D}{N(1-\lambda^2)}} + b < 0$$

863 Second, to ensure there are no units active for a familiar stimulus (true positive),

$$864 \quad Kn - \eta\lambda^{R-1}D + \xi\eta \sqrt{\frac{f_{\text{eff}}D}{N(1-\lambda^2)}} + b < 0$$

865 For sufficiently large R , i.e. if $\eta\lambda^{R-1}D < 2K$, the third of these conditions implies the second. Since we are
 866 interested in maximizing R , we only consider the first and third conditions. Furthermore, note that these
 867 conditions are overparameterized. If we divide all three equations by η (e.g. let $k = \frac{K}{\eta}, B = \frac{b}{\eta}$), we can
 868 eliminate this free parameter. In other words, for any value of η we can scale K and b proportionally to
 869 satisfy the conditions, so for simplicity we choose $\eta = 1$. Similarly, the term $Kn + b$ can be replaced by a
 870 single parameter since for any choice of K we can rescale b to keep this sum constant. To ensure that the
 871 condition $\eta\lambda^{R-1}D < 2K$ holds for all R , we can choose $K = D$. For convenience, we also let $b = \beta D$ and
 872 $\sqrt{\frac{f_{\text{eff}}D}{N(1-\lambda^2)}} = \alpha_\lambda D$, the subscript indicating explicit dependence on λ . Dividing both inequalities by D , the
 873 conditions simplify to

$$874 \quad n + \alpha_\lambda \xi + \beta > 0, n + \beta + \alpha_\lambda \xi - \lambda^{R-1} < 0$$

875 The accuracy, i.e. probability of a correct response, is given by $P_{\text{correct}} = (1-f)P_{TP} + fP_{TN}$. For
 876 convenience, we compute the false positive instead of the true negative rate, noting that $P_{TN} = 1 - P_{FP}$. The
 877 false positive and true positive rates are given by

$$878 \quad P_{FP} = \mathbb{P}\left[\xi < -\frac{n + \beta}{\alpha_\lambda}\right], P_{TP} = \mathbb{P}\left[\xi < -\frac{n + \beta - \lambda^{R-1}}{\alpha_\lambda}\right]$$

879 Since ξ is a standard Normal random variable, $\mathbb{P}[\xi < z] = \frac{1}{2}\text{erfc}\left(-\frac{z}{\sqrt{2}}\right)$, so

$$880 \quad P_{FP} = \frac{1}{2}\text{erfc}\left(\frac{n + \beta}{\alpha_\lambda \sqrt{2}}\right), P_{TP} = \frac{1}{2}\text{erfc}\left(\frac{n + \beta - \lambda^{R-1}}{\alpha_\lambda \sqrt{2}}\right)$$

881 We would now like to set the optimal values of λ and β which maximize R , given a desired true
 882 positive and false positive probability P_{FP}^*, P_{TP}^* . Note that fixing these probabilities also fixes $f_{\text{eff}} = f^* =$
 883 $(1 - P_{FP}^*)f + (1 - P_{TP}^*)(1 - f)$. Rearranging the previous equations, we get

$$884 \quad \frac{n + \beta}{\alpha_\lambda} = \sqrt{2}\text{erfc}^{-1}(2P_{FP}^*), \frac{n + \beta - \lambda^{R-1}}{\alpha_\lambda} = \sqrt{2}\text{erfc}^{-1}(2P_{TP}^*)$$

885 The first equality sets the value for β . To determine λ , we substitute β into the second equality to get

886
$$\sqrt{2}\operatorname{erfc}^{-1}(2P_{FP}^*) - \frac{\lambda^{R-1}}{\alpha_\lambda} = \sqrt{2}\operatorname{erfc}^{-1}(2P_{TP}^*)$$

887 For notational convenience, let $E = \sqrt{2}[\operatorname{erfc}^{-1}(2P_{FP}^*) - \operatorname{erfc}^{-1}(2P_{TP}^*)]$. Using the definition of α_λ and $f_{\text{eff}} = f^*$, we have $\lambda = \sqrt{1 - \frac{f^*}{\alpha_\lambda^2 ND}}$. Rearranging, we have

889
$$\left(1 - \frac{f^*}{\alpha_\lambda^2 ND}\right)^{\frac{R-1}{2}} = \alpha_\lambda E$$

890 Assuming N and D are large (so λ is close to 1), we can use the first-order Taylor expansion $\exp(-z) \approx 1 - z$ for the term in parentheses (this will be necessary to get a closed-form expression for the optimal λ)
891 and solve for R

893
$$\exp\left(-\frac{f^*}{\alpha_\lambda^2 ND} \cdot \frac{R-1}{2}\right) = \alpha_\lambda E \Rightarrow R = 1 + \frac{2ND\alpha_\lambda^2}{f^*} \ln\left(\frac{1}{\alpha_\lambda E}\right)$$

894 Setting $\frac{dR}{d\lambda} = 0$ and solving for λ gives the optimum

896
$$\lambda = \sqrt{1 - \frac{eE^2 f^*}{ND}}$$

895 Thus, the capacity of the optimized network is

897
$$R_{\max} = 1 + \frac{ND}{eE^2 f^*}$$

898 where f^* and E are constants that depend on P_{FP}^* and P_{TP}^* (f^* also depends on the true fraction of novel
899 stimuli f). For instance, if we impose that $P_{FP}^* = 0.01$ and $P_{TP}^* = 0.99$, with our value of $f = \frac{2}{3}$, we get

900
$$R_{\max} = 1 + \frac{ND}{e \cdot 2 \cdot [\operatorname{erfc}^{-1}(2P_{FP}^*) - \operatorname{erfc}^{-1}(2P_{TP}^*)]^2 \cdot [(1 - P_{FP}^*)f + (1 - P_{TP}^*)(1 - f)]}$$

901
$$= 1 + \frac{ND}{2e[1.645 - (-1.645)]^2[0.98f + 0.01]}$$

902
$$= 1 + \frac{0.017ND}{0.98f + 0.01}$$

903
$$= 1 + 0.026ND$$

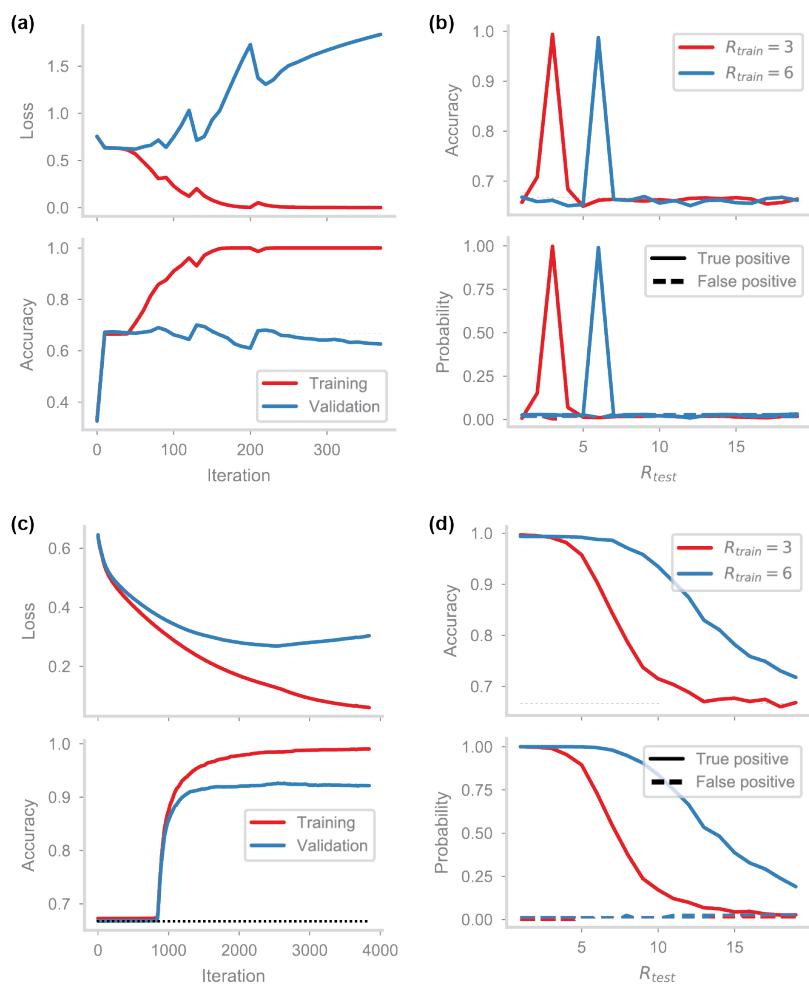
904 It is clear that the capacity scales in proportion to the number of plastic synapses in the network.
905 Furthermore, since $d = n + D$, i.e. $D = d - \log_2(N)$, the capacity scales in proportion to the total number of
906 synapses d , as long as $D \gg n$.

907
$$R_{\max} = O(ND) = O(N(d - n)) = O(Nd - N\log N) = O(Nd)$$

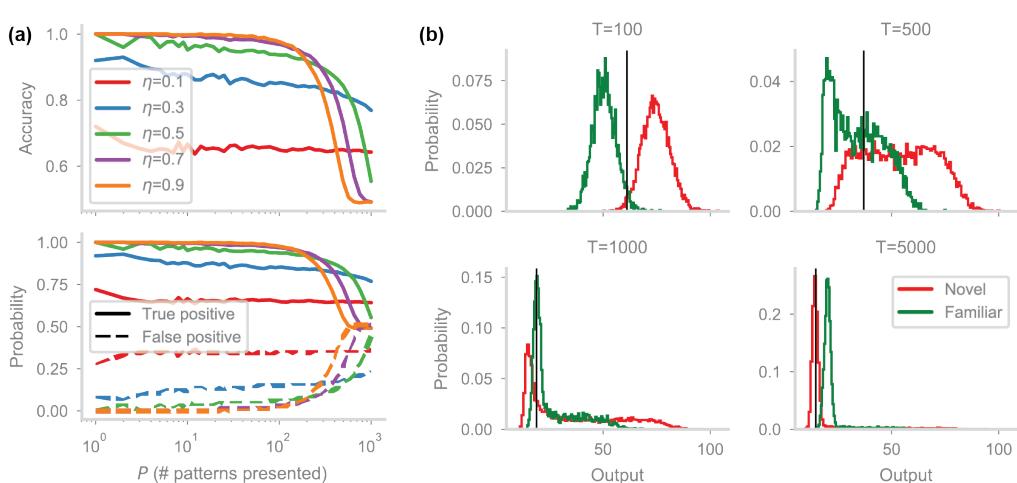
908 Finally, note that the equations for P_{FP} and P_{TP} are a function of f_{eff} due to the α_λ parameter, and
909 therefore recursively depend on P_{FP} and P_{TP} . We cannot compute the closed-form solution for these, but
910 we can approximate the values with arbitrary accuracy by iterating through this recurrence until
911 convergence to the fixed point. As the initial value for the recurrence, we use P_{FP} and P_{TP} computed using
912 $f_{\text{eff}} = f$, i.e. assuming no errors.

913

914 **Supplementary figures**

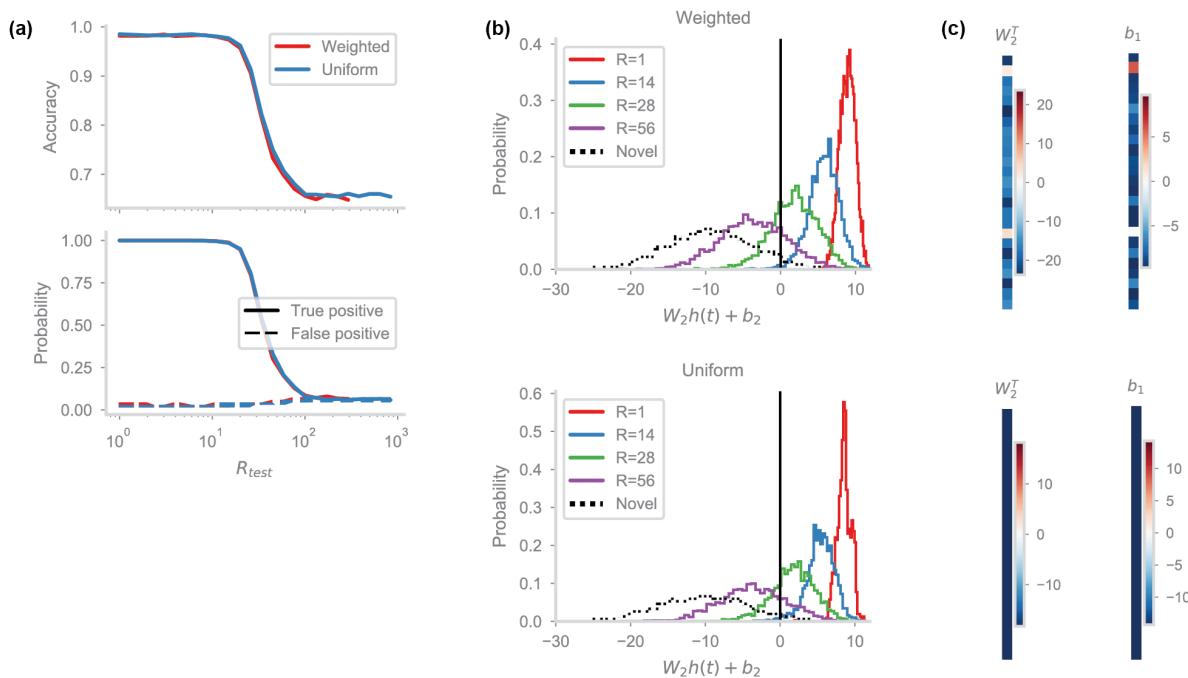


915
916 **Figure S1. HebbFF and RNN comparison, matching total number of dynamic variables.** (a,b) RNN performance as in Fig 2a-b,
917 with $d = 25$ and $N = 625$. (c,d) HebbFF performance as Fig 3a-b, with $d = 25$ and $N = 25$. The number of plastic synapses in the
918 HebbFF network $N * d = 625$ is the same as the number of recurrent units in the RNN, matching the number of dynamic variables
919 between the networks rather than the number of neurons. HebbFF still shows better generalization in both training scenarios.
920

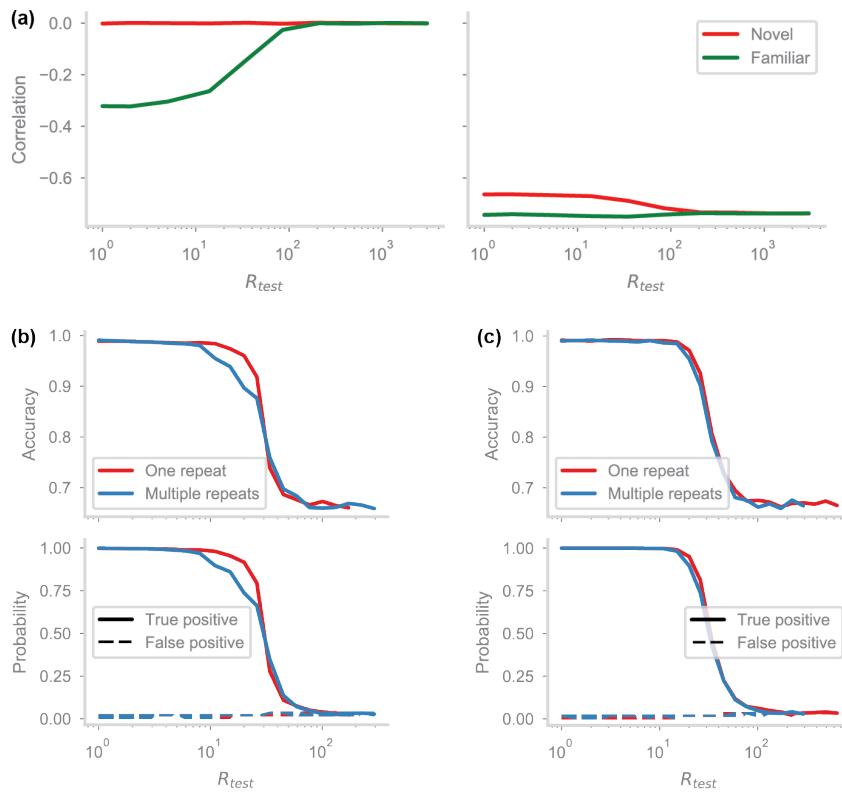


921
922 **Figure S2. Validation of network from Bogacz and Brown (2003).** (a) Performance of the anti-Hebbian network from (Bogacz and

923 Brown, 2003) ($d = N = 100$) on the non-continual familiarity detection task. The network is shown P uncorrelated randomly generated
 924 patterns, and tested on all the presented patterns as well as P novel ones. If the network's readout falls below a threshold, the
 925 corresponding input is classified as "familiar," otherwise "novel." The threshold is set such that exactly half of the test inputs are
 926 classified as "familiar." The plot shows the network's performance for a range of learning rates η as a function of the number of patterns
 927 presented. (b) Performance of the same network ($\eta = 0.7$) on the continual task. Plots show the probability distribution of the network's
 928 readout for novel (red) and familiar (green) stimuli. Threshold for familiarity (vertical black line) is set such that the top f of outputs are
 929 classified as novel (fraction f equal to proportion of novel stimuli in dataset), analogous to that in (Bogacz and Brown, 2003). Network
 930 performance declines as the distributions get closer together for increasing dataset size T . For very large dataset sizes, e.g. $T = 5000$,
 931 the readout for familiar stimuli is actually higher than that for novel, which causes overall performance (Fig. 3d) to fall below chance.
 932

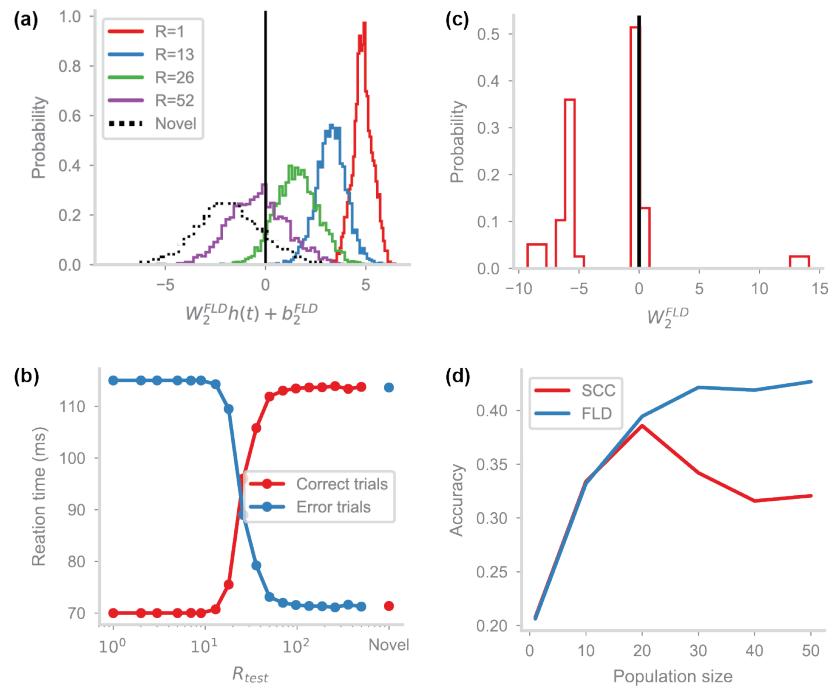


933
 934 **Figure S3. Comparison of uniform and trained readout.** (a) Generalization performance of a trained HebbFF network ($d = N =$
 935 25) with a fully trainable readout matrix, i.e. a weighted average of the hidden layer activity (red) and a version where all the entries
 936 of the readout matrix are constrained to be equal, i.e. a scaled average of the hidden layer (blue). (b) Distributions of the output unit
 937 activity (prior to applying the nonlinearity) for a trained network, for several test values of R . The output distributions, as well as
 938 performance, are almost identical for both versions of the readout. (c) Examples of the W_2 readout matrix, as well as the bias b_1 for
 939 the weighted (top) and uniform (bottom) readouts. Values for both are negative, indicating a qualitatively similar readout mechanism.
 940 Note that the weight matrices are plotted transposed for visualization.
 941



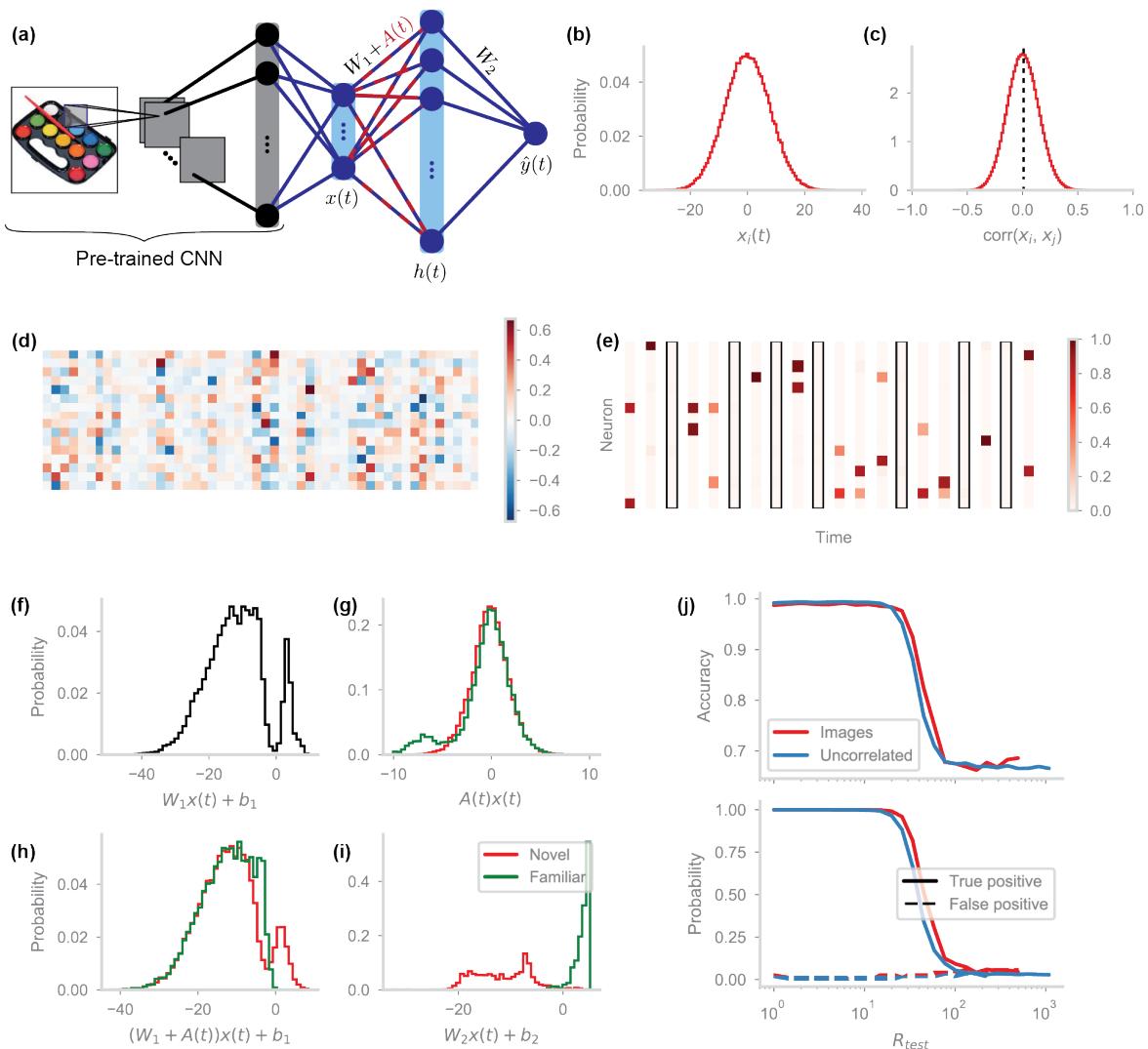
942
943
944
945
946
947
948
949
950
951
952

Figure S4. Idealized and HebbFF model differences. (a) Correlation between the hidden layer input currents due to static and plastic synapses for the idealized model (left) and HebbFF (right) for novel (red) and familiar (green) stimuli as a function of delay interval ($d = 25, N = 32$). In the idealized model, due to the split static and plastic synapses, the currents are uncorrelated for novel stimuli and familiar stimuli at long delay intervals. At short intervals, the two input currents are anti-correlated, which enables repetition suppression. In the HebbFF model, there is anti-correlation in both cases, although there is still less anti-correlation for novel stimuli. (b) Performance of the idealized network on the continual familiarity detection task with familiar stimuli repeated either exactly once, as used throughout this work (red), or multiple times (blue). Since there is exactly zero hidden unit activity for a familiar stimulus it does not get reinforced in memory, and less likely to be recognized on its second and subsequent repetitions. (c) The HebbFF network, trained to maximum capacity on the single-repeat task does not suffer any loss in performance due to multiple repeats.



953
954
955
956

Figure S5. Behavior on augmented task. Panels correspond to the right-hand side plots of Fig 7(a-d), but for the classifier-augmented HebbFF network ($d = 25, N = 50$) performing binary classification and familiarity detection simultaneously.



957
958
959
960

Figure S6. Performance on real-world images. Subplots as in Fig 8, but using a trained fully-connected linear layer to transform the activations of the CNN's penultimate layer into the inputs of HebbFF ($R_{\text{train}} = R_{\text{test}} = 19$).