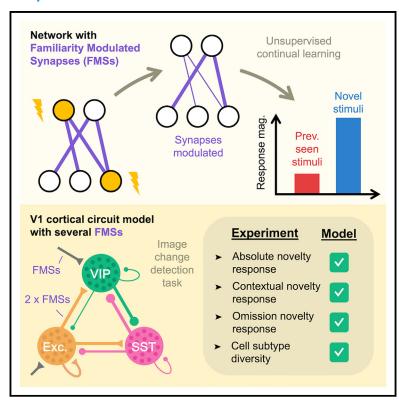
# Simple synaptic modulations implement diverse novelty computations

#### **Graphical abstract**



#### **Authors**

Kyle Aitken, Luke Campagnola, Marina E. Garrett, Shawn R. Olsen, Stefan Mihalas

#### Correspondence

kyle.aitken@alleninstitute.org

#### In brief

Aitken et al. introduce a simple, biologically inspired model for synaptic plasticity that leads to distinct responses to novel versus familiar stimuli. Using an experimentally constrained model of a cortical circuit with plasticity at specific synapses, multiple types of complex novelty effects recently observed in experiment are simultaneously reproduced.

#### **Highlights**

- A biology-motivated model of synaptic plasticity can encode several types of novelty
- Relative to other models, its simplicity allows for flexibility and understandability
- Absolute, oddball, and omission novelty are reproduced in a cortical circuit network
- The network develops cell subtype diversity, leading to testable connectivity predictions







#### **Article**

# Simple synaptic modulations implement diverse novelty computations

Kyle Aitken. 1,5,\* Luke Campagnola. Marina E. Garrett. Shawn R. Olsen. and Stefan Mihalas 1,4

<sup>1</sup>Center for Data-Driven Discovery for Biology, Allen Institute, Seattle, WA 98109, USA

\*Correspondence: kyle.aitken@alleninstitute.org https://doi.org/10.1016/j.celrep.2024.114188

#### **SUMMARY**

Detecting novelty is ethologically useful for an organism's survival. Recent experiments characterize how different types of novelty over timescales from seconds to weeks are reflected in the activity of excitatory and inhibitory neuron types. Here, we introduce a learning mechanism, familiarity-modulated synapses (FMSs), consisting of multiplicative modulations dependent on presynaptic or pre/postsynaptic neuron activity. With FMSs, network responses that encode novelty emerge under unsupervised continual learning and minimal connectivity constraints. Implementing FMSs within an experimentally constrained model of a visual cortical circuit, we demonstrate the generalizability of FMSs by simultaneously fitting absolute, contextual, and omission novelty effects. Our model also reproduces functional diversity within cell subpopulations, leading to experimentally testable predictions about connectivity and synaptic dynamics that can produce both population-level novelty responses and heterogeneous individual neuron signals. Altogether, our findings demonstrate how simple plasticity mechanisms within a cortical circuit structure can produce qualitatively distinct and complex novelty responses.

#### **INTRODUCTION**

Brains of complex organisms contain internal representations of the world that are shaped by stimuli they have become familiar with over time. Since their environment can change rapidly, an organism's survival can be dependent upon its ability to quickly identify novel stimuli. Indeed, over decades of study, effects of stimulus novelty have been found throughout the brain and are known to occur over many timescales 1-5 These effects vary from internal changes, such as promoting learning and memory, to behavioral adjustments, including changes to perception, attention, and exploration.<sup>2-4</sup> Across sensory modalities and species, novel stimuli are generally associated with an increased response relative to their familiar counterparts.<sup>2,3</sup> Such novelty responses (or their inverse, familiarity-responses) have been observed in cortical, subcortical, and neuromodulatory areas of the brain at both an individual cell level<sup>6-9</sup> and across macroscopic cell populations. 10-12 Additionally, studies have distinguished responses to distinct types of novelty. For example, absolute novelty, when an organism is exposed to a previously unobserved stimulus, 10,13 is distinguished from contextual (or oddball) novelty, where a previously observed stimulus is novel only in the context of recently observed stimuli that may also occur from the omission of an expected stimulus. 14-16

The mammalian neocortex is believed to play an especially important role in modeling the world around us and thus how it

responds to these various types of novel stimuli is of great interest. Within the cortex, what is believed to be a general purpose disinhibitory circuit is repeated across different brain regions and species, and many recent experimental studies have elucidated the properties of the cells within this circuit. 17-20 Specifically, the structure of this cortical circuit is defined by connectivity between somatostatin (SST) and vasoactive/intestinal peptide (VIP) expressing inhibitory interneurons as well as pyramidal excitatory neurons.<sup>21</sup> This circuit is thought to facilitate novelty responses through mutual inhibition between the VIP and SST populations that provides a disinhibitory pathway from VIP to excitatory cells.<sup>22</sup> Recent experimental studies have found that novelty responses vary significantly across these distinct cell populations.<sup>23,24</sup> These studies suggest that the enhanced response of VIP cells to novel stimuli suppresses the SST population's response, releasing the local excitatory population from inhibition and leading to an increased excitatory novelty

Although broad cell classes are a useful simplification to understand the function of the cortical circuit, each class can be further divided into subclasses or types that differ in gene expression patterns, synaptic connectivity, electrical properties, and morphology. <sup>19,25–28</sup> Indeed, within the excitatory, SST, and VIP cell populations, subpopulations that have distinct feature-coding across familiar and novel stimuli have been recently identified. <sup>23,24</sup> Given these recent results, an open question is what



<sup>&</sup>lt;sup>2</sup>Allen Institute for Brain Science, Seattle, WA 98109, USA

<sup>&</sup>lt;sup>3</sup>Allen Institute for Neural Dynamics, Seattle, WA 98109, USA

<sup>&</sup>lt;sup>4</sup>Applied Mathematics, University of Washington, Seattle, WA 98195, USA

<sup>5</sup>Lead contact



biological mechanisms might allow populations to have such diversity in experience-dependent coding, and how this coding diversity relates to changes in the population's macroscopic response to novel stimuli.

Since the observation of the brain's ability to rapidly detect novel stimuli, computational models have been used to investigate how the brain might distinguish familiar representations and evoke distinct responses to unfamiliar stimuli. 16,29,30 Many of these models rely on modifications of synaptic connections to encode stimuli. For example, Hopfield networks can encode familiar stimuli via lateral connections and are capable of recalling said stimuli using recurrent activity.31 However, many of these computational models require carefully placed synaptic connections to encode distinct memories<sup>32,33</sup> or strict training and testing phases that do not reflect an organism's natural behavior, 32,34 both of which limit their ability to be implemented into more general models. Additionally, some models rely on complex non-local credit assignment mechanisms that are biologically unrealistic to develop their novelty responses. 35,36

In this work, we introduce a mechanism that implements simple plasticity rules via synaptic modulations and is capable of adapting to stimuli through biologically realistic local, unsupervised learning. Broadly, it relies on modulating the synapses that play a role in producing the output responses of familiar stimuli, and as such we refer to the mechanism as familiaritymodulated synapses (FMSs). A strength of FMSs is their simplicity and thus generality; we show FMSs can broadly represent various synaptic plasticity effects that occur over different timescales. We focused on parameterizing the FMSs such that they represent biologically realistic plasticity mechanisms such as long-term potentiation/depression (LTP/D)<sup>37,38</sup> or short-term synaptic plasticity (STSP).39 FMSs can be implemented on a set of excitatory or inhibitory synapses feeding from one cell population to another whose strengths and connections are randomly drawn, meaning it requires essentially no specialized architecture and is thus straightforward to implement into more complex neural network models. The mechanism also requires no specific training regimen, simply becoming adapted to stimuli it has seen in recent history under continuous learning, similar to how biological organisms learn.

We first establish properties of the FMSs in the simplest possible feedforward setting. Afterward, we incorporate several distinct FMS mechanisms into a model of the visual cortical circuit, with connectivity properties constrained from multi-patch synaptic physiology studies, 20 relative cell counts from in situ hybridization experiments, 17,18 and additional cell properties from electrophysiology recordings. 19 We demonstrate the generalizability of the FMSs by modeling three distinct novelty effects: absolute, 10,13 contextual (oddball), 14,15 and omission novelty. 16 Although each of these novelty effects has been studied in isolation, recent studies in the visual cortical circuit of mice investigate how distinct cell populations respond to all three types of novelty.<sup>23,24</sup> The flexibility of FMSs allows for us to simultaneously capture the three novelty effects within our experimentally constrained model of the cortical circuit, while also reproducing the diverse subpopulation coding seen in the same experiments.24

#### **Related works**

Many existing models of novelty detection rely on modifications of synaptic connections in order to encode familiar stimuli, but often require specialized connection architectures in order to encode distinct memories, 32,33 do not operate under a continual learning setting, <sup>22,32,34,40,41</sup> or rely on complex non-local credit assignment, 35,36,42 all of which the FMSs avoid. Lim et al.40 and Sukbin<sup>41</sup> consider how a firing-rate dependent learning rule, directly derived from passive and dimming-detection experiments, can match time-averaged and time-dependent responses. Feedforward adaptation as a means of repetition suppression has been previously studied previously 22,29,35,42-44 and is advantageous because it does not require convergence to a steady-state or feedback-dependent activity to distinguish stimuli.45 Novelty responses on an image change detection task were reproduced using STSP-like synaptic modulations.<sup>29,42</sup> The specific form of the synaptic modulations used in this work are an unsupervised version of those described in Tyulmankov et al.<sup>35</sup> and Aitkin and Mihalas<sup>36</sup> that originated in the learning-how-to-learn machine learning literature.

Many other computational models of the visual cortical circuit have been built to understand the individual cell population effects of disinhibition and how the circuit's activity might change over learning. 18,22,47-54 While many models of cortical circuits treat inhibitory interneurons as a unitary population, 47,53 more recent models have incorporated the diversity of interneuron populations, including the VIP-SST-Excitatory disinhibitory circuit. 18,22,48-52,54 Keller et al.51 studied and modeled the VIP-SST-Exc. disinhibitory circuit in L2/3 of mice in the setting of visual context modulation and found contextual modulation is unlikely to be inherited from L4 and thus may rely on local circuitry. A computational model of the cortical circuit constrained by electrophysiological studies that incorporates population diversity and inhibitory plasticity was recently used to study prediction errors in Hertäg and Sprekeler<sup>52</sup> and Hertäg and Clopath. 54 Although they also investigate how connectivity influences the development of neuron subpopulations, the training/ testing stimulus sequences are different from the ones we investigate here.

#### Setup: FMSs

In this work we consider networks of firing-rate, point-like excitatory, and inhibitory neurons that can influence one another through synapses that we represent using weight matrices. Let W represent a set of fixed synapses that connect a presynaptic population of neurons to a postsynaptic population, with firing rates at time t represented by the vectors  $\mathbf{x}_{t}^{\text{pre}}$  and  $\mathbf{y}_{t}^{\text{post}}$ , respectively (Figure 1A, left). For example, the postsynaptic population's firing rates may be related to the presynaptic population's activity via  $\mathbf{y}_t^{\mathrm{post}} = \phi \left( \mathbf{W} \mathbf{x}_t^{\mathrm{pre}} \right)$ , where  $\phi(\cdot) \geq 0$  is a nonlinear function that accounts for the postsynaptic neurons' properties such as their firing threshold and maximum firing rate. We take W to be sparse and, for simplicity, take the nonzero weights to be drawn from a normal distribution. Furthermore, the sign of the nonzero elements of W are fixed by the cell type of the presynaptic population: excitatory neurons only have positive-weight synapses so that they increase postsynaptic potentials and inhibitory neurons only have negative-weight synapses.





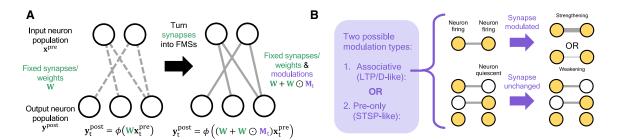


Figure 1. Familiarity modulated synapses

(A) On the left, an exemple feedforward firing-rate network, where a population of (firing-rate) input neurons,  $\mathbf{x}^{\text{pre}}$ , influences a population of output neurons,  $\mathbf{y}^{\text{post}}$ , through a set of fixed synaptic connections,  $\mathbf{W}$ . On the right, the fixed synaptic connections are modified to become familiarity modulated synapses (FMSs), i.e.,  $\mathbf{W} \to \mathbf{W} + \mathbf{W} \odot \mathbf{M}_t$ , allowing each synapse's strength to be modulated over time via the matrix  $\mathbf{M}_t$ .

(B) The two types of modulations we consider in this work: (1) associative and (2) pre-only dependent. See Equation 2 for explicit expressions. For an associative update rule, examples of how the behavior of neurons influences the way their synapses are modulated (see Figure S1A for equivalent pre-only diagram). In short, the modulations will either strengthen  $(\eta > 0)$  or weaken  $(\eta < 0)$  the neuron connections if both the pre- and postsynaptic neuron are firing and a synaptic connection already exists between said neurons.

We modify the fixed weights to be FMSs by taking

$$\mathbf{W} \rightarrow \mathbf{W} + \mathbf{W} \odot \mathbf{M}_t$$
, (Equation 1)

where  $\mathbf{M}_t$  represents time-dependent modulations to the synapses represented by  $\mathbf{W}$  and  $\odot$  is the element-wise product. In our exemplar network, the relation between pre- and postsynaptic activity would be  $\mathbf{y}_t^{\mathrm{post}} = \phi \left( (\mathbf{W} + \mathbf{W} \odot \mathbf{M}_t) \mathbf{x}_t^{\mathrm{re}} \right)$  (Figure 1A, right). We investigate two distinct modulation mechanisms throughout this work that determine how  $\mathbf{M}_t$  evolves in time,

$$\mathbf{M}_{t} = \lambda \mathbf{M}_{t-1} + \eta \mathbf{y}_{t}^{\text{post}} (\mathbf{x}_{t}^{\text{pre}})^{T}$$
, (associative), (Equation 2a)

$$\mathbf{M}_t = \lambda \mathbf{M}_{t-1} + \eta \mathbf{1} \left( \mathbf{x}_t^{\mathrm{pre}} \right)^T \bigg/ \sqrt{n} \;, \quad \text{(pre-only)} \;. \;\; \text{(Equation 2b)}$$

Both rules are completely unsupervised and modulated based on only information locally available to the synapse. The associative update, Equation 2A, is the more general modulation rule dependent upon both the post- and presynaptic neuron firing rates at time t (Figure 1B). The parameter  $0 < \lambda < 1$  controls how quickly the modulations return to their baseline values, while  $|\eta|$  determines the size of the updates. Importantly, the sign of  $\eta$  controls the sign of  $\mathbf{M}$  and thus whether synapses are strengthened or weakened by the modulations, i.e., if their magnitude increases or decreases, respectively. The pre-only modulation update expression, Equation 2B, is only dependent on the presynaptic firing rate, in which case the dependence on  $\mathbf{y}_t^{post}$  is replaced with  $\mathbf{1}$ , the all 1's vector, and normalized by the square root of the number of output neurons n.

Throughout this work, all **W** are fixed and thus the total synapse strength is only modified through the  $\mathbf{M}_t$  term. "Training" will refer to the time period where a network is exposed to certain stimuli and its synapses are modified solely via the unsupervised *FMSs* described above. Crucially, we do not allow the modulations to change whether a synapse is excitatory or inhibitory, i.e., if  $W_{ij} \geq 0$  then  $W_{ij} + W_{ij}M_{t,ij} \geq 0$  for all time. For simplicity, we also do not allow for new synapses to form, i.e., a synapse that does not exist at initialization cannot be modulated.

Biologically, we envision the modulations as various mechanisms leading to changes in the synapses that occur over varied timescales and biological mechanisms. The associative mechanism, Equation 2A, could broadly represent long timescale synaptic changes resulting from LTP/D mechanisms. Long-term potentiation or depression of said synapses can be implemented by changing the sign of the learning rate,  $\eta$ . Meanwhile, the modulations that are only presynapse-dependent, Equation 2B, could represent faster modulation mechanisms such as STSP. With these biological mechanisms in mind, we limit the size of the modulations such that they do not exceed synaptic changes that have been observed in the experiment (see STAR Methods for additional details).

#### **RESULTS**

#### A simple, unsupervised, feedforward novelty-detector

To explore some basic properties of the FMSs, we first investigate their effect in a simple feedforward network that we show develops distinct responses to stimuli it has been exposed to before, what we refer to as *familiar* stimuli throughout this work. Many of the results we establish in the simple network with a single FMS mechanism generalize to the visual cortical circuit model we discuss afterward in the section "Cortical microcircuit novelty response in a stimulus change task" with several distinct FMS mechanisms.

We represent the neuronal encodings of stimuli using distinct sparse random binary vectors (Figure 2A, STAR Methods). Prior to training, we draw two sets of eight stimuli from this distribution. During training, the stimuli from what becomes the familiar set will be exposed to the network while its weights undergo unsupervised updates via an FMS mechanism. After training, we compare the network's response to the familiar set and the other set that was held out during training, what we refer to as the *novel* set. Noise is added to all input stimuli throughout this work (STAR Methods).

The simple network consists of only two populations of neurons, an excitatory input population and an arbitrary output



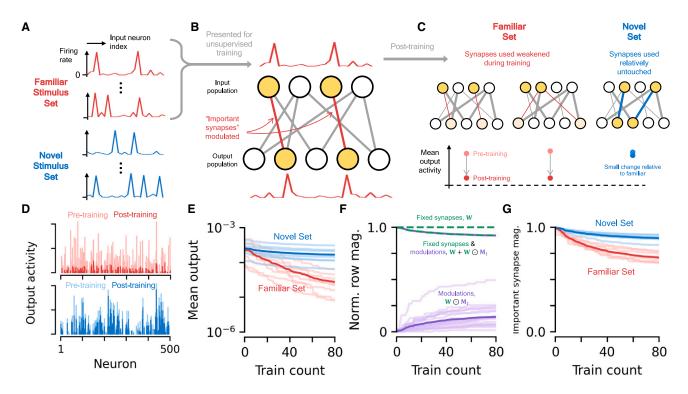


Figure 2. Familiar modulated synapses in a simple network

(A-C) Schematic of network behavior and exposure to familiar set.

(A) The familiar (red) and novel (blue) sets of stimuli that excite the input neuron population are drawn from the same distribution, random sparse binary vectors with added noise.

(B) We consider a simple two-layer network with FMSs connecting an excitatory input population to an output neuron population. At each time step, a randomly chosen familiar stimulus excites the input population and, through the modulated synapses, causes the output population to fire in some pattern. For the example considered in this figure, the associative modulations weaken any synapses that connect a pre- and postsynaptic neuron that both fired for the given familiar stimulus, e.g., an effect that could arise from LTD.

- (C) After training, many of the network's synapses have been modulated, changing its output behavior. The familiar set's mean output activity is reduced relative to its pre-training activity. The post-training mean output activity of the novel set is relatively unchanged.
- (D-G) Results from example network and training. In this example, there are eight familiar and eight novel stimuli. Each familiar stimulus has been input into the network 10 times (shuffled order) for 80 training steps total.
- (D) Example raw output response activity for a familiar (red) and novel (blue) stimulus pre- and post-training.
- (E) Change in mean output activity of the familiar and novel sets over training. Mean output activity across each stimulus set (dark) and individual stimuli (light) shown.
- (F) Normalized mean row magnitude of the modulation term, **W** ⊙ **M**<sub>ℓ</sub> (purple), the unmodulated weight matrix, **W** (green), and total synaptic strength (green and purple) over training. Mean (dark) and individual rows (light) shown.
- (G) Change in important synapse magnitude for familiar and novel inputs as a function of training time (STAR Methods). See also Figure S1.

population (in this setup, cell type [excitatory versus inhibitory] only influences the sign of weights leaving a population; since the activity of the output neuron population is directly measured, results here hold for either excitatory or inhibitory output neurons; an excitatory input population was chosen for simplicity, see STAR Methods for the equivalent setup with inhibitory neurons as well) that are sparsely connected by synapses represented by the weight matrix  ${\bf W}$  and with nonlinearity  $\varphi(\,\cdot\,)$  providing the output population activity (Figure 2B). For brevity, we refer to this network as the *familiar modulated synapse network* (FMSN). Before training, the synapse strengths are randomly initialized, but they are subject to modulations via an FMS mechanism, represented by the matrix  ${\bf M}_t$ . The two modulation types of Equation 2 and the possibility of strengthening or weakening synapses (i.e., the sign of  $\eta$ ) gives four qualitatively

distinct FMSs. For the example we explicitly consider here, we take the FMS's modulations to be associative and weakening, meaning a synapse/weight is weakened if both its pre- and post-synaptic neurons are firing, e.g., it is LTD-like (Figure 2C). This corresponds to updates via Equation 2a with  $\eta\!<\!0$ . Equivalent plots for the pre-only rule, e.g., STSP-like, and synapses that are strengthened by the modulations, e.g., LTP-like, are provided in the supplemental figures (Figure S1). We will later return to how these choices affect the results presented here.

### The FMSN develops distinct responses to familiar and novel stimuli

We use a training schedule where the FMSN is sequentially passed stimuli from the familiar training set several times in a random order. That is, at each time step, a stimulus is randomly

#### **Article**



drawn from the familiar set, noise is added to it, and it is input into the network. After each pass through the network, the FMSs are updated according to Equation 2a. For the example considered here, each familiar stimulus is presented to the network 10 times, for a total of 80 training steps. Post-training, we observe that the familiar output activity is significantly suppressed relative to its pre-training activity (Figure 2D). Comparatively, the novel output activity changes little from the modulations, and so post-training its activity is large relative to the familiar set. (To evaluate the novel activity pre-training without it becoming "familiar" to the network, we treat as we would a test set and do not modulate the synapses from its activity via Equation 2a. The FMSN then has no memory of being exposed to it. We emphasize this is done solely for the sake of comparison to the familiar set and is not a necessary step in training.) We can understand how the network's response changes during training by comparing the output activity of the familiar and novel sets had we stopped training after a certain number of familiar stimulus exposures. Over the course of training, we see the network's response to all eight familiar stimuli quickly weakens while its response to the eight stimuli of the novel set remains relatively unchanged (Figure 2E). This happens concurrently with a growth in the size of the synaptic modulations and, since the modulations in this example are weakening, a smaller total synaptic magnitude (Figure 2F). Eventually, the changes to the network stabilize as additional examples continue to be presented. The reduction of output activity for the familiar stimuli occurs concurrently with a sparser response to the familiar stimuli over time as well as decreased decodability of stimulus identity, consistent with experimental results of familiarization (STAR Methods, Figures S2A-S2C). 23,24

## **Distinct "important synapses" lead to distinct responses**

What about the pattern of synapse modulation causing this significant change in response for stimuli in the familiar set? Although almost all synapses undergo some modulation during training (a byproduct of the noise added to inputs), only a small percentage are modulated significantly (Figure S2D). Intuitively, a reason for the distinct output behavior could be that different synapses have large contributions to the output activity for members of the familiar and novel sets, so changing a subset of them only affects certain stimuli (Figure 2C). For a given stimulus, we define its important synapses as those synapses that would be modulated according to Equation 2a from passing the stimulus through the network, before any training has occurred (STAR Methods). With this definition, for the setup we consider here, each (nonzero) synapse has an approximately 2.5% chance of being an important synapse for a given stimulus. Prior to training, we can check that the important synapses of distinct stimuli have little overlap: a familiar and novel stimulus share on average only 0.14% of their important synapses. We can then track how the update rule of Equation 2a affects the important synapses of the familiar and novel sets differently. The total strength of the important synapses of the familiar set changes drastically, while those of the novel set remain relatively unchanged because of the small overlap of important synapses (Figure 2E). It is the greater weakening of important synapses associated to the familiar stimuli, often bringing the neurons' activity below firing thresholds, that leads to their significantly smaller responses relative to the novel stimuli.

The idea of targeted synaptic modulations as a means of encoding familiarity has been known for quite some time, most famously in Hopfield networks.<sup>31</sup> In the STAR Methods, we argue the FMSN can be approximately viewed as a feedforward Hopfield network, i.e., the weight modulations that encode the memory of the familiar inputs are on feedforward synapses and not lateral connections. A stimulus forward pass through the FMSN is similar to measuring its energy in the equivalent Hopfield network. Thus, familiar stimuli having a low mean response is similar to them being low-energy states.

## Synapse modulations change responses to stimuli in the subspace spanned by familiar stimuli

Since we draw the familiar and novel stimuli from the same distribution, between-stimulus correlations are relatively uniform across all stimuli. How would the FMSN respond to a stimulus that is more correlated with a familiar stimulus than the novel stimuli? More generally, one may consider what characteristics of stimuli determine how much they are suppressed by the learned modulations.

In the STAR Methods, we argue that the approximate M learned over the FMSN training causes any stimulus that lies in the subspace spanned by the familiar set to have a decreased response relative to its pre-training magnitude. (For this approximation, we have assumed that all familiar inputs are presented roughly the same number of times in a randomized order, as is done for the FMSN training. For cases where familiar stimuli are presented in an uneven manner, the network will respond most weakly to inputs it has been exposed to the most and those most recently presented, see the STAR Methods.) This includes the familiar stimuli themselves but also their linear combinations (Figure 3A). Furthermore, since any stimulus can be decomposed into parts that lie within and perpendicular to said subspace, the less any stimulus lies within this familiar subspace the less its response will be suppressed by modulations (Figures 3A and S3A). In other words, the more a stimulus is correlated with the familiar inputs, the more its response will be suppressed in the FMSN. Part of the success of the FMSN we investigate here relies on the fact that the familiar subspace is small relative to the full space of possible stimuli. Stimuli randomly drawn from the distribution that are not exposed to the network, e.g., the novel inputs, are likely to lie approximately perpendicular to this subspace and thus have their response relatively unchanged by training.

### Learning and decay rates strongly influence magnitude of modulation effects

For training, we have assumed that one stimulus is presented at each time step and time steps are separated by some  $\Delta t$  that could be a characteristic timescale of the input stimulus sequence. Of course, biological effects such as STSP and LTP/D can affect synapses over significantly different timescales. How can the FMSs be adjusted to account for such effects? To investigate this, it is useful to recast the FMS's decay rate,  $\lambda$ , as  $decay\ timescale$ ,  $\tau_{decay} = \Delta t/(1-\lambda)$ . Modifying



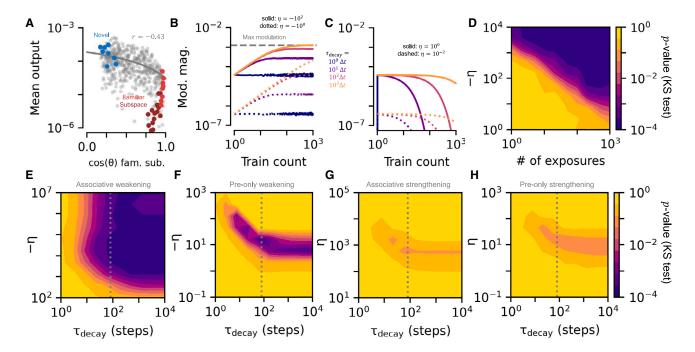


Figure 3. Additional properties of familiarity modulated synapses

(A) Cosine distance of stimuli to the subspace spanned by the familiar stimuli ("familiar subspace") versus mean output from the FMSN. Gray dots show sparse random binary vectors (STAR Methods). Familiar stimuli (light red), their linear combinations (dark red), and the novel stimuli (blue) are highlighted. Gray line shows linear regression fit.

- (B) Growth of modulation magnitude while being repeatedly exposed to a single familiar stimulus as a function of  $\tau_{\text{decay}} = 1/(1-\lambda)$ , in units of time steps, and  $\eta$ . Dashed gray line shows maximum modulation strength imposed by biological constraints (STAR Methods).
- (C) Decay of modulation magnitude after a single familiar stimulus exposure as function of  $\lambda$  and  $\eta$ .
- (D) Ability to distinguish output magnitude distributions of familiar and novel sets (KS-test p value) as a function of learning rate,  $\eta$ , and the number of times each familiar stimulus has been exposed.
- $(E-H)\ KS-test\ to\ distinguish\ post-training\ output\ magnitude\ distributions\ of\ familiar\ and\ novel\ sets\ for\ the\ four\ types\ of\ modulations\ as\ a\ function\ of\ \tau_{decay}\ and\ \eta.$
- (E) FMSN with associative weakening modulations. The gray vertical line shows the timescale of the task, 80 time steps.
- (F) Same as (E), for pre-only weakening modulations.
- (G) Associative strengthening modulations.
- (H) Pre-only strengthening modulations. See also Figures S2 and S3.

 $au_{
m decay}$  affects the time to saturation of the modulations, allowing one to tune both the number of stimuli and time it takes to see the modulations stabilize as well as their steady-state magnitude (Figure 3B). Varying the size of the FMS's other parameter, the learning rate  $\eta$ , affects the size of the modulations and thus the speed and magnitude of the FMSN's change in response. For large enough  $\eta$ , the modulations encounter the biological bounds, which limit their growth in size (Figure 3B). Relatedly, how long a given input influences the modulations, or, how long the FMSN "remembers" a past stimulus, is also affected by the decay timescale and learning rate (Figure 3C). A single familiar input can influence responses for only a few time steps or thousands, a fact that will play an important role later on when we model novelty effects of significantly different timescales.

The modulation learning rate can also influence how many exposures to the familiar set are needed in order for the network to develop distinct responses relative to the novel set. The larger the modulations, the greater the change to the FMSN from a single input stimulus, leading to distinct responses in a fewer number of stimulus presentations (Figure 3D). Notably, in the setup

we consider here, distinct responses can develop after *just* one exposure to each familiar stimulus. Although large learning rates can lead to quicker response changes, when one has noisy input stimuli, a large learning rate causes the modulations to also fit the noise. Indeed, for fixed training time, there exists optimal learning rates for distinguishing the familiar and novel sets that balance this trade-off between modulations that quickly capture the stimulus signal but not the noise (Figure 3E).

### What FMSN properties lead to significant differences in familiar and novel responses?

So far, we have specifically considered the case of an FMS that has associative updates that weaken the network's excitatory synapses. Of course, this covers a small subset of biological mechanisms—there are synapse modulations that strengthen connections, are only presynaptic dependent, and/or act on inhibitory synapses. The FMSs of Equation 2 are general enough to model all these cases.

Much of what we discussed above also holds for the presynapse-only update mechanism of Equation 2b that also weakens the excitatory synapses of the FMSN (Figure S1). However,

#### **Article**



because of its lack of postsynaptic dependence to pinpoint which synapses to update, the pre-only weakening mechanism is much more susceptible to noise. Too large of a learning rate can overfit the noise and quickly cause all inputs to be suppressed (Figure 3F). Surprisingly, we observe that distinguishing the familiar and novel outputs using modulations that strengthen the excitatory connections of the FMSN is significantly less effective for both associative and pre-only dependent FMSs (Figures 3G and 3H). Note that the strengthening of excitatory synapses enhances the response of familiar stimuli relative to their pre-training magnitudes (Figure S1B). We investigate what causes the differences between the strengthening and weakening FMSs in more detail in the STAR Methods (Figures S3B-S3J). In short, we find two major contributions to the relatively poorer performance of the strengthening mechanisms: (1) tighter modulation bounds for strengthening imposed by experiment and (2) neurons' nonlinear behavior that causes firing to cut off below certain potentials and saturate at higher potentials, built into  $\varphi(\cdot)$ . The latter of these effects can be partially overcome by considering an FMS that strengthens inhibitory synapses.<sup>22</sup> Stronger inhibition causes the output neurons' responses to get smaller, a similar effect as the weakening of excitation we found to be the most effective above (STAR Methods, Figures S3B-S3J]).

There are many other properties of the FMSN that can be explored that we only briefly touch upon here. For example, allowing modulations to further weaken or strengthen synapses beyond the bounds imposed by associating these modulations with LTP/D and STSP leads to even larger differences between the FMSN's response to familiar and novel stimuli (Figures S2E and S2F). Increasing the noise makes it harder for the FMS mechanism to isolate the signal, making it more difficult to distinguish novel and familiar responses (Figure S2G). However, effects from noise can be overcome by exposing the network to the familiar stimuli more times, giving it more observations to isolate the signal. Increasing both the number of input and output neurons also increases the distinguishability between the familiar and novel sets (Figure S2H). Though we leave a full investigation of FMS capacity for future work, we also see the FMSN is capable of becoming familiar with much more than eight stimuli while still having a distinct response to novel stimuli (Figures S2K and S2L). Last, we can use the FMSN to predict the most efficient coding of the sparse binary input vectors for distinguishing the familiar and novel sets. Lower sparsity reduces the variance in neuronal responses, and thus makes it easier to distinguish familiar and novel inputs, but also increases the similarity of any two stimuli because each one has more nonzero components. Thus, optimal sparsity is not too high or low (Figure S2J).

## Cortical microcircuit novelty response in a stimulus change task

We now implement the FMSs in a visual cortical circuit model to capture three distinct novelty responses recently observed in mice while they perform an image change detection task. <sup>23,24</sup> We note that the primary purpose of this model is to demonstrate the flexibility of FMSs and their ability to simultaneously produce three novelty effects observed in the VIP population recordings, <sup>23,24</sup> but do not attempt to constrain this as the

only types of plasticity that could lead to the experimentally observed results.

### Review of image change detection task and measurement

The stimuli used in the experimental task consist of a set of eight familiar training images and a held out set of eight novel images (Figure 4A). The task consists of image presentations from these sets at quick, regular intervals that are separated by a gray screen (Figure 4B). The same image is presented several times in a row before switching to another image within the set and mice are rewarded for responding to the image change by licking a waterspout. During this time, neuronal responses from the visual cortex are recorded in hour-long sessions using two-photon calcium imaging. Mice are trained on what becomes a familiar set of eight images and their neuronal responses are recorded in a "familiar" imaging session after achieving a performance threshold (Figure 4C). Shortly after, neuronal responses are also gathered over multiple sessions when the mice are exposed to the same task using the novel set of eight images. The mice's initial exposure and exposure after at least one session to this novel set of images are referred to as the "novel" and "novelplus" imaging sessions. Additionally, only during the imaging sessions, image omissions can occur, i.e., gray screen is displayed in place of a single image presentation (Figure 4B).

The responses to various novelty effects are recorded across several transgenic lines to capture excitatory, SST, and VIP population responses in the visual cortex. These cell populations form the cortical microcircuit discussed in the introduction whose connection probabilities and strengths have been carefully studied (Figure 4D). Experimental analyses show that the effects of novelty give rise to significantly different responses in these three populations, <sup>23,24</sup> which we discuss in more detail below.

#### **Cortical microcircuit model**

Given that we have observed the FMS mechanism yields distinct responses to familiar and novel stimuli, we built a model of the cortical microcircuit to study if it can develop the several experimentally observed novelty responses when exposed to stimulus sequences similar to that of experiment. Our firing-rate model consists of three groups of neurons, representing the SST, VIP, and excitatory neuron populations (Figure 4E). (Parvalbumin [PV] expressing inhibitory neurons is not included in our cortical circuit model directly, though the inhibition it provides to the other populations is partially accounted for from the threshold adjustments at the model's initialization [STAR Methods]. This important simplification is driven by the desire to build a minimal model of the data from Garrett et al., 24 where excitatory, VIP, and SST neurons were recorded and furthermore from the fact that VIP cells do not receive strong input from the PV population [Figure S4A]. The blanket inhibition in our model is in part supported by the general lack of specificity of PV to excitatory connections, 55 though more recent evidence points to some levels of specificity.<sup>56</sup>) The excitatory population receives inputs representing the bottom-up encoding of the raw stimulus sequence while the VIP population receives inputs representing topdown information about the history of the sequence (specifically



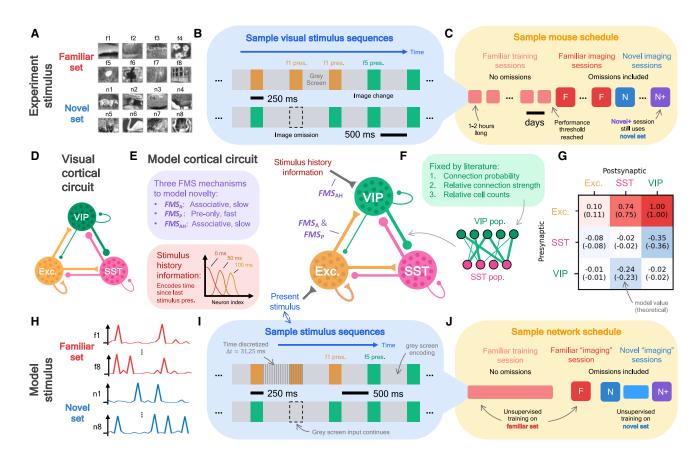


Figure 4. Image change task and visual cortical circuit: Experiment and model setup

- (A-C) Experimental stimulus details.
- (A) Example familiar and novel image sets (reproduced, with permissions from Garrett et al. <sup>24</sup>).
- (B) Sample stimulus sequences showing an image change (top) and omission (bottom).
- (C) Typical training/imaging schedule. Boxes represent sessions that occur on different days, each lasting an hour or two.
- (D) Diagram of a subset of the visual cortical circuit showing the cell populations that were recorded in experiment.
- (E) Diagram of the cortical circuit model we study in this work. The SST, VIP, and Exc. circles each represent populations of neurons connected by weights fixed from experimental data. <sup>17,18,20</sup> Three FMS mechanisms (purple), *FMS*<sub>A</sub>, *FMS*<sub>P</sub>, and *FMS*<sub>AH</sub>, are added to the network to model novelty responses. At each time step, the network receives inputs representing an encoding of the "present stimulus" being shown (blue) as well as "stimulus history" information (red) in the form of an encoding of the time since the last image presentation (STAR Methods).
- (F) Many features of the cortical circuit model are fixed by experimental literature. 17,18,20
- (G) Mean inter-population connection strengths. Values from an exemplar model (top) and analytically computed values (bottom) are shown (STAR Methods). (H–J) Model stimulus details.
- (H) Exemplar familiar and novel stimuli sets, drawn from a sparse random binary vector distribution.
- (I) Sample present stimulus sequences showing a stimulus change (top) and omission (bottom, STAR Methods).
- (J) Model training scheduling consisting of training session on familiar stimuli, familiar "imaging" session, and then novel/novel-plus "imaging" sessions. At all points of training and imaging, all FMSs are continuously updated via their unsupervised rules as stimuli are passed through the network. See also Figure S4.

the "timing" of recent stimuli and not image identity, see below for additional details). We estimate connection properties between populations by aggregating results from several recent experimental studies. Relative cell counts are estimated from *in situ* hybridization experiments<sup>17,18</sup> (STAR Methods). An estimation of inter-population connection probabilities comes from multi-patch synaptic physiology<sup>20</sup> and the relative strength of neuron connections between populations is estimated from the same study, supplemented with additional cell dynamical properties from electrophysiology recordings<sup>19</sup> (Figures 4F and 4G; STAR Methods). In particular, fits of measured postsynaptic potentials are used to estimate unmodulated, individual synapse

strengths as a function of the pre- and postsynaptic cell type (Figure S4). Due to the unprecedented detail of recent experiments, <sup>17–20</sup> coupled with necessary corrections from the experimental to the *in vivo* setting, we believe the "skeleton" of the cortical circuit model represents one of the most accurate estimates of this system to date.

We allow the connections in our microcircuit model to change by introducing several FMS mechanisms into the synapses connecting the various populations of the network. Since it is observed that the VIP cells drastically change their response across all three types of novelty in the experiment,<sup>24</sup> in this work we focus on adding FMS mechanisms to capture their

#### **Article**



specific novelty responses. The purpose of focusing only on the VIP response is to demonstrate how several FMS mechanisms may collectively model distinct novelty responses within a single population. We leave a complete modeling of the distinct cell type responses and related plasticity mechanisms for future work. To capture the VIP novelty responses, we add three separate FMS mechanisms to the synapses onto the VIP neurons: FMS<sub>A</sub>, FMS<sub>P</sub>, and FMS<sub>AH</sub> (Figure 4E).

- (1) FMS<sub>A</sub> (Associative, Exc. → VIP) is added on the synapses going from the excitatory to the VIP cells. Its learning and decay rate (η and λ) are tuned to learn and retain familiarity over a timescale of hours to days. Since it operates on a slow timescale and is pre- and postsynaptic dependent, FMS<sub>A</sub> could model LTD-like effects on said synapses
- (2) FMS<sub>P</sub> (Pre-only, Exc. → VIP) is also added to the set of synapses between the excitatory and VIP populations, but unlike FMS<sub>A</sub> it is tuned to learn and forget on a timescale of seconds. The fast timescale over which it operates and its presynaptic dependence makes FMS<sub>P</sub> a natural model for STSP-like effects on the synapses.
- (3) FMS<sub>AH</sub> (Associative, Stimulus history → VIP) is added to the synapses feeding into the VIP population from the stimulus history input neurons (see below). Its learning and decay rates are tuned to operate over long timescales, similar to the LTD-like FMS<sub>A</sub>.

Motivations for adding these particular modulations within the circuit are discussed below.

#### Model "image" change stimulus

As we saw in the FMSN, modulations are entirely driven by the stimuli being passed to the network, so we reproduce the pattern of stimuli from the image change detection experiment. To represent neuronal encodings of the images used in the experiment, we again use random sparse binary vectors as the distinct stimuli (Figure 4H). An "image" presentation is represented by a stimulus encoding being passed to the network for several time steps along with time-varying noise (Figure 4I). The image presentation is followed by a proportional number of gray screen time steps, where the network receives only noisy input (Figure 4I). This pattern repeats with a similar distribution of image change times used in the experiment (Figure S5A). Stimulus omissions are represented by additional gray screen time steps (Figure S5B). We assume the excitatory population receives this bottom-up present stimulus input and drives the other populations (Figure 4E). Additionally, we assume the microcircuit receives top-down inputs representing information about the recent history of the stimulus (STAR Methods, Figure S5C). In particular, the stimulus history input is an encoding of the time since the last stimulus presentation, with encodings of similar times more correlated than disparate times. (The primary purpose of this input is to give the microcircuit information about the recent stimulus history. A simple neuronal circuit that counts the time steps since the last stimulus presentation, e.g., an RNN, could represent the higher cortical areas that may produce this additional input directly from the bottom-up present stimulus input.) This information is passed directly to the VIP cells, which are known to receive feedback inputs from higher cortical areas.<sup>8,20</sup> Finally, time-correlated noise is injected into all neuron populations to represent activity from sources neglected in this model, e.g., activity from behavior (STAR Methods, Figure S5D).

Similar to the training schedule used in the experiment, we first expose the network to the familiar stimuli over a long training session, then gather cell responses to the task using the familiar set in what we continue to call an "imaging" session. Immediately afterward, we gather responses to the stimulus change task using the novel stimulus set, and, after additional exposure to the novel image set, gather the novel-plus responses (Figure 4J, STAR Methods). The familiar, novel, and novel-plus imaging session stimulus sequences are statistically identical. Importantly, the neuronal responses presented here are gathered in a continuous learning setting, i.e., the network continues to modulate its weights via FMSA, FMSP, and FMSAH at all steps of training and imaging. We scan over three parameters, the learning rates for all three FMS mechanisms, to determine modulation rates that best match experimental observations (STAR Methods, Figure S5O). We emphasize that, other than minor adjustments to the network at initialization to ensure realistic responses, the cortical circuit model only undergoes unsupervised adjustments via the various FMS mechanisms from exposure to stimulus sequences that closely match the stimuli on which the mice were trained (STAR Methods).

For the purposes of comparing our model to the experiment, we first focus on three distinct novelty responses that our model captures seen in mean VIP population responses of the experimental data<sup>24</sup>: (1) absolute, (2) contextual, and (3) omission novelty (see Figure S6 for SST and Exc.).

### Absolute novelty: familiar modulation occurs despite irregular stimulus sequence.

The change between the familiar and novel image sets represents absolute novelty-up until the novel imaging session the mice have never observed the set of images now used in the image change task. In both the experiment and our model, the VIP cells respond weakly to image presentations in the session that uses the familiar set relative to image presentations in the session that uses the novel set (Figure 5A). As we confirm below, for the cortical circuit model, the change in response is caused by FMSA, the slow-learning FMS mechanism on the excitatory to VIP synapses. FMSA functions almost identically to the FMSN discussed earlier: over training, exposure to the familiar stimuli causes the network to develop a suppressed response to them relative to the novel stimuli (Figure 5B). The stimulus sequence here is quite different from that of the FMSN: a single stimulus is repeatedly input to the network and is often separated by noisy gray screen. Additionally, the postsynaptic population of FMSA, the VIP cells, receive input from several additional sources such as the SST population and the recurrent VIP connections. Nevertheless, over the long familiar training period, FMSA gradually modulates the important synapses of the familiar set more than the novel set, leading to a distinct response across sessions (Figure 5B). Notably, the additional synaptic inputs and noise make modulating only those synapses



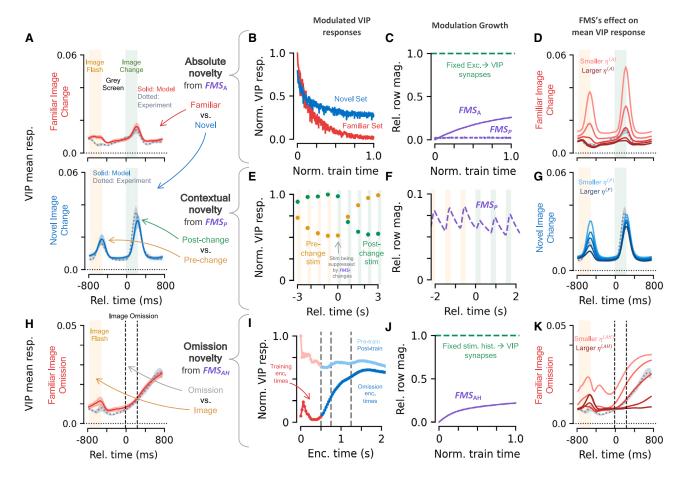


Figure 5. FMSs implement three distinct novelty effects in a cortical circuit model

(A) Mean VIP responses to image changes of the cortical circuit model (solid colored line) and experiment (gray dotted line).<sup>24</sup> Top shows mean response in the familiar imaging session and bottom in the novel imaging session. Green-shaded background represents time where changed image is presented, yellow-shaded is pre-change image.

- (B-D) Absolute novelty and FMSA.
- (B) Exemplar mean VIP image change responses to the familiar (red) and novel (blue) sets over training (STAR Methods).
- (C) Change in FMS<sub>A</sub> (purple solid) and FMS<sub>P</sub> (purple dashed) modulations over training. Green dotted line shows fixed portion of excitatory to VIP synapses.
- (D) Same as top of (A), for different FMS<sub>A</sub> learning rates (STAR Methods).
- (E-G) Contextual novelty and FMS<sub>P</sub>.
- (E) Normalized VIP responses around an image change event. Yellow is pre-change image, green is changed image.
- (F) Same as (C) but zoomed in to show change in FMSP over an example image change.
- (G) Same as bottom of (A), for different  $FMS_P$  learning rates.
- (H) Same as (A), for mean VIP response to image omission in a familiar session. Area between vertical dotted lines represents times where image would normally be presented.
- (I-K) Omission novelty and FMSAH.
- (I) Change in VIP response to various encoded times over training. Red encoded times are seen during training, blue are times when omissions are present.
- (J) Change in FMSAH (purple) modulations over training. Green dotted line shows fixed portion of stimulus history to VIP synapses.
- (K) Same as (H), for different FMSAH learning rates. See also Figures S5-S7.

important for the familiar set more difficult, leading to a fair amount of suppression to novel inputs as well. However, after these changes stabilize, we still observe distinct responses to the familiar and novel stimuli. Just like the FMSN, the change in stimulus response occurs concurrently with a growth in the modulations of  $FMS_A$  over training and, since the modulations are once again weakening, an overall decrease in the strength of the synapses connecting the excitatory population to the VIP population (Figure 5C).

To confirm that it is the modulations from  $FMS_A$  that cause the large difference in VIP response across the familiar and novel sessions, we can isolate its effect by training identical microcircuit models with different  $FMS_A$  learning rates. Indeed we see that, as we decrease (or increase)  $FMS_A$ 's learning rate and thus its modulation magnitude, the response of the VIP population in the familiar sessions grows (or shrinks) as the overall strength of the excitatory to VIP synapses changes (Figures 5D and S7A). Once again there is a trade-off between modulations



with too large of a learning rate that suppresses all responses and those with too small of a learning rate that suppresses none.

### Contextual novelty: fast familiarization and forgetting captures local oddball effects.

In the experimental task, image changes represent contextual novelty-since images are repeatedly presented at least 10 times, when the image identity changes it represents a local oddball and is contextually novel. In the novel session, we observe an increased response of VIP cells to image changes relative to the pre-change image in both our model and the experimental data (Figure 5A, bottom). Although smaller, the effect is also present in the familiar session (Figure 5A, top). Notably, this is a very different effect from the absolute novelty we discussed above; it only takes seconds for mice to establish an image as their baseline and this information is quickly updated to the current image being presented.<sup>24</sup> To model a novelty effect that learns and forgets quickly, FMSP is introduced. For FMS<sub>P</sub>, the presynaptic-dependent modulations make the most recent stimulus presentations become familiar, leading to smaller VIP response on repeats of the same stimulus. An image after a change is novel to FMSP, meaning the VIP response is larger because it is not familiarity suppressed (Figure 5E). After the change occurs, FMS<sub>P</sub> begins suppressing the important synapses of the current stimulus, while those that were important for the pre-change image are gradually released from suppression. This rapid turnover is reflected in the significantly quicker growth and decay of the FMS<sub>P</sub> modulations relative to those of FMS<sub>A</sub> (Figure 5F). Finally, as we did for FMSA, we see that varying the learning rate of FMS<sub>P</sub> isolates its effects on the cortical circuit response. A weaker learning rate changes the relative heights of the pre-change and post-change VIP responses because the image that has been presented several times in a row is less suppressed by modulations (Figures 5G and S7B).

The operation of both  $FMS_A$  and  $FMS_P$  on the excitatory to VIP synapses demonstrates a remarkable property of the FMSs: distinct types of modulations, e.g., slow and fast, can function on the same set of synapses simultaneously. This matches biology, where effects of LTP/D and STSP can affect the same set of synapses. The distinct FMSs can encode different types of novelty present in stimuli, allowing them to model the various novelty responses that are observed in certain neuronal populations. This allows for a single synaptic population to affect the postsynaptic population's response in a way that compounds the various novelty effects. For example, the largest responses of the VIP cells occur when both absolute and contextual novelty occurs, i.e., a novel image change, which is a result of the minimal suppression from  $FMS_A$  and  $FMS_P$  simultaneously.

### Omission novelty: a decrease in familiar correlation causes omission ramping.

Due to the temporal structure of the task during training, images are expected to be separated by 500 ms of gray screen. *Omission novelty* occurs when, instead of an image, additional gray screen is displayed, representing a global oddball. We observe a ramping response in the VIP cells of both the model and experiment when images are omitted (Figure 5H). In the

model, this is a result of FMSAH, the FMS mechanism on the stimulus history input synapses to the VIP cells. Recall that the stimulus history signal is a neuronal encoding of the time since the last stimulus presentation, where the encodings of similar times are more correlated than disparate times (Figure 5E, STAR Methods). Over training, FMSAH becomes familiar with encoded times-since-last-image that occur with no omissions, suppressing the corresponding VIP responses to the stimulus history signal. When omissions occur, longer-time encodings are passed to the network and the familiarity suppression is lost, leading to an increased response in VIP cells. The ramping occurs because the longer-time encodings have a less correlated representation to the familiar short-time encodings. That is, similar to what was seen in the FMSN, the network has formed a familiar subspace of the short-time encodings and the longertime encodings that gradually get farther from this subspace cause a gradual increase in the VIP response (Figure 5I). The encoded times that are novel but still close to the encoded times that are familiar, e.g., 510 ms, have their outputs quite suppressed. The longer encoded times, e.g., 1,000 ms, have outputs barely suppressed at all. As with the other FMSs we've investigated, this change in response occurs concurrently with a gradual growth in the modulations on FMSAH's synapses over training (Figure 5J). Additionally, the size and time of onset of the VIP ramping can be changed by adjusting the magnitude of FMS<sub>AH</sub>'s learning rate (Figures 5K and S7C).

The omission novelty responses occurring concurrently with the absolute and contextual demonstrates the ability to have multiple inputs into the same postsynaptic cell population with distinct synaptic dynamics.  $FMS_A$  and  $FMS_P$  operate on the excitatory to VIP synapses, while  $FMS_{AH}$  acts on the stimulus history to VIP synapses and all three can produce their corresponding novelty effect in the VIP cells when the corresponding stimulus occurs.

#### Novel images become familiar with exposure over time

Although the novel image set is initially unfamiliar to the mice and evokes distinct novelty-related responses across cell populations, over many exposures one would expect the images to gradually become familiar to the mice. Indeed, the enhanced VIP response to the novel images persists throughout the entire novel imaging session, but gradually disappears as the mice become accustomed to the novel set over many sessions of exposure. 23,24 Since our model is evaluated in a continuous learning setting, the FMS mechanisms are actively modulating the network's response, allowing it to also adapt to the novel stimulus set over time in the same way it adapted to the familiar set during training. Hence, our model also exhibits a gradual change in response to the novel image set over sessions, eventually returning to a suppressed VIP response to novel set images in the novel-plus imaging session (Figure S6A, top right; STAR Methods). In the experiment, even after being exposed to the novel set of images, the familiar set of images still evoke a response consistent with them being familiar stimuli.<sup>24</sup> The FMS<sub>A</sub> modulations decay slowly enough that the modulatory effects of both the familiar and novel images can persist simultaneously (Figure S5P). Additionally, as in the experiment, the image omission response does not change considerably between



the familiar and novel-plus sessions (Figure S6B, top right). (We do not attempt to model the suppressed omission ramping that is observed in the VIP population during the novel session that gradually returns to familiar levels in the novel-plus session.<sup>24</sup> See the discussion for potential mechanisms that can model this effect.)

### Cortical circuit model's cell subpopulations have diverse coding

Although the changes in mean response of our model's cell populations is dependent upon the behavior of individual cells within these populations, we observe a significant variation in each neuron's response over both stimulus features and experience level. Indeed, a key finding of Garrett et al.<sup>24</sup> is the emergence of functional cell subpopulations within the VIP, SST, and excitatory cell populations. Within each population, subpopulations are identified with similar changes in coding features over experience level, as measured by a coding score metric that we briefly review. To determine the coding score, each cell's response is fit using a kernel regression model. 57-60 Several input features that may influence a cell's response, including image presentations, image omissions, and task-relevant triggers such as image changes, are convolved with fit kernels to reproduce the cell's activity (Figure 6A, top; STAR Methods). To determine the importance of a given input feature category in explaining a cell's activity, the regression model is refit while omitting each category and its corresponding kernel(s). The coding score of a cell to a given feature category is defined as the relative amount of variance explained that the regression model loses by removing the feature category (Figure 6A, bottom; STAR Methods). This procedure is repeated across the three distinct sessions/ experience levels for four feature categories, resulting in a 12-dimensional coding score vector for each cell. Last, the resulting coding score vectors of a given cell type are collected across mice and run through an unsupervised clustering algorithm (Figure 6B, STAR Methods).<sup>24</sup>

We use the same analysis pipeline to analyze cell subpopulation diversity in our cortical circuit model. We again focus on investigating the VIP population's activity in particular (Figure S8 for Exc. and SST). Since we have no explicit behavioral effects in the network, we only code for three input feature categories: image presentations, omissions, and changes ("task"). Repeating the aforementioned fitting procedure to determine coding scores and then clustering the data across 10 network initializations, we again observe diverse coding across features and experience levels in the VIP population (Figure 6C). Notably, the resulting kernel fits of the cortical microcircuit qualitatively resemble the fits on the experimental data (Figures 6D and S8). Membership of the clusters is shared across the different networks, demonstrating the diversity is not due to the different initialization parameters or training sequences (Figure S8F).

Many of the same subpopulation motifs observed in the experiment are also present in the microcircuit model. Although our overall cell coding scores are larger, we observe subpopulations that have very little coding to any feature or can be coded to one or more features (STAR Methods). Several clusters of cells have weak coding to images in the familiar session only to gain said coding in later sessions and vice versa. As in the experiment,

since the VIP population responds more strongly to image changes in the novel session, its average image coding score increases relative to that in the familiar and novel-plus sessions (Figure S8C, right). Nevertheless, there are many features of the experimental VIP subpopulations that we do not observe in our model: solely novel-plus image coded clusters, a clear overrepresentation of novel image coding, and significantly less diversity in omission coding.

The FMS mechanisms we have inserted into our model evidently affect cell diversity in addition to the mean responses. Since the familiar and novel stimulus trains are statistically identical, without the FMS modulations, the coding scores would be distributed across the two sessions evenly. That is, the fact that the excitatory to VIP synapses are not as familiarity suppressed by FMS<sub>A</sub> causes many cells to become image-driven in the novel session. The effects of the FMS modulations can also be seen in the substantial difference in coding scores between the novel and novel-plus sessions. Since these sessions are driven by the same stimuli trains, there should be no statistical difference in coding scores without changes in connection properties due to modulations. There are several qualitative features of the experimental data we do not capture that are outside the scope of this work, e.g., clusters with experience-dependent omission coding. In the discussion, we comment on how additional FMS mechanisms could be added to the model to produce such effects.

### Cell-specific synapse properties strongly influence cell coding

Having observed a diversity in VIP cell coding in the cortical circuit model, we analyzed what cell-specific network properties may be responsible for the heterogeneous coding. The pointlike neurons in our model primarily differ in their connectivity properties and how said connectivity is acted upon by the FMS mechanisms. To determine what differences individual VIP neurons may have that explain their diverse coding scores. we take many cell-specific network properties and see how well each of these correlates with either individual cell coding scores or cluster-averaged coding scores (STAR Methods). For example, if we plot the cluster-averaged novel image coding scores as a function of the total synaptic input a VIP cell receives from the excitatory to VIP synapses during the image presentations of a novel session, we observe a statistically significant positive correlation (p = 0.002; Figure 6E). Since the amount of input a given VIP cell receives is influenced by the heterogeneous synaptic connections and image encoding signals, this quantity is different for each VIP cell and is evidently reflected in the cell's image coding scores. We repeat this procedure across 16 cell-specific network properties of our cortical circuit model and all nine coding score values (Figure S9). Additionally, we analyze how said network properties vary as a function of each VIP cluster (Figure S10).

Across all the cell-specific properties we consider, the strength of the modulated excitatory to VIP synapses mentioned above has the largest correlation with the novel image coding scores (Figures 6F and S9). Similarly, if a given VIP cell happens to have strong input from excitatory synapses during familiar image presentations, it tends to have a larger familiar image coding



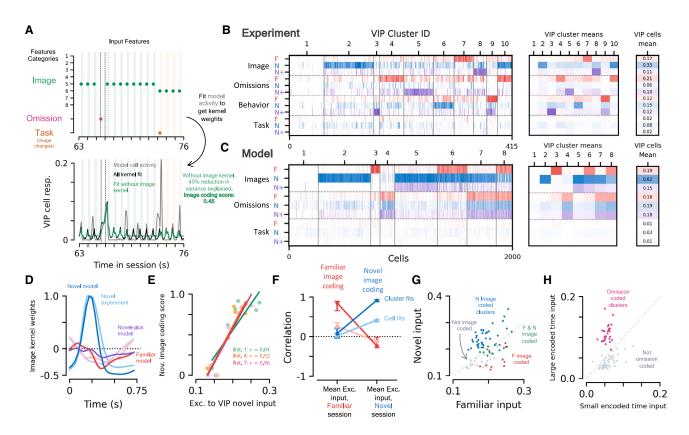


Figure 6. Coding diversity in the cortical circuit model

(A) Kernel regression model: input features are convolved with learned kernels and summed to predict individual cell activity in a kernel regression model (STAR Methods). Top: The cortical circuit model is driven by three primary feature categories: image presentations, omissions, and task-relevant image changes. The experimental data share these same feature categories as well as behavioral effects.<sup>24</sup> Bottom: Feature categories are removed one at a time and the kernel regression model is refit to determine each feature's contribution to a cell's activity, summarized by the category's coding score.

(B and C) VIP coding scores.

(B) Experimental results. <sup>24</sup> Left shows clustered VIP coding scores for all four feature categories across the familiar, novel, and novel-plus sessions; middle shows cluster means; and right shows mean over all cells (STAR Methods).

(C) Same as (B) for the cortical circuit model.

(D) Normalized mean image feature kernels for the familiar (red), novel (blue), and novel-plus sessions (purple). Dark lines show fits to model data, light lines to experimental data.<sup>24</sup>

(E-H) Cell-specific cortical circuit properties influence VIP cell coding.

(E) Cluster-averaged novel image coding scores as a function of the mean input each cell receives from the excitatory population during the novel session. Colors show three different initializations, dots are cluster-averaged values, lines are linear regression fits.

(F) Change in correlation of familiar (red) and novel (blue) image coding scores with network properties. Correlation over all VIP cells (light) and cluster-averaged values (dark). Dots are median, error bars are Q1 to Q3 over initializations.

(G) Cluster-averaged network properties across 10 different initializations. Colored dots (red, blue, green) show clusters with distinct types of image coding, gray dots are clusters not coded to images. Gray dotted line shows equal familiar and novel input.

(H) Network properties for omission-coded clusters (pink) versus omission-agnostic clusters (gray). See also Figures S8-S10.

score (Figure 6F). As might be expected, we see the amount of input a cell receives in a familiar session has relatively little correlation with the novel image coding score and vice versa. These coding score correlations can be used to determine properties that cells in a given cluster may share. For example, the cluster-averaged excitatory input from the familiar and novel sets roughly separates clusters into those that are familiar and/or novel image coded across network initializations (Figure 6G).

Interestingly, the image coding scores for the familiar and novel-plus sessions also have a strong correlation with how much stimulus history input the cells receive (Figure S9). That

is, in the absence of significant within-layer excitatory input, a result of *FMS*<sub>A</sub>'s suppression, the VIP cells that happen to have strong synapses from stimulus history sources are the most coded to images during the familiar and novel-plus sessions. An important distinction about the stimulus history input the VIP cells receive is that it is image-presentation-correlated but is uniform across all images, up to effects from noise. This lines up with the experimental observation that VIP cells increase image decodability during the novel session, but do not see a significant increase in decodability during presentations of the familiar and novel-plus sessions.<sup>24</sup> Note that since these inputs



Table 1. Relative cell counts and inter-population connection strengths of the various cell populations

	Rel. count	Post Exc.	Post SST	Post VIP	Post PV
Pre Exc. (L2/3)	27.35	0.105	0.750	1.000	0.908
Pre SST	1.00	- 0.081	$-\ 0.024$	$-\ 0.356$	- 0.060
Pre VIP	1.67	- 0.008	- 0.227	- 0.020	- 0.004
Pre PV	1.38	- 0.262	- 0.107	N/A	- 0.353

Second column shows relative cell counts for Layer 2/3, and it is assumed that 70% of Htr3a cells are VIP. <sup>18</sup> Third through sixth columns show population connection strengths for various pre- and postsynaptic combinations. See STAR Methods for details on how relative population connection strengths are computed. Note, although PV neurons are not included in our microcircuit model, they are included in this table for completeness. The "N/A" entry corresponds to synapses for which there are no data in Campagnola et al.<sup>20</sup>

do not change across images or sessions significantly, this could alternatively be interpreted as cells coded to familiar/novel-plus images are those closer to firing thresholds during all image presentations.

Although the model's omission coding exhibits far less experience-level diversity than what is observed in the data, we see how strongly connected a given VIP cell is to those cells excited during omission-encoded times, i.e., >500 ms, strongly influences the omission coding (Figure S9). Additionally, there is a slight negative correlation of the omission coding scores with the amount of input a VIP cell receives during non-omission-encoded times, i.e.,  $\leq 500$  ms. Any VIP cell that receives stimulus history input during the training session will have their presynaptic connections suppressed by  $FMS_{AH}$ , meaning said inputs will not influence the omission response. We see the clusters that are omission-coded tend to have larger input during large encoded times compared with those that do not (Figure 6H).

#### **DISCUSSION**

In this work we have introduced FMSs, a simple familiaritydetection mechanism that relies solely on local, unsupervised synaptic modulations to encode exposure to past stimuli. The individual modulations of the FMS mechanisms evolve via wellcharacterized dynamics: Hebbian or anti-Hebbian associative or presynaptic only dependence. We first investigated the basic properties of the FMS mechanism in a simple feedforward network, what we refer to as the FMSN. There we saw that, unlike several other familiarity-detection models, FMSs can detect novelty in a single forward pass, which is supported by evidence showing such stimulus distinctions can occur rapidly in humans. 61,62 We then demonstrated the generalizability of FMSs by modeling three distinct novelty novelty effects recently observed in a cortical disinhibitory circuit containing excitatory, VIP, and SST neurons. The connectivity of the cortical circuit model we develop is constrained by an aggregate of recent experimental results. 17,18,20 The three separate VIP novelty effects were reproduced in a continual learning setting with experimentally realistic stimulus sequences. Finally, due largely to the

**Table 2. Threshold voltage estimates** 

Cell type	Cre-line	Cell count	$V_{\text{threshold}}$ (mV, mean $\pm$ SD)
Exc.	Cux2-CreERT2	81	$-47.4 \pm 6.0$
SST	Sst-IRES-Cre	123	$-\ 41.0 \pm 7.6$
VIP	Vip-IRES-Cre	97	$-\ 47.2 \pm 8.7$
PV	Pvalb-IRES-Cre	217	$-\   35.0\pm 8.1$

Values are computed across an entire session sweep and then averaged across a given specimen ID and then Cre-line. From the Allen Cell Types Database, found at https://celltypes.brain-map.org/data.<sup>19</sup>

modulations that change the network's response over time, we found significant cell subpopulation diversity in our model, reproducing results that have been recently highlighted in the cortical disinhibitory circuit.<sup>24</sup>

Although we do not explicitly model all the novelty effects observed in the experiment (see limitations below), given the results we have seen here we can speculate on how the generalizability of FMSs would allow them to model such effects. The inhibition from the VIP population alone is not enough to produce the change in SST response seen across the familiar and novel sessions (Figure S6). A slow strengthening FMS on the the excitatory to SST synapses, i.e., similar FMSA with positive learning rate, would result in a larger SST population response in the familiar session, similar to what is observed in the experimental data. 23,24 A fast FMS, analogous to FMSP, on synapses that carry the bottom-up signal into the excitatory population could drive the observed increased excitatory response to image changes, very similar to what we produced in the VIP cells. We also do not attempt to model the suppressed omission ramping that is observed in the VIP population during the novel session.<sup>2</sup> A population of VIP neurons that have increased mean activity in a novel session, via an FMS mechanism similar to FMSA, could act to gate the ramping signal. If this population inhibits the excitatory neurons that produce the stimulus history input signal, it would lead to an overall smaller input and thus a smaller ramping only during the novel session. Furthermore, if this VIP population activity was different across experience levels, like the VIP population in our model and experiment, the inhibition to the history input would change between sessions, potentially leading to experience-level-dependent omission coding that has been observed<sup>24</sup> (Figure 6B). We also note that it may be possible to observe the omission ramping response using a fast, STSPlike FMS in place of the LTD-like FMS<sub>AH</sub>, i.e., the  $\leq$ 500-ms encoded times become familiar on a timescale of seconds rather than hours/day.

We cannot rule out that part of the novelty responses observed in the cortical circuit may be driven by signals outside the visual cortex (though see Keller et al. <sup>51</sup>). Regardless, since it has been confirmed that it is not the specific stimuli in the familiar and novel image sets that evoke the novelty responses, <sup>24</sup> they must be generated somewhere in the brain and we have demonstrated a plausible plasticity mechanism that could produce said responses in cell populations.

As mentioned above, the FMSs' simplicity, generality, and effectiveness in producing novelty effects makes them an ideal

#### **Article**



candidate for studying plasticity in future work. For example, modeling projects could scan over various FMS configurations and parameterizations within the cortical circuit to match the experimentally observed novelty responses and see if the resulting fits match or may constrain experimentally observed plasticity. On the experimental end, understanding how neuronal responses change throughout all of training would allow us to further characterize the types of plasticity that modulate experience-dependent activity. From our cortical circuit model, predictions about how connectivity and plasticity influence the observed cell subpopulation diversity could be tested by pairing physiological and learning studies together. Finally, the FMS's similarity to Hopfield networks means modern extensions to said networks could be used in the FMSs to increase their effectiveness and/or capacity. 63

Altogether, the effectiveness of FMSs highlights the role simple modulations within large synapse populations may have in shaping neuronal responses to stimuli. We demonstrated the FMS relies on no specialized training and testing schedules and requires no carefully placed excitatory or inhibitory connections to operate. Our cortical circuit demonstrates two important features of the FMS mechanism: (1) its ability to operate with several distinct types of modulations on the same synapses, even at significantly different timescales, as well as (2) the ability to have multiple inputs with distinct synaptic dynamics influencing a single cell population. Crucially, these mechanisms allowed us to model the novelty effects that have been observed to occur over significantly different timescales (seconds to days) and from different sets of information on the same set of cells.<sup>24</sup> Other than a few parameters adjusted at initialization to ensure realistic input and firing rates (< 10), the cortical circuit's response is driven by the FMS mechanisms, themselves only containing one free parameter apiece: their learning rate/size of modulations. It is surprising to the authors that what seem like complicated novelty responses can be captured by such straightforward modulation mechanisms. yet speaks toward the influence that simple synaptic changes can have on our brains.

#### Limitations of the study

In the cortical circuit model, we specifically focused on reproducing three distinct novelty effects seen in the VIP population responses. The goal of this modeling study was to demonstrate how several FMS mechanisms could be combined to produce various novelty responses within a single cell population. With this in mind, the specific choices of where we have added FMS mechanisms and their parameterization are not tested against other potential plasticity configurations that could also give rise to the novelty responses. With the addition of plasticity elsewhere in the cortical circuit and/or further neuronal contributions to the model that we have neglected, plasticity that weakens the synapses feeding into the VIP neurons may not be necessary to produce novelty effects.<sup>22</sup> We leave an extensive study of how various plasticity configurations could give rise to all the novelty effects observed within the cortical circuit for future work (see above).<sup>24</sup> We do not attempt to model the heterogeneous learning rates across the same synaptic population that has been observed in short-term plasticity.<sup>20</sup> However, since synapses continue to either be strengthening or weakening on average, we do not believe this will affect the mean population activity significantly. We also neglect effects of neuron scaling, e.g., adjusting those synapses not strengthened/weakened via, say, heterosynaptic LTP or LTD.

#### **STAR**\*METHODS

Detailed methods are provided in the online version of this paper and include the following:

- KEY RESOURCES TABLE
- RESOURCE AVAILABILITY
  - Lead contact
  - Materials availability
  - o Data and code availability
- METHOD DETAILS
  - Familiarity modulated synapses (FMSs)
  - o Familiarity modulated synapse network (FMSN)
  - Response sparsity adjustment
  - o Decoding accuracy and dimensionality
  - Strengthening versus weakening
  - Cortical microcircuit network
  - Cortical microcircuit connectivity
  - Response sparsity adjustment
  - $\,\circ\,$  Response rates and variance explained
  - Noise injection and matching baseline responses
  - o Response smoothing
  - o Familiarity-novelty task
  - Experimental image change detection task
  - Model stimulus change task
  - Present stimulus
  - o Stimulus history information
  - o Time-correlated noise
  - o Training schedule details
  - Analytical intuition and additive modulations
- QUANTIFICATION AND STATISTICAL ANALYSIS
  - $\ \, \circ \ \, \text{Fitting FMS learning rates} \\$
  - $\bigcirc \ \, \text{Cell subpopulation analysis}$
  - Experimental data
  - Model fitting
  - Coding scores
  - Cell clustering
  - Additional Figure details

#### SUPPLEMENTAL INFORMATION

Supplemental information can be found online at https://doi.org/10.1016/j.ceirep.2024.114188.

#### **ACKNOWLEDGMENTS**

We thank Anton Arkhipov, Stefan Berteau, Darrell Haufler, Shinya Ito, Lukasz Kusmierz, Zhixin Lu, Alex Piet, and Iryna Yavorska for feedback on this paper. K.A. and S.M. have been in part supported by NSF 2223725, NIH R01EB029813, and RF1DA055669 grants. We also wish to thank the Allen Institute for Brain Science founder, Paul G. Allen, for his vision, encouragement, and support.

#### **AUTHOR CONTRIBUTIONS**

Conceptualization, K.A., M.G., and S.M.; methodology, K.A., L.C., and S.M.; software, K.A.; formal analysis, K.A.; investigation, K.A.; writing – original draft, K.A.; writing – review & editing, K.A., L.C., M.G., S.O., and S.M.; visualization, K.A.; supervision, M.G., S.O., and S.M.; funding acquisition, S.M.





#### **DECLARATION OF INTERESTS**

The authors declare no competing interests.

Received: November 6, 2023 Revised: February 9, 2024 Accepted: April 17, 2024 Published: May 6, 2024

#### REFERENCES

- Manahan-Vaughan, D., and Braunewell, K.-H. (1999). Novelty acquisition is associated with induction of hippocampal long-term depression. Proc. Natl. Acad. Sci. USA 96.15, 8739–8744.
- Ranganath, C., and Rainer, G. (2003). Neural mechanisms for detecting and remembering novel events. Nat. Rev. Neurosci. 4, 193–202.
- Schomaker, J., and Meeter, M. (2015). Short-and long-lasting consequences of novelty, deviance and surprise on brain and cognition. Neurosci. Biobehav. Rev. 55, 268–279.
- Jaegle, A., Mehrpour, V., and Rust, N. (2019). Visual novelty, curiosity, and intrinsic reward in machine learning and the brain. Curr. Opin. Neurobiol. 58, 167–174.
- Rust, N.C., and Cohen, M.R. (2022). Priority coding in the visual system. Nat. Rev. Neurosci. 23, 376–388.
- Li, L., Miller, E.K., and Desimone, R. (1993). The representation of stimulus familiarity in anterior inferior temporal cortex. J. Neurophysiol. 69, 1918–1929
- Xiang, J.-Z., and Brown, M.W. (1998). Differential neuronal encoding of novelty, familiarity and recency in regions of the anterior temporal lobe. Neuropharmacology 37, 657–676.
- Zhang, S., Xu, M., Kamigaki, T., Hoang Do, J.P., Chang, W.-C., Jenvay, S., Miyamichi, K., Luo, L., and Dan, Y. (2014). Long-range and local circuits for top-down modulation of visual cortex processing. Science 345, 660–665.
- Meyer, T., and Rust, N.C. (2018). Single-exposure visual memory judgments are reflected in inferotemporal cortex. Elife 7, e32259.
- Courchesne, E., Hillyard, S.A., and Galambos, R. (1975). Stimulus novelty, task relevance and the visual evoked potential in man. Electroencephalogr. Clin. Neurophysiol. 39, 131–143.
- Daffner, K.R., Mesulam, M.M., Scinto, L.F., Calvo, V., Faust, R., and Holcomb, P.J. (2000). An electrophysiological index of stimulus unfamiliarity. Psychophysiology 37, 737–747.
- Hawco, C., and Lepage, M. (2014). Overlapping patterns of neural activity for different forms of novelty in fMRI. Front. Hum. Neurosci. 8, 699.
- Bunzeck, N., and Düzel, E. (2006). Absolute coding of stimulus novelty in the human substantia nigra/VTA. Neuron 51, 369–379.
- Tulving, E., and Schacter, D.L. (1990). Priming and human memory systems. Science 247, 301–306.
- Polich, J., and Comerchero, M.D. (2003). P3a from visual stimuli: typicality, task, and topography. Brain Topogr. 15, 141–152.
- Braga, A., and Schönwiesner, M. (2022). Neural substrates and models of omission responses and predictive processes. Front. Neural Circ. 16, 799581
- Lee, S.H., Hjerling-Leffler, J., Zagha, E., Fishell, G., and Rudy, B. (2010).
   The largest group of superficial neocortical GABAergic interneurons expresses ionotropic serotonin receptors. J. Neurosci. 30, 16796–16808.
- Billeh, Y.N., Cai, B., Gratiy, S.L., Dai, K., Iyer, R., Gouwens, N.W., Abbasi-Asl, R., Jia, X., Siegle, J.H., Olsen, S.R., et al. (2020). Systematic integration of structural and functional data into multi-scale models of mouse primary visual cortex. Neuron 106, 388–403.e18.
- Gouwens, N.W., Sorensen, S.A., Berg, J., Lee, C., Jarsky, T., Ting, J., Sunkin, S.M., Feng, D., Anastassiou, C.A., Barkan, E., et al. (2019). Classification of electrophysiological and morphological neuron types in the mouse visual cortex. Nat. Neurosci. 22, 1182–1195.

- Campagnola, L., Seeman, S.C., Chartrand, T., Kim, L., Hoggarth, A., Gamlin, C., Ito, S., Trinh, J., Davoudian, P., Radaelli, C., et al. (2022). Local connectivity and synaptic dynamics in mouse and human neocortex. Science 375, eabj5861.
- Pfeffer, C.K., Xue, M., He, M., Huang, Z.J., and Scanziani, M. (2013). Inhibition of inhibition in visual cortex: the logic of connections between molecularly distinct interneurons. Nat. Neurosci. 16, 1068–1076.
- Schulz, A., Miehl, C., Berry, M.J., 2nd, and Gjorgjieva, J. (2021). The generation of cortical novelty responses through inhibitory plasticity. Elife 10, e65309.
- Garrett, M., Manavi, S., Roll, K., Ollerenshaw, D.R., Groblewski, P.A., Ponvert, N.D., Kiggins, J.T., Casal, L., Mace, K., Williford, A., et al. (2020). Experience shapes activity dynamics and stimulus coding of VIP inhibitory cells. Elife 9, e50340.
- Garrett, M., Groblewski, P., Piet, A., Ollerenshaw, D., Najafi, F., Yavorska, I., Amster, A., Bennett, C., Buice, M., Caldejon, S., et al. (2023). Stimulus novelty uncovers coding diversity in visual cortical circuits. Preprint at bio-Rxiv. https://doi.org/10.1101/2023.02.14.528085.
- Tremblay, R., Lee, S., and Rudy, B. (2016). GABAergic interneurons in the neocortex: from cellular properties to circuits. Neuron 91, 260–292.
- Zeng, H., and Sanes, J.R. (2017). Neuronal cell-type classification: challenges, opportunities and the path forward. Nat. Rev. Neurosci. 18, 530–546.
- Tasic, B., Yao, Z., Graybuck, L.T., Smith, K.A., Nguyen, T.N., Bertagnolli, D., Goldy, J., Garren, E., Economo, M.N., Viswanathan, S., et al. (2018). Shared and distinct transcriptomic cell types across neocortical areas. Nature 563, 72–78.
- Gouwens, N.W., Sorensen, S.A., Baftizadeh, F., Budzillo, A., Lee, B.R., Jarsky, T., Alfiler, L., Baker, K., Barkan, E., Berry, K., et al. (2020). Integrated morphoelectric and transcriptomic classification of cortical GABAergic cells. Cell 183, 935–953.e19.
- Misha, V.T., and Markram, H. (1997). The neural code between neocortical pyramidal neurons depends on neurotransmitter release probability. Proc. Natl. Acad. Sci. USA 94.2, 719–723.
- Montgomery, D.P., Hayden, D.J., Chaloner, F.A., Cooke, S.F., and Bear, M.F. (2021). Stimulus-selective response plasticity in primary visual cortex: progress and puzzles. Front. Neural Circ. 15, 815554.
- Hopfield, J.J. (1982). Neural networks and physical systems with emergent collective computational abilities. Proc. Natl. Acad. Sci. USA 79, 2554–2558.
- Bogacz, R., Brown, M.W., and Giraud-Carrier, C. (2001). Model of familiarity discrimination in the perirhinal cortex. J. Comput. Neurosci. 10, 5–23.
- Bogacz, R., and Brown, M.W. (2002). The restricted influence of sparseness of coding on the capacity of familiarity discrimination networks. Netw. Comput. Neural Syst. 13, 457–485.
- Bogacz, R., and Brown, M.W. (2003). Comparison of computational models of familiarity discrimination in the perirhinal cortex. Hippocampus 13, 494–524.
- Tyulmankov, D., Yang, G.R., and Abbott, L.F. (2022). Meta-learning synaptic plasticity and memory addressing for continual familiarity detection. Neuron 110, 544–557.e8.
- Aitken, K., and Mihalas, S. (2022). Neural Population Dynamics of Computing with Synaptic Modulations. Preprint at bioRxiv. https://doi. org/10.1101/2022.06.27.497776.
- Collingridge, G.L., Peineau, S., Howland, J.G., and Wang, Y.T. (2010).
   Long-term depression in the CNS. Nat. Rev. Neurosci. 11, 459–473.
- Nicoll, R.A. (2017). A brief history of long-term potentiation. Neuron 93, 281–290.
- McFarlan, A.R., Chou, C.Y.C., Watanabe, A., Cherepacha, N., Haddad, M., Owens, H., and Sjöström, P.J. (2023). The plasticitome of cortical interneurons. Nat. Rev. Neurosci. 24, 80–97.

#### **Article**



- Lim, S., McKee, J.L., Woloszyn, L., Amit, Y., Freedman, D.J., Sheinberg, D.L., and Brunel, N. (2015). Inferring learning rules from distributions of firing rates in cortical neurons. Nat. Neurosci. 18, 1804–1810.
- Sukbin, L. (2019). Mechanisms underlying sharpening of visual response dynamics with familiarity. Elife 8, e44098.
- Hu, B., Garrett, M.E., Groblewski, P.A., Ollerenshaw, D.R., Shang, J., Roll, K., Manavi, S., Koch, C., Olsen, S.R., and Mihalas, S. (2021). Adaptation supports short-term memory in a visual change detection task. PLoS Comput. Biol. 17, e1009246.
- Vogels, R. (2016). Sources of adaptation of inferior temporal cortical responses. Cortex 80, 185–195.
- 44. Rust, N.C., and Palmer, S.E. (2021). Remembering the past to see the future. Annu. Rev. Vis. Sci. 7, 349–365.
- Yakovlev, V., Amit, D.J., Romani, S., and Hochstein, S. (2008). Universal memory mechanism for familiarity recognition and identification. J. Neurosci. 28, 239–248.
- Ba, J., Hinton, G.E., Mnih, V., Leibo, J.Z., and Ionescu, C. (2016). Using fast weights to attend to the recent past. Adv. Neural Inf. Process. Syst. 29.
- Potjans, T.C., and Diesmann, M. (2014). The cell-type specific cortical microcircuit: relating structure and activity in a full-scale spiking network model. Cerebr. Cortex 24, 785–806.
- Fu, Y., Tucciarone, J.M., Espinosa, J.S., Sheng, N., Darcy, D.P., Nicoll, R.A., Huang, Z.J., and Stryker, M.P. (2014). A cortical circuit for gain control by behavioral state. Cell 156, 1139–1152.
- Jung, H.L., Koch, C., and Mihalas, S. (2017). A computational analysis of the function of three inhibitory cell types in contextual visual processing. Front. Comput. Neurosci. 11, 28.
- Wilmes, K.A., and Clopath, C. (2019). Inhibitory microcircuits for top-down plasticity of sensory representations. Nat. Commun. 10, 5055.
- Keller, A.J., Dipoppa, M., Roth, M.M., Caudill, M.S., Ingrosso, A., Miller, K.D., and Scanziani, M. (2020). A disinhibitory circuit for contextual modulation in primary visual cortex. Neuron 108, 1181–1193.e8.
- 52. Hertäg, L., and Sprekeler, H. (2020). Learning prediction error neurons in a canonical interneuron circuit. Elife 9, e57541.
- Oldenburg, I.A., Hendricks, W.D., Handy, G., Shamardani, K., Bounds, H.A., Doiron, B., and Adesnik, H. (2024). The logic of recurrent circuits in the primary visual cortex. Nat. Neurosci. 27, 137–147. https://doi.org/10. 1038/s41593-023-01510-5.
- Hertäg, L., and Clopath, C. (2022). Prediction-error neurons in circuits with multiple neuron types: Formation, refinement, and functional implications. Proceedings of the National Academy of Sciences USA 119.13, e2115699119.
- Hofer, S.B., Ko, H., Pichler, B., Vogelstein, J., Ros, H., Zeng, H., Mrsic-Flogel, T.D., Lein, E., and Lesica, N.A. (2011). Differential connectivity and response dynamics of excitatory and inhibitory neurons in visual cortex. Nat. Neurosci. 14, 1045–1052.

- Chettih, S.N., and Harvey, C.D. (2019). Single-neuron perturbations reveal feature-specific competition in V1. Nature 567, 334–340.
- Pillow, J.W., Shlens, J., Paninski, L., Sher, A., Litke, A.M., Chichilnisky, E.J., and Simoncelli, E.P. (2008). Spatio-temporal correlations and visual signalling in a complete neuronal population. Nature 454, 995–999.
- Simon, M., Matthew, T.K., Juavinett, A.L., Gluf, S., and Churchland, A.K. (2019). Single-trial neural dynamics are dominated by richly varied movements. Nat. Neurosci. 22, 1677–1686.
- Engelhard, B., Finkelstein, J., Cox, J., Fleming, W., Jang, H.J., Ornelas, S., Koay, S.A., Thiberge, S.Y., Daw, N.D., Tank, D.W., and Witten, I.B. (2019). Specialized coding of sensory, motor and cognitive variables in VTA dopamine neurons. Nature 570, 509–513.
- Steinmetz, N.A., Zatka-Haas, P., Carandini, M., and Harris, K.D. (2019).
   Distributed coding of choice, action and engagement across the mouse brain. Nature 576, 266–273.
- Seeck, M., Michel, C.M., Mainwaring, N., Cosgrove, R., Blume, H., Ives, J., Landis, T., and Schomer, D.L. (1997). Evidence for rapid face recognition from human scalp and intracranial electrodes. Neuroreport 8, 2749–2754.
- Hintzman, D.L., Caulton, D.A., and Levitin, D.J. (1998). Retrieval dynamics in recognition and list discrimination: Further evidence of separate processes of familiarity and recall. Mem. Cognit. 26, 449–462.
- Krotov, D., and Hopfield, J.J. (2016). Dense associative memory for pattern recognition. Adv. Neural Inf. Process. Syst. 29.
- Cho, K., Kemp, N., Noel, J., Aggleton, J.P., Brown, M.W., and Bashir, Z.I. (2000). A new form of long-term depression in the perirhinal cortex. Nat. Neurosci. 3, 150–156.
- Sjöström, P.J., and Häusser, M. (2006). A cooperative switch determines the sign of synaptic plasticity in distal dendrites of neocortical pyramidal neurons. Neuron 51, 227–238.
- 66. Purves, D., Augustine, G., Fitzpatrick, D., Katz, L., LaMantia, A., McNamara, J., and Williams, S. (2001). Neuroscience 2nd edition. sunderland (ma) sinauer associates. In Types of Eye Movements and Their Functions.
- Zhu, J.J., and Lo, F.S. (1999). Three GABA receptor-mediated postsynaptic potentials in interneurons in the rat lateral geniculate nucleus. J. Neurosci. 19, 5721–5730.
- MacKenzie, G., and Maguire, J. (2015). Chronic stress shifts the GABA reversal potential in the hippocampus and increases seizure susceptibility. Epilepsy Res. 109, 13–27.
- Siegle, J.H., Jia, X., Durand, S., Gale, S., Bennett, C., Graddis, N., Heller, G., Ramirez, T.K., Choi, H., Luviano, J.A., et al. (2021). Survey of spiking in the mouse visual system reveals functional hierarchy. Nature 592, 86–92.
- 70. Ernst, H.W. (1978). De subtilitate tactu. In The Sense of Touch.
- Stanislas, D. (2003). The neural basis of the Weber–Fechner law: a logarithmic mental number line. Trends Cognit. Sci. 7.4, 145–147.
- Akre, K.L., and Johnsen, S. (2014). Psychophysics and the evolution of behavior. Trends Ecol. Evol. 29, 291–300.





#### **STAR**\*METHODS

#### **KEY RESOURCES TABLE**

REAGENT or RESOURCE	SOURCE	IDENTIFIER
Deposited data		
Image change detection data	Garrett et al., 2023 <sup>24</sup>	https://portal.brain-map.org/circuits-behavior/visual-behavior-2p
Software and algorithms		
Code for models and data analysis	This paper	https://github.com/kaitken17/fms
		https://doi.org/10.5281/zenodo.10914528

#### **RESOURCE AVAILABILITY**

#### **Lead contact**

Further information and requests for resources should be directed to and will be fulfilled by the lead contact, Kyle Aitken (kyle.aitken@alleninstitute.org).

#### **Materials availability**

This study did not generate new unique reagents.

#### **Data and code availability**

- This paper analyzes existing, publicly available data. These accession numbers for the datasets are listed in the key resources table.
- All original code has been deposited at <a href="https://github.com/kaitken17/fms">https://github.com/kaitken17/fms</a> and is publicly available as of the date of publication.
   DOIs are listed in the key resources table.
- Any additional information required to reanalyze the data reported in this paper is available from the lead contact upon request.

#### **METHOD DETAILS**

Many quantities we consider throughout this work are sequence time-dependent, and their time dependence is generally denoted by a subscript t (and sometimes s), i.e.,  $x_t$ . Time is uniformly discretized in our setup, so the quantities  $x_t$  and  $x_{t+1}$  are separated by  $\Delta t$ . When unambiguous, we use I, J = 1, ..., d to index the neurons of the presynaptic layer, and i, j = 1, ..., n to index the neurons of the postsynaptic layer.

#### Familiarity modulated synapses (FMSs)

As in the main text, we take  $\mathbf{x}_t \in \mathbb{R}^d$  to represent the presynaptic (i.e., input) neuron population firing rates at time t and  $\mathbf{y}_t \in \mathbb{R}^n$  the postsynaptic (output) equivalent. The synaptic modulation matrix,  $\mathbf{M}_t \in \mathbb{R}^{n \times d}$ , represents changes to the network's connections induced by some general biological mechanism through unsupervised learning. To incorporate changes in synapses due to various modulation mechanisms, we allow the modulation matrix to change otherwise *fixed* synapses generally represented by some randomly initialized matrix  $\mathbf{W}$ .(An alternative form of the modulations considered in Refs. <sup>35,36</sup> uses an additive modulation, rather than the multiplicative one we consider here. Essentially all results used for the FMSN generalize to this form of modulations as well [Figure S11].) We consider two distinct update rules for the modulations in this work. The *associative* mechanism, Equation 2a, dependent upon both the pre- and postsynaptic firing rates, is given by

$$\mathbf{M}_{t+1} = \lambda \mathbf{M}_t + \eta \mathbf{y}_t \mathbf{x}_t^T.$$
 (Equation 3)

The simpler pre-only modulation update that is only dependent upon the presynaptic firing rates, Equation 2b, is

$$\mathbf{M}_{t+1} = \lambda \mathbf{M}_t + \eta \mathbf{1} \mathbf{x}_t^T / \sqrt{n},$$
 (Equation 4)

where  $\mathbf{1} \in \mathbb{R}^n$  is the 1s vector. In the above expressions,  $\eta \in \mathbb{R}$  is the learning rate of the modulations that controls the rate at which modulations are learned. Its sign determines the sign of the modulations and thus whether they *strengthen* or *weaken* the corresponding synapses. (Since we use the FMSs to model several distinct types of synapse modulations that have their own vocabulary for synapse changes (e.g., depression and facilitation for STSP versus depression and potentiation for more long-term effects), we use "strengthening" and "weakening" as a general terminology that applies across the individual mechanisms the FMSs may model.) The other parameter,  $0 < \lambda < 1$ , represents the gradual decay of changes to the weight matrix. Occasionally it will be useful to discuss



the decay timescale  $au_{\text{decay}}$ , which is related to the decay rate via  $\lambda=1-\Delta t/ au_{\text{decay}}$ . Throughout this work, at the beginning of training, the modulations are initialized to be zero, i.e.,  $\mathbf{M}_0=\mathbf{0}$ .

We distinguish neurons in the networks we consider between excitatory and inhibitory. A neuron that is excitatory is defined to have only positive weights leaving it so that it can only enhance the response of the postsynaptic neurons it feeds into. Similarly, an inhibitory neuron has all negative weights leaving it, ensuring it can only depress the postsynaptic neurons it feeds into. For this definition of excitation/inhibition to be meaningful, we also limit all firing rates in the network to be positive definite, including all input firing rates. Explicitly, this means the weights of our network are subject to the constraints

Neuron *I* Excitatory : 
$$W_{il} \ge 0$$
, (Equation 5a)

Neuron / Inhibitory : 
$$W_{ij} \le 0$$
. (Equation 5b)

for all i = 1,...,n. Importantly, we do not allow for modulations to change the sign of weights leaving a given neuron. Per our definition above, this ensures that an excitatory neuron can never inhibit and vice-versa. Explicitly, this means that

Neuron *I* Excitatory : 
$$W_{ii} + W_{ii}M_{ii,t} \ge 0$$
, (Equation 6a)

Neuron / Inhibitory : 
$$W_{il} + W_{il}M_{il,t} \le 0$$
, (Equation 6b)

for all i = 1, ..., n and all t.

Several experimental studies have investigated the amount of change a given neuron can undergo through mechanisms such as STSP or LTP/D. To bound the modulations to realistic values, we further restrict the modulations so they cannot enhance or depress weights beyond what has been observed in experimental settings. Since such changes can differ depending on the mechanism, we enforce different bounds for the associative and pre-only dependent modulations. Explicitly, we limit

$$M_{\rm b}^{\rm A,lower} \le M_{il.t} \le M_{\rm b}^{\rm A,upper},$$
 (Equation 7a)

$$M_{\rm b}^{\rm P,lower} \le M_{il.t} \le M_{\rm b}^{\rm P,upper},$$
 (Equation 7b)

We take  $M_b^{A,upper} = M_b^{P,upper} = 1.0$ , so that both types of modulations can at most double the strength of the corresponding synapses. Meanwhile, we take  $M_b^{A,lower} = -0.8$ , i.e., a synapse can at most be reduced to 20% of its original strength, similar to values observed in several long-term plasticity experiments. STSP has been recently observed to almost completely suppress certain synapses, so we take  $M_b^{P,lower} = -1.0$ . Note with these chosen values of  $M_b$ , the bounds are stricter than those of Equation 6, so in practice the enforcement of Equation 7 implies the bounds of Equation 6 are automatically met.

#### Familiarity modulated synapse network (FMSN)

As mentioned in the main text, it is useful to first understand the properties of FMSs in a simple setting. To this end, we investigate FMS properties in a simple two-layer neural network, what we call the familiarity modulated synapse network (FMSN). We take the input and output layers to have d and n neurons, respectively. The input and output layers are connected by weights, representing the synapses of a biological neural network. We will assume the synapses of the network have some underlying strength at initialization, that we denote by the weight matrix  $\mathbf{W} \in \mathbb{R}^{n \times d}$ . We take  $\mathbf{W}$  to be sparse such that its elements have magnitude

$$|W_{il}| = |w_{il}|b_{il}, w_{il} \sim \mathcal{N}(0, w^2), b_{il} \sim \text{Bernoulli}(p_W)$$
 (Equation 8)

where the parameters w and  $p_W$  determine the magnitude of the nonzero elements and the sparsity, respectively. The signs of the nonzero elements are determined by the input neuron cell type (see above).

Like the synapses in the brain, we allow the individual weights in our network to be modulated over time. We denote the modulations at time step t by  $\mathbf{M}_t \in \mathbb{R}^{n \times d}$ , i.e., the same size as the weight matrix  $\mathbf{W}$ . The combined weights and modulation matrix yield the output neuron preactiviation values.

$$\tilde{\mathbf{y}}_t = (\mathbf{W} + \mathbf{W} \odot \mathbf{M}_t) \mathbf{x}_t + \mathbf{b},$$
 (Equation 9)

where  $\mathbf{b} = b\mathbf{1}$  with  $\mathbf{1} \in \mathbb{R}^n$  is a uniform bias term that can represent neuron firing-rate thresholds as well as other network factors neglected in this simple model (see below). The parameter b is adjusted at initialization to ensure realistic response sparsity in the output population, but is otherwise fixed throughout training, see below. Finally, the output preactiviations are passed through a nonlinear function,  $\varphi(\,\cdot\,)$ , representing the output neurons' properties such as their firing rate threshold and maximum firing rate. Thus, the output population's activity at time t is given by

$$\mathbf{y}_t = \varphi(\tilde{\mathbf{y}}_t),$$
 (Equation 10)





where  $\varphi(\cdot)$  is applied piecewise. Throughout this work, we use a rectified tanh as the output neuron's nonlinearity,

$$\varphi(x) = \operatorname{ReTanh}(x) \equiv \begin{cases} \tanh(x) & \text{for } x \ge 0, \\ 0 & \text{for } x < 0. \end{cases}$$
 (Equation 11)

This activation was chosen since it has three desirable properties: (1) it is positive definite which was required for cell types in our model to make sense, (2) it is bounded, and (3) it is approximately linear for  $0 < x \ll 1$ , so small positive preactivation values become small activity.

The FMSN's behavior can be changed considerably by modifying the distribution of excitatory/inhibitory neurons in the input neuron population. Note that since the neuron types only influence the sign of the weights/synapses leaving the neurons, the cell types of the output neurons do not affect the FMSN's behavior. When we investigate the FMSN in the main text, for simplicity, we consider a setup that only has excitatory neurons in its input layer. The bias adjustment within the output population that we perform to ensure realistic firing rates helps balance excitation/inhibition in the output population that only explicitly receives excitatory input. We find that the bias term is always negative for the output population firing rates we consider, which can be thought of a uniform inhibitory input into the output neuron population that provides E-I balance.

In the Methods, to understand the effect of FMSs applied to inhibitory synapses, we also consider cases where input population of the FMSN has both excitatory and inhibitory neurons in it. In these cases, the unmodulated weight matrix, **W**, has some columns with all positive nonzero elements and some columns with all negative nonzero elements. We study a simple case where the FMS mechanism applies to synapses belonging to either the excitatory or inhibitory input neuron populations. Of course, it is also possible to have modulations acting on both populations, but we leave such investigations to future work. All other properties of the FMSN, including the firing rate adjustment, remain unchanged.

#### Response sparsity adjustment

With only excitatory synapses in the input later of the FMSN, all inputs to the output neuron population are positive and so, without any threshold/bias term, the output population would have a response rate of close to 100% for every possible stimulus. Of course, such a response is not realistic over an entire population of neurons in the visual cortex and such excitation should be balanced by inhibition to achieve realistic response sparsities. This could be achieved by introducing inhibitory neurons with appropriate synaptic strengths into the input layer, in which case output neurons would receive somewhat balanced levels of excitation and inhibition. Indeed, we study such a network when we want to consider modulations on inhibitory neurons. However, for the FMSN investigated in the main text we choose a simpler solution that we discuss here that generalizes to the method we use to balance the cortical microcircuit below.

The output neuron population's response rate can also be adjusted by changing the population's firing threshold/bias, **b** in Equation 9. In this work, we only consider  $\mathbf{b} = b\mathbf{1}$  with  $\mathbf{1} \in \mathbb{R}^n$  so that the point neurons we study only differ in connectivity and noise injection. For a given neuron, a negative b effectively acts as a uniform inhibition across all possible inputs. Across the entire output population, a negative bias allows the neurons to have more realistic response sparsities despite only receiving excitatory inputs.

To adjust b to get the desired response rate, we first draw a validation set of size 100 from the same distribution that generates the familiar and novel sets. This entire validation set is passed through the network at initialization with b = 0, with no adjustment to the modulation matrix  $\mathbf{M}$ . Given the known activation function, Equation 11, from the validation set's preactivation values the bias needed to have the desired response sparsity across the validation set can be exactly computed. In short, all preactivation values (across the stimuli of the validation set and output neurons) are sorted by value, and the bias is chosen such that the desired percentage of these values are above 0. Since the familiar and novel sets are drawn from the same distribution as the validation set, this yields a similar response rate over said sets without being directly fit to them. Note this procedure means that the network has the desired response sparsities at initialization, but the induced modulations during training can change the response rate of the novel and familiar sets (Figure S2A).

#### **Decoding accuracy and dimensionality**

For the associative weakening FMSN example considering in the main text, the change in output activity of the familiar set significantly affects the decodability of the output signal. Post-training, decodability of stimulus identity within the familiar set is significantly lower  $(0.46 \pm 0.05)$  while that of the novel set is perfect  $(1.00 \pm 0.0)$ , mean  $\pm$  std). The difficultly in decodability is reflected in the effective dimensionality of each set's output activity: the novel outputs occupy a low-dimensional space  $(D = 6.3 \pm 1.5)$  while the familiar outputs are small enough that their signal is hard to distinguish from noise and thus the space they occupy is significantly higher dimensional  $(D = 48.5 \pm 7.1)$ , mean  $\pm$  std, Figures S2B and S2C). Both the above properties are a function of the amount of modulation within the network, so the decoability/dimensionality of the familiar and novel sets can vary significantly by, say, changing the modulation learning rate (Figures S2M and S2N).

To compute the decoding accuracy of the familiar and novel sets, we use the same input noise that is used during training (see above) to create 1000 noisy versions of each stimulus. Each noisy stimulus is then passed through the trained network, resulting in a total of 8000 output responses across the entire familiar or novel set. Said responses are then labeled by their index within their set and a linear SVC is used to decode them using 10-fold cross validation to compute test accuracies. Specifically, we use



sklearn.svm.LinearSVC with default parameters other than max\_iterations=1e5. The approximate dimensionality of the representations reported in the main text are computed using the participation ratio of the ratios of variance explained of the resulting PCA fits.

#### Strengthening versus weakening

Here we briefly discuss differences in the effectiveness of developing distinct responses to the familiar and novel sets in the FMSN from strengthening and weakening modulation mechanisms. We first study the FMSN in the setting we investigated in the main text, namely networks with only excitatory synapses in the input layer. As a result, a strengthening (or weakening) of the synapses results in a larger (or smaller) output neuron response.

There is a significant difference in the evolution of the  $\mathbf{M}_t$  as a function of training time in identical FMSNs that simply differ in the sign of  $\eta$  (Figure S3B). Notably,  $\mathbf{M}_t$  at roughly the same rate for the first 8 training steps, but for the weakening mechanism the growth of  $\mathbf{M}_t$  drops significantly after these first few times steps. This is a result of the output activity being smaller, which means the updates to  $\mathbf{M}_t$  are smaller (Figures S3C and S3D).

We can also see the effect the biological bounds have on the modulation growth. In particular, for the strengthening mechanism, the evolution of  $\mathbf{M}_t$  differs significantly with and without the biological bounds (Figures S3E and S3F). Notably, the weakening mechanism isn't as strongly affected by said bounds. The fact that weakening has been observed down to 20% for associative and only 200% for strengthening makes the former much more effective. This is a simple comparison of ratios of novel to familiar: for weakening, 20% leads to a ratio of 1/0.2 = 5 whereas strengthening to 200% leads to a ratio of only 2/1 = 2.

Another major difference between the strengthening and weakening mechanisms is the way the biologically motivated nonlinearity, Equation 11, acts on preactivation values. Since strengthening excitatory connections can only increase a neuron's output, but said output is bounded by the nonlinearity, eventually the strengthening yields diminishing returns in terms of how much a given output can change. Meanwhile, weakening excitatory connections can push a neuron below its firing threshold, completely cutting off a neuron's response. We see that the evolution of preactivation values is fairly comparable between the two mechanisms (Figure S3G).

Thus we have seen two major factors that cause the weakening and strengthening of excitatory synapses to differ: (1) a difference in the bounds of said changes from experiment and (2) an asymmetry in the FMSN of how larger/smaller output activity is handled through the neuron's nonlinearity as well as the modulation updates. Of course, we have only considered plasticity mechanisms on synapses belonging to excitatory neurons thus far. For inhibitory synapses, a strengthening (or weakening) of the synapses results in smaller (or larger) output neurons response, the opposite effect of the excitatory synapses. Thus we can investigate if the FMSN has different behavior when introducing inhibitory neurons in the input population and then make the inhibitory synapses FMSs.

For direct comparison to the FMSN with only excitatory synapses, we assume inhibitory plasticity obeys similar bounds to what we use for STSP and LTP/D effects. <sup>20,39</sup> We compare the behavior of an FMSN with both excitatory and inhibitory neurons in its input when either the excitatory or inhibitory neurons have a strengthening FMS mechanism on them (Figure S3H). Consistent with our findings above, we find that strengthening of the inhibitory neurons is more effects of at separating the novel and familiar distributions than a strengthening of excitatory neurons (Figures S3I and S3J). Since the bounds of the two FMS mechanisms are identical, this difference is caused by the neuron's nonlinear behavior discussed in point (2) above.

#### **Cortical microcircuit network**

A rough schematic of the cortical microcircuit network is shown in Figure 2C. The three primary populations we consider in the network are the excitatory, SST, and VIP neuron populations. We will index variables belonging to these three populations using p = E,S,V, respectively. Lastly, we also include an additional population of excitatory neurons that drive the stimulus history inputs into the VIP population and represent a subset of the top-down input into the cortical layer we explicitly model. We denote these additional excitatory neurons by the superscript 'hist'. We do not make any attempt to model behavioral effects related to the image change task, including the licking response to the task.<sup>42</sup>

We begin by introducing the cortical microcircuit without any FMS mechanisms added to its synapses. The preactivation response of the excitatory, VIP, and SST populations are respectively given by

$$\tilde{\mathbf{y}}_{t}^{E} = \mathbf{W}^{E,E} \mathbf{y}_{t-1}^{E} + \mathbf{W}^{E,V} \mathbf{y}_{t-1}^{V} + \mathbf{W}^{E,S} \mathbf{y}_{t-1}^{S} + \mathbf{b}^{E} + \mathbf{n}_{t}^{E} + \mathbf{x}_{t}^{stim},$$
 (Equation 12a)

$$\tilde{\mathbf{y}}_{t}^{S} = \mathbf{W}^{S,E} \mathbf{y}_{t}^{E} + \mathbf{W}^{S,V} \mathbf{y}_{t-1}^{V} + \mathbf{W}^{S,S} \mathbf{y}_{t-1}^{S} + \mathbf{b}^{S} + \mathbf{n}_{t}^{S}.$$
 (Equation 12b)

$$\tilde{\mathbf{y}}_{t}^{V} = \mathbf{W}^{V,E}\mathbf{y}_{t}^{E} + \mathbf{W}^{V,V}\mathbf{y}_{t}^{V}, + \mathbf{W}^{V,S}\mathbf{y}_{t}^{S} + \mathbf{b}^{V} + \mathbf{n}_{t}^{V} + \mathbf{W}^{V,hist}\mathbf{x}_{t}^{hist},$$
 (Equation 12c)

where  $\mathbf{W}^{p,p'}$  represents the synapses connecting presynaptic population p' to postsynaptic population p,  $\mathbf{b}^p$  is the bias vector of population p, and  $\mathbf{n}^p_t$  represents additional noise injection (see below). Note the three populations do not update in sync: at time t the excitatory population's activity is updated first, followed by the SST population, and then VIP. Asynchronous updates were found to help numerical stability. This order is also biologically motivated since the canonical input to layer 2/3 from layer 4 pyramidal neurons is much weaker to VIP and SST than pyramidal neurons.



All the preactivation responses pass through a nonlinearity,

$$\mathbf{y}_t^p = \varphi(\tilde{\mathbf{y}}_t^p), \text{ for } p = E, S, V,$$
 (Equation 13)

where  $\varphi(\cdot) = \text{ReTanh}(\cdot)$ , ensuring the rate remains positive definite, see Equation 11.

The excitatory and VIP populations both receive additional external input related to the stimulus change task. Specifically,  $\mathbf{x}_t^{\text{stim}}$  and  $\mathbf{x}_t^{\text{hist}}$  represent the present stimulus input and activity of the stimulus history excitatory neurons, respectively (see below for details). Note the activity of the excitatory neurons that drive the stimulus history signal is an input to the network, so we have denoted it by  $\mathbf{x}_t^{\text{hist}}$  rather than  $\mathbf{y}_{t-1}^{\text{hist}}$  for clarity. Both stimulus inputs,  $\mathbf{x}_t^{\text{stim}}$  and  $\mathbf{x}_t^{\text{hist}}$ , are fed directly into an excitatory population that in turn drives the rest of the circuit through sparse synapses. The primary difference between them is that the former drives an excitatory population inside the microcircuit while the latter drives an excitatory population that is assumed to reside in a higher cortical area that we do not explicitly model. The stimulus history excitatory neurons represent a subset of top-down information fed into the microcircuit and were chosen to feed into the VIP population since they are known to receive feedback input.<sup>20</sup>

The  $\mathbf{W}^{p,p'}$  are sparse matrices whose elements are also drawn according to Equation 8, i.e., in the same way as the weight matrix of the FMSN. Once again, the cell type of the presynaptic population p determines the sign of the nonzero  $\mathbf{W}^{p,p'}$  elements, and thus  $\mathbf{W}^{p,E} \geq 0$ ,  $\mathbf{W}^{p,S} \leq 0$ , and  $\mathbf{W}^{p,V} \leq 0$  for all p. For a given  $\mathbf{W}^{p,p'}$ , the sparsity of the synapses (i.e.,  $p_W$ ), magnitude of the nonzero elements (i.e., w), and relative number of cells in each population are all set by experimental literature,  $^{17,18,20}$  see below. In short, for all  $\mathbf{W}^{p,p'}$ , we take  $p_W$  to be the corresponding entry in Figure S4D and w to be the corresponding entry in Figure S4B, up to a multiplicative constant c. Importantly, c is the same for all  $\mathbf{W}^{p,p'}$ , and thus the relative connection strengths between populations are completely fixed by experimental results, up to changes from FMSs.

Since  $\mathbf{W}^{V,\text{hist}}$  represents an unknown subset of top-down excitatory to VIP connections, we simply set its sparsity equal to the within-layer excitatory to VIP connections. Its synaptic strength is set to a value comparable to the within-layer excitatory connections so that the omission ramping response has a comparable magnitude to image responses.

All three biases are parameterized similarly to the FMSN, i.e.,  $\mathbf{b}^{\rho} = b^{\rho} \mathbf{1}$  where  $\mathbf{1}$  is the all 1's vector in the corresponding space. Once again, the  $b^{\rho}$  are adjusted at initialization to ensure realistic response sparsities in all three neuron populations, see STAR Methods.

We add the following three FMS mechanisms onto synapses feeding into the VIP cells,

$$\mathbf{W}^{V,E} \to \mathbf{W}^{V,E} + \mathbf{W}^{V,E} \odot \left( \mathbf{M}_t^{(A)} + \mathbf{M}_t^{(P)} \right),$$
 (Equation 14a)

$$\mathbf{W}^{V,\text{hist}} \to \mathbf{W}^{V,\text{hist}} + \mathbf{W}^{V,\text{hist}} \odot \mathbf{M}_{t}^{(AH)},$$
 (Equation 14b)

where the superscripts in  $(\cdot)$  refer to distinct FMS mechanisms. Specifically, (A), (P), and (AH) respectively correspond to what we refer to as the  $FMS_A$ ,  $FMS_P$ , and  $FMS_{AH}$  mechanisms. Note we have added two FMS mechanisms to the same set of synapses, those going from the excitatory to the VIP population. When multiple sets of FMSs are present on the same synapses, we still enforce the cell-type bounds of Equation 6. The three distinct modulation correspond to the three novelty responses we aim to model. They are respectively subject to the following update expressions,

FMS<sub>A</sub>: 
$$\mathbf{M}_{t+1}^{(A)} = \lambda^{(A)} \mathbf{M}_{t}^{(A)} + \eta^{(A)} \mathbf{y}_{t}^{V} (\mathbf{y}_{t}^{E})^{T}$$
, (Equation 15a)

$$\mathsf{FMS}_{\mathsf{P}}: \mathbf{M}_{\mathsf{t}+1}^{(\mathsf{P})} \ = \ \lambda^{(\mathsf{P})} \mathbf{M}_{\mathsf{t}}^{(\mathsf{P})} + \eta^{(\mathsf{P})} \mathbf{1} \left( \mathbf{y}_{\mathsf{t}}^{\mathsf{E}} \right)^{\mathsf{T}} \bigg/ \sqrt{\mathsf{n}^{\mathsf{E}}}, \tag{Equation 15b}$$

$$\mathsf{FMS}_{\mathsf{AH}}: \mathbf{M}_{\mathsf{t+1}}^{(\mathsf{AH})} = \lambda^{(\mathsf{AH})} \mathbf{M}_{\mathsf{t}}^{(\mathsf{AH})} + \eta^{(\mathsf{AH})} \mathbf{y}_{\mathsf{t}}^{\mathsf{V}} (\mathbf{x}_{\mathsf{t}}^{\mathsf{hist}})^{\mathsf{T}}. \tag{Equation 15c}$$

Note that the updates are distinct, but are all of the same fundamental form we have used throughout this work, see Equation 2. That is, the associative updates are dependent upon both the pre- and postsynaptic firing rates of the populations they connect, while the pre-only is only dependent on the presynaptic firing rates since we want it to represent STSP-like modulations that occur at time-scales on the order of seconds. The three FMS mechanisms are subject to the corresponding bounds motivated from experiment discussed below Equation 7. In practice, during training, the  $FMS_A$  and  $FMS_P$  modulations rarely come close to saturating the bounds imposed by experiment, while the modulations of  $FMS_{AH}$  come close to their bounds at a much higher rate. For the exemplar network shown in Figure 5, for the modulation matrix terms corresponding to nonzero synapses of  $\mathbf{M}^{(A)}$ ,  $\mathbf{M}^{(P)}$ , and  $\mathbf{M}^{(AH)}$ , only 1.5%, 0.09%, and 42% come within 50% of their bound and 0%, 0%, and 23% come within 10% of their bound, respectively. Note for the slower modulation mechanisms, these rates were only calculated during roughly the last quarter of training time.

To adjust the parameters of the three FMSs shown above, we scan over learning rates of all three FMS mechanisms and determine which of these yields the best mean response fits, see below for additional details. We take  $\tau_{\rm decay}^{(A)} = 10^5$  seconds,  $\tau_{\rm decay}^{(P)} = 1$  second, and  $\tau_{\rm decay}^{(AH)} = 10^6$  seconds based on the timescale of the corresponding biological mechanisms we wish to match onto and experimental observations (though see STAR Methods for how these may change to expedite training).



#### **Cortical microcircuit connectivity**

In our model, the total strength of the collection of synapses from a given presynaptic population into a given postsynaptic neuron depends on three major factors.

- (1) **Connection probability:** The probability for a given synapse to exist between any two cells of the given pre- and postsynaptic type (Figure S4D).
- (2) Relative cell counts: Along with the mean connection probability, the number of presynaptic cells in a given population affects the mean number of inputs a given postsynaptic neuron receives. Since our network does not explicitly model true cell counts found in the visual cortex, this enters as a relative value that we can then normalize by some baseline number (Figure S4E).
- (3) **Synapse strength:** The strength *per synapse*, which we measure by the mean time-integrated postsynaptic potential (PSP), see below (Figures S4B and S4C).

When layer-specific experimental observations are available, we take the values given for L2/3 of the visual cortex. The three factors above directly affect both the population count in our microcircuit model as well as the explicit form of the  $\mathbf{W}^{p,p'}$ . It can also be helpful to track the mean population strength, a product of the three factors above

$$r_{\rho \leftarrow \rho'} = \frac{n^p}{n_{\text{baseline}}} \times P_{\rho \leftarrow \rho'}^{\text{conn.}} \times Z_{\rho \leftarrow \rho'}^{\text{PSP}}.$$
 (Equation 16)

In practice, we take  $n_{\text{baseline}} = n^{\text{S}}$ , which gives us the inter-population connection strengths shown in Table 1 and Figure S4A. Although they are not explicitly included in our microcircuit model, we include data for PV neurons here as well (see discussion for further details).

(1) Connection probability: The connection probabilities between pre- and postysnaptic populations are computed from the fully adjusted connection propabilities of ref. <sup>20</sup> (Figure S4D). In said work, it was found that connection strength was independent of connection probability. Additionally, excitatory connections most strongly distinguished by postsynaptic connection, while inhibitory by presynaptic connection.<sup>20</sup>

The adjusted connection probabilities of ref.  $^{20}$  are reported as fits that are dependent upon the distance between cells in addition to the dependence on pre- and postsynaptic neuron type. Since we do not explicitly simulate the spatial distribution of cells in our model, we use length scales from the experimental measurements to set distance-dependent connection probabilities. Specifically, since the imaging field of the two-photon experiment was  $400~\mu\text{m}\times400~\mu\text{m}^{24}$ , we randomly generated cell locations within a two-dimensional box of this size and then computed the average connection probability between all possible pairs. The connection probability decay lengths were taken to be  $100~\mu\text{m}$  for E  $\rightarrow$  I or I  $\rightarrow$  E and  $125~\mu\text{m}$  for E  $\rightarrow$  E or I  $\rightarrow$  I.  $^{20}$  From the randomly generated cell locations, this resulted in a reduction of  $p_{\text{max}}$ , i.e., the connection probability if the cells were right on top of one another,  $^{20}$  by 0.25 for E  $\rightarrow$  I or I  $\rightarrow$  E and 0.34 for E  $\rightarrow$  E or I  $\rightarrow$  I. Taking into account this distance-dependent reduction yields the connection probabilities shown in Figure S4D.

- (1) **Relative cell counts:** We assume the microcircuit has a ratio of cell counts of Excitatory: VIP: SST found in the investigation of L2/3 of the visual cortex of mice from ref. <sup>18</sup> (reproduced in Table 1, Figure S4E). However, in order to maintain a reasonable cell counts for numerical simulation, we instead use the ratio  $n^E: n^S: n^V=2:1:1$ , and adjust each population's outgoing synapses to account for any discrepancy in their simulated cell count relative to their experimental cell count. For example, since in simulation there are only twice as many excitatory to SST cells, but from experimental data their ratio is closer to 27.35: 1, we strengthen each excitatory synapse by a factor of 27.35/2 = 13.675 to account for the missing simulated cells. From Figure 4G, we see this scaling of synaptic strengths to account for the differences in cell counts in our microcircuit maintains the relative population strengths.
- (2) **Synapse strength:** We take the time-integrated voltage over a typical postsynpatic potential (PSP) pulse fit as a measure of the synaptic strength, where

$$Z^{PSP} = PSP'_{unmod} \times T_{eff}$$
 (Equation 17)

where  $PSP'_{unmod}$  is the adjusted PSP amplitude when the neuron is not being facilitated or depressed from STSP effects<sup>20</sup> and  $T_{eff}$  is the effective time of the PSP pulse.

We compute an adjusted PSP amplitude that accounts for potential differences of the *in vitro* measurements versus what we assume to be a cell's *in vivo* operating potential.<sup>20</sup> These differences are distinct across cell types, and thus can affect the relative strengths of excitation and inhibition within the cortical circuit. To arrive at the adjustment factor, we assume the experimental current is proportional to the difference in the experimental reversal potential and the resting potential. Furthermore, we assume the *in vivo* current is proportional to the difference in reversal potential and the potential where we presume neurons are generally close to





operating, which we take to be the threshold potential. The constant of proportionality in both cases is taken to be the neuron's conductance, and thus the ratio of these differences gives the adjustment factor of the PSP value, namely

$$PSP'_{unmod} = \left(\frac{V_{rev} - V_{thresh}}{V_{rev}^{(exp)} - V_{rest}^{(exp)}}\right) \times PSP_{unmod}$$
 (Equation 18)

where  $V_{\text{rev}}$  is the estimated reversal potential of relevant channels in the presynaptic neuron,  $V_{\text{thresh}}$  is the estimated threshold potential of the postsynpatic neuron,  $V_{\text{rev}}^{(\text{exp})}$  is the experimentally measured reversal potential that is dependent on neurotransmitters of the presynaptic neuron, and  $V_{\text{rest}}^{(\text{exp})}$  is the experimentally targeted resting potential that is presynaptic dependent.

From the literature, we use  $V_{\text{rev}} = 0 \text{ mV}$  for excitatory <sup>66</sup> and  $V_{\text{rev}} = -82 \text{ mV}$  for inhibitory presynaptic populations. <sup>67,68</sup> We determine  $V_{\text{thresh}}$  from electrophysiology data from the Allen Cell Types Database, found at <a href="https://celltypes.brain-map.org/data.">https://celltypes.brain-map.org/data.</a> Specifically, the  $V_{\text{thresh}}$  for each sweep and averaging over all sweeps for a given specimen identification, then averaging these values across the Cre-line. We only vary  $V_{\text{thresh}}$  by postsynaptic cell identity. The Cre-lines, total cell count, and computed  $V_{\text{thresh}}$  are shown in Table 2.

Since only a small subset of synapses have experimentally measured  $V_{\rm rev}^{(\rm exp)}$ , we take the median value across synapses of a given pre- and postsynaptic neuron type and use this across all cells. We do not find this impacts the resulting  $V_{\rm rest}'$  significantly. Finally, for  $V_{\rm rest}^{(\rm exp)}$ , we use the targeted holding potential from experiment, which are -70 mV for excitatory presynaptic cells and -55 mV for inhibitory presynaptic cells. Junction potential corrections of -14 mV are accounted for at all steps of this calculation. Altogether, these above computation yields the PSP $_{\rm unmod}'$ PSP $_{\rm unmod}$  ratios shown in Figure S4H.

The effective time of the PSP pulse is computed by integrating the PSP fits over time.<sup>20</sup> Up to an amplitude correction, synapse PSPs were fit using the following function

$$F_{\text{PSP}}(t) = \begin{cases} \frac{1}{A_{\text{norm}}} (1 - e^{-t/\tau_{\text{rise}}})^2 e^{-t/\tau_{\text{fall}}} & t \ge 0, \\ 0 & t < 0, \end{cases}$$
 (Equation 19)

where  $A_{\text{norm}} = F_{\text{PSP}}(T_{\text{max}})$  with  $T_{\text{max}} = \tau_{\text{rise}} \ln(1 + \tau_{\text{fall}} / \tau_{\text{rise}})$  is a normalization factor to ensure the maximum of  $F_{\text{PSP}}$  is equal to 1.0.<sup>20</sup> Integrating this expression over time, we find the effective time of the PSP fit,

$$T_{\text{eff.}} = \int F_{\text{PSP}}(t) dt = \frac{1}{A_{\text{norm}}} \frac{2\tau_{\text{fall}}^3}{(\tau_{\text{rise}} + \tau_{\text{fall}})(\tau_{\text{rise}} + 2\tau_{\text{fall}})}.$$
 (Equation 20)

This procedure yields the values shown in Figure S4I.

We computed  $Z^{PSP}$  for each synapse and then averaged across all synapses of the given pre- and postsynaptic cell type (see Figure S4F for count). For certain pre- and postsynaptic populations, we found the fits of  $\tau_{\text{fall}}$  were exceedingly high, and so any synapse with a  $\tau_{\text{fall}}$  > 300 ms was omitted. In practice, this only resulted in a small decrease in the number of synapses for each pre- and post-synaptic cell combination (Figure S4F).

#### **Response sparsity adjustment**

Similar to the FMSN, the biases/threshold of the various populations are adjusted in order to set baseline response sparsity at initialization. Since we consider 3 primary populations of neurons in this work, this procedure amounts to the fitting the 3 parameters of the network at initialization, namely  $b^E$ ,  $b^S$ , and  $b^V$  that determine the biases in Equation (12). Again like the FMSN, since our model neglects many influences that affect the firing rates of the various populations, e.g., from inputs from other layers or from PV neurons, we assume that this bias adjustment partially accounts for the mean activity of other possible inputs. In particular, since some populations receive fairly unbalanced inputs from excitatory or inhibitory populations (i.e., the VIP population), this threshold adjustment is assumed to at least partially account for excitatory-inhibitory balance. Unlike the FMSN, our microcircuit model has recurrent connections, and so any adjustment to response sparsity at one time step affects the inputs and thus the response sparsity of subsequent time steps. This leads us to a different bias fitting procedure to account for this additional complication. Finally, note this entire procedure is performed at the network's initialization, prior to any unsupervised training, and is thus insensitive to FMS mechanism placement or parameters.

The neuron population thresholds ( $b^S$ ,  $b^S$ , and  $b^V$ ) are adjusted using supervised training to reach a certain population response sparsity over a validation set prior to training. The validation set consists of the 8 familiar input vectors as well as 504 additional vectors (for a total of 512) drawn from the same distribution. In this work, all neurons of the same population share the same firing threshold parameter, meaning particular neurons within a given population may fire more/less over the given validation set.

In particular, the response rate is adjusted with respect to the loss function

$$L_{\text{MSE}}(\mathbf{g}, \widehat{\mathbf{g}}) = \sum_{p = \text{E,V,S}} \left( \frac{g^p - \widehat{g}^p}{g^p + \epsilon} \right)^2,$$
 (Equation 21)



where  $\mathbf{g} = (g^E, g^V, g^S)$  are the experimental response rates to be matched,  $\hat{\mathbf{g}} = (\hat{g}^E, \hat{g}^V, \hat{g}^S)$  are the network's approximate response rates, and  $\epsilon = 10^{-4}$  provides numerical stability. Note here we use a weighted mean squared error loss so that populations with smaller response rates are treated on even footing. See Figure S5M for exemplar fit results. For a given network population p, the approximate response sparsity is computed as

$$\widehat{g}^{\rho} = \frac{1}{n^{\rho}} \sum_{i=1}^{n_{\rho}} \sigma(\gamma \widetilde{y}_{i}^{\rho}),$$
 (Equation 22)

where  $\sigma(\cdot)$  is the sigmoid function and serves as smoothed version of the step-function to enable backpropagation,  $\gamma=100$  controls the rate of smoothing, and  $\tilde{y}_i^{\rho}$  is the preactiviation of neuron i, see STAR Methods.

As mentioned above, since the microcircuit we investigate is recurrently connected, adjusting the response rate of one population influences the response rate of the other populations, so a self-consistent solution across all neurons must be met. To do so, the validation set is repeatedly passed to the network and the thresholds of all neuron populations are adjusted simultaneously until a self-consistent solution is found. To best resemble the stimulus sequence the network will be exposed to during the stimulus change detection task, the validation sequence is smoothed with a ramping function that matches the deconvolved signal. Response rates are computed for the populations at the peak of the ramping function, and the thresholds of the various populations are adjusted using backpropagation through time. Backpropagation is truncated to 10 time steps backwards. We used ADAM with default parameters, a batch size of 128, and shuffled the validation set every 10 network steps. Note we neglect the effect of the time-correlated noise that is injected into all neuron populations in adjusting the firing rates. Additionally, we assume the stimulus history information during images changes will be strongly suppressed, so within the validation set the corresponding part of the input is just noise.

We use  $\mathbf{g} = (0.05, 0.4, 0.2)$  throughout this work to represent target firing rates during an image response in the novel session. We do not observe a large difference of the parameters  $b^{\mathrm{E}}$ ,  $b^{\mathrm{V}}$ , and  $b^{\mathrm{S}}$  values across different microcircuit initializations.

#### Response rates and variance explained

We note that we do not aim to exactly reproduce response rates or variance explained across cells that are observed in experiment due to the significantly smaller neuron populations used throughout this work. As mentioned above, for the parameters used in our cortical circuit model, individual cells receive synaptic input from on order 10 to 100 other cells. Simulating larger populations of neurons would allow significantly more noise to be injected into each individual cell since each cell, on average, has a smaller effect on output behavior of other neurons.

#### Noise injection and matching baseline responses

From the experimental data, we see that all neuron populations exhibit a nonzero baseline mean response between image stimuli (Figure S6). Said baseline responses carry almost no information about image identity or task information except within a small window after the image stimulus turns off, indicating they may represent neuronal activity unrelated to the image change detection task.<sup>24</sup> To model the effects of these baseline responses, we inject time-correlated noise directly into each neuron population. We adjust the variance of this time-correlated noise to match the baseline response values observed in each population using supervised training. Similar to the firing rate adjustment considered above, this again corresponds to only one number per population, so this procedure fits a total of 3 parameters at the network's initialization (and occurs after the firing rate adjustment). Once again, since this procedure occurs at network initialization, it is completely independent of FMS placement or parameterization.

We define the experimental baseline responses to be the mean population response halfway between the pre-change image and the change image. With this definition, for a given population, the mean population response doesn't change much between the familiar and novel sessions, so we average across the two sessions. Taking the novel session values yields baseline targets of  $1.5 \times 10^{-3}$ ,  $6.1 \times 10^{-3}$ , and  $3.4 \times 10^{-3}$  for the excitatory, VIP, and SST populations, respectively. Once again, we weight how well our model fits the experimental data using a weighted MSE loss,

$$L_{\mathsf{MSE}}(\boldsymbol{\delta}, \widehat{\boldsymbol{\delta}}) = \sum_{\rho = \mathsf{E.V.S}} \left( \frac{\delta^{\rho} - \widehat{\delta}^{\rho}}{\delta^{\rho} + \epsilon} \right)^{2}, \tag{Equation 23}$$

where now  $\delta = (\delta^E, \delta^V, \delta^S)$  are the experimental baseline mean responses to be matched,  $\hat{\delta} = (\hat{\delta}^E, \hat{\delta}^V, \hat{\delta}^S)$  are the network's baseline mean responses. The network's baseline responses are simply the mean over each population response,

$$\delta^{p} = \frac{1}{n^{p}} \sum_{i=1}^{n_{p}} y_{i}^{p}.$$
 (Equation 24)

For simplicity, to initially fit the amount of noise injected into each population, we inject uncorrelated noise with standard deviation  $\sigma^p$ . In practice, we find the addition of time-correlation causes a negligible change in the networks' baseline responses (and the kernel used to generate uncorrelated noise is chosen such that it has approximately the same variance as its uncorrelated counterpart). See Figure S5N for exemplar fit results.





Similar to the above firing rate adjustment, we pass uncorrelated to the network as the input stimulus repeatedly until it reaches a self-consistent solution. We once again use ADAM with default parameters and truncate backpropagation through time to 3 network passes backward. Both the present stimulus and stimulus history information parts of the input are assumed to be just noise for the validation set. Due to additional contributions from the stimulus history input during gray screen times, the VIP baseline was found to overestimate the noise needed to fit the data, so the noise injection was reduced after fitting by a fixed percentage across all initializations.

#### **Response smoothing**

To get the raw cell responses, we convolve the output function with the same half-normal function to match the responses of ref. <sup>24</sup>. Explicitly,

$$g(t) = \begin{cases} \sqrt{\frac{2}{\pi\sigma^2}} & \exp\left(-\frac{t^2}{2\sigma^2}\right) & t \ge 0, \\ 0 & t < 0. \end{cases}$$
 (Equation 25)

with  $\sigma = 60$  ms (Figure S5A). We discretize the kernel long time steps separated by  $\Delta t$  and normalize such that the summed amplitude is equal to 1.0.

#### **Familiarity-novelty task**

The familiarity-novelty task is used to train the FMSN discussed in the main text. The neuronal encoding of different stimuli are represented by distinct random binary vectors,  $\mathbf{x}^{\alpha} \in \mathbb{R}^{d}$ , where  $\alpha$  indexes the distinct stimuli. The random binary vectors are chosen to be sparse, i.e., the elements of stimuli  $\alpha$  are given by

$$\chi_L^{\alpha} = Ab_L^{\alpha}, b_L^{\alpha} \sim \text{Bernoulli}(p_{\text{stim}}),$$
 (Equation 26)

where A is some normalization factor. Since we generally use small  $p_{\text{stim}}$ , we also ensure that there are a minimum number of nonzero elements for each  $\mathbf{x}^{\alpha}$ .

Prior to training,  $N_F$  stimuli are generated and defined to be the *familiar set*,  $S_F = \{\mathbf{x}^1,...,\mathbf{x}^{N_F}\}$ . The unsupervised training consists of passing the network a sequence whose elements are randomly drawn from the set  $S_F$  (with replacement). After each sequence step, the network's modulations are updated according to Equation (2). Random Gaussian noise (iid to each element/sequence step)(Results do not differ significantly from bit-flipped noise, both methods increase the dot product between two randomly drawn stimuli, making the familiar stimuli harder to distinguish from the novel stimuli.),  $\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma_{\epsilon} \mathbf{1})$ , is added to the inputs before each element is passed through a ReLU function to ensure all elements are positive,

$$\mathbf{x}_t = \mathsf{ReLU}(\mathbf{x}^\alpha + \boldsymbol{\epsilon}).$$
 (Equation 27)

As mentioned in the main text, the sequence of familiar stimuli is ordered such that each element of the familiar set is seen every  $N_F$  sequence steps. The order of the familiar set is shuffled within every  $N_F$  window.

Implicit in this training is that the time difference between successive stimuli is constant, a feature we relax in stimulus change detection task. The parameter  $\lambda$  that controls the decay can be thought of as corresponding to a decay length relative to number of examples.

During and after training, we test the network's response to both the familiar set of stimuli as well as a novel set of stimuli,  $S_N = \{\mathbf{x}'^1, ..., \mathbf{x}'^{N_F}\}$ , where  $\mathbf{x}^{\alpha} \neq \mathbf{x}'^{\beta}$  for all  $\alpha$  and  $\beta$ . For these test responses, the network's modulations are not updated after being passed through the network, so they do not affect the network's response to future inputs. Measuring the network's response at these steps is simply done for the sake of comparison and is not a necessary step in training.

#### **Experimental image change detection task**

The image change detection stimulus sequence consists of image presentations in quick, regular succession. Stimuli are presented for 250 ms and then followed by a gray screen for 500 ms (see ref. <sup>24</sup> for significantly more details on this task). The same image is presented several times in a row before a new image is chosen and the process repeats (Figure 4B). In the experiment, mice are tasked with licking in response to an image change. The number of times an image is presented in a row is between 4 and 11, with the count being drawn from a truncated exponential distribution so that 4 image presentations in a row is the most likely. Additionally, there is a 3 s grace period after an image presentation before the trial restarts. Note the image after an image change is drawn from the entire set of possible images and as such there is a 1/8 chance that the *same* image is drawn again. These cases are not included when measuring the network response to image changes.

The mice are first exposed to static grating and trained on a grating change detection task, which was found to improve the training time on the subsequent image change detection task. Mice are trained on the image change detection task from a set of 8 images that gradually become the *familiar set* (Figure 4C). During these *training sessions* mice learn to perform the task. Once mice reach a particular performance threshold on the image change detection task using the familiar set, their neuronal responses are recorded over several *familiar imaging sessions* that are separated by at least one night of rest. Generally, the second of these sessions is a



"passive" imaging session where they do not need to perform the task to obtain rewards (see ref. <sup>24</sup> for exact training sequences). Afterward, a novel set of 8 images are introduced into the same image change task. Without any additional training, the mices' neuronal responses while performing the task are recorded over several *novel imaging sessions*. Once again, generally the second of these is a "passive" session and is omitted from this analysis. After at least one session of exposure to the novel imaging set, the mice's responses while performing the same task on the novel imaging set are gathered in what is called a *novel-plus imaging session*.

During only the imaging sessions, i.e., not included in the training sessions, each image has a 5% chance of being omitted. For an omission, the gray screen continues to be displayed for the 250 ms where the image would have been presented (Figure 4A, middle). There is no limit on the number of omissions that can occur in a row, though longer chains become increasingly rare. Omissions cannot occur for the image presentation that would be a change or the pre-change image. This means that all omissions, including sequences of multiple omissions, are surrounded on either side by the same image.

#### Model stimulus change task

The task we train our cortical circuit model upon is meant to imitate the image change task used in the experimental data we are modeling. <sup>23,24</sup> At any given time, the input to the network consists of three distinct parts.

- (1) **Present stimulus:** A stimulus input vector,  $\mathbf{x}_t^{\text{stim}}$ , representing a neuronal encoding of the current visual input at time t (see Figures S5A and S5 for examples).
- (2) **Stimulus history information:** Information about the recent history of the stimulus sequence, specifically an encoding of the amount of time that has elapsed since the last stimulus presentation,  $\mathbf{x}_{t}^{hist}$  (see Figure S5C for example).
- (3) **Time-correlated noise:** Additional noise input into each population representing contributions to the neuronal activity from factors neglected in our model (e.g., behavioral effects),  $\mathbf{n}_p^p$  for p = E, S, V (Figure S5D).

All input sequences are discretized to a time length of  $\Delta t = 1/32$  s, or 31.25 ms. This time difference is chosen to match the experimental sampling rate.

#### **Present stimulus**

The present stimulus sequence is constructed to represent a neuronal encoding of the equivalent visual stimulus of the experiment. It is constructed to closely match the statistics of the image change detection task the mice are trained upon. Since image presentations last 250 ms and are separated by 500 ms of gray screen (ignoring the possibility of omissions for the moment), there are 250 ms  $\Delta t = 8$  time steps of the neuron encoding of the image followed by  $500 \text{ ms}/\Delta t = 16$  time steps of the neuronal encoding of gray screen (though see below for additional details). This sequence then repeats, with the image identity of each 750 ms window being chosen so that image changes and omissions occur at frequencies described above. Different images encodings are represented by distinct random binary vectors drawn in an identical manner to that described in familiarity-novelty task above. Similar to experiment, the familiar and novel sets are chosen to have 8 distinct stimuli in them. All inputs have random Gaussian noise added to them. As observed in the experimental data, neuronal activity is low during stimulus times where the gray screen is displayed, so the present stimulus inputs representing gray screen encodings only consist of the added Gaussian noise discussed above. When an image omissions occurs, the image input is simply replaced by additional gray screen input. We allow for at most two omissions to occur in a row.

During the time steps representing an image presentation, there are three additional contributions to the stimulus sequence used to mimic the responses of experimental studies. First, to best match mice response data, we delay the mean onset of the image presentation stimulus by two time steps, corresponding to  $2\Delta t \approx 62$  ms, relative to when we consider the image stimulus has begun display. This has the net effect of shifting the neuronal responses later in time relative to the image presentation time period (for example, see Figure 5) and approximately matches known delays of the visual cortex to visual stimuli<sup>69</sup> as well as the experimental data. <sup>23,24,69</sup> Second, we also smooth the input stimulus with a smoothing kernel to represent the ramping and decay of the image response to the input sequence. <sup>69</sup> The smoothing kernel is the normalized experimental mean excitatory response, deconvolved with the experimental stimulus smoothing function (see above for details). (In practice, smoothing the present stimulus signal from the L4 excitatory response would have been more realistic. However, the depth differences between L2/3 and L4 did not change the excitatory response significantly, so we have just used L2/3 for simplicity.) The resulting smoothing kernel from this process is shown in Figure S5F. Third, for each cell we allow the onset of the image stimulus to vary by  $\pm \Delta t$  so that not all cells receive input representing the image presentation at the same time step. This incorporates effects of lag times of stimulus responses across a population and is also useful for numerical stability so that all cells do not respond in unison. We incorporate all three of the above effects on a cell-by-cell basis into a stimulus kernel  $\mathbf{k}_t^{\rm stim}$ . If  $\mathbf{x}_t$  is the random binary vector representing the current image being presented, then the full raw stimulus stimulus input is given by

$$\mathbf{x}_{t}^{\text{stim}} = \text{ReLU}(\mathbf{k}_{t}^{\text{stim}} \odot \mathbf{x}_{t} + \epsilon),$$
 (Equation 28)





where  $ReLU(\cdot) = max(\cdot, 0)$  and  $\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma_{\epsilon}\mathbf{1})$  Z (HTML translation failed). See Figures S5A and S5B for exemplar present stimulus input for an image change and image omission.

Unless otherwise stated, the present stimulus inputs are taken to have average sparsity p = 0.05 with a minimum sparsity of 0.025 for any given stimulus. The input normalization was chosen so that nonzero elements had size 0.2. Since time-correlated noise is already injected into each population (see below), the present stimulus noise was taken to only be  $\sigma_{\epsilon} = 0.03 \times 0.02$ .

Since mice can lick mid-sequence and this can reset the task mid-trial, the experimental distribution of number of presentations between a given image is thus fairly heavy tailed. Across the entire experimental dataset, we find there to be on average 20.4 image presentations between changes. To ensure shorter trial times, we truncate the maximum number of image presentations between a change to 75. This only omits the 2.3% of image changes on the tail, and shifts the average number of images between a change to 18.4. See Figures S5G and S5H for a plot of the true number of image presentations until the next presentation. We do not find using the experimental image-change-distribution versus the idealized one that assumes no licks affects the mean response results significantly. However, in the cell subpopulation analysis the fitting metrics are dependent upon the global distribution of change occurrences and so the truncated experimental distribution was used in said analysis.

#### **Stimulus history information**

As mentioned above, we also pass the network information about the recent history of the stimulus, in particular the time that has elapsed since the last image presentation. This information is assumed to be encoded in a subset of the top-down input to the cortical circuit. This additional input into the network is necessary to observe responses that are dependent on the relative time between image presentations, e.g., the omission ramps. The top-down inputs could be produced from the present stimulus sequence described above using, say, a simple recurrent network that counts the time steps since the last stimulus and encodes said information in output neuron responses that match known stimulus tuning properties. As our goal for this study is the effects for FMS in the local circuit, we avoid an explicit implementation of such history encoding and simply input it directly into the network.

In this section, we denote the time that has elapsed since the last image presentation at time t by s, which is measured in seconds. For example, with no omissions, s=500 ms immediately before the onset of the next image presentation. For times when the stimulus is currently being presented, s=0. The time since the last image presentation is maximized after omissions, and since at most we allow for two successive omissions,  $0 \le s \le 2$  seconds.

We denote the neuronal encoding of the time s by  $\mathbf{r}(s)$ . We assume that encodings of times that are close together are more similar than times further apart, as measured the dot product between the two representations. That is, if |s - s'| < |s - s''|, then  $\mathbf{r}(s) \cdot \mathbf{r}(s') > \mathbf{r}(s) \cdot \mathbf{r}(s'')$ .

The temporal encoding input is generated by creating a population of neurons that are each tuned to a particular s. For simplicity, we take the neuronal tuning curves to take the shape of a Gaussian, though our results easily generalize to other tuning curves such as cosine bumps. To match experimental results, we assume the population of neurons' tuning curves are centered at times that are logarithmically distributed and that the width of each tuning curves is proportional to their center. To Epecifically, we take the neuron tuning centers to be logarithmically distributed between  $10^{-2}$  to  $10^{1}$  seconds. For neuron *I* centered at time  $s_{I}$ , its width is directly proportional to the size of the center of its tuning,  $\sigma_{I} = \frac{1}{3}s_{I}$ . Altogether, the tuning curve of neuron *I* is

$$r_l(\mathbf{s}) = \exp\left(-\frac{(\mathbf{s} - \mathbf{s}_l)^2}{2\sigma_l^2}\right).$$
 (Equation 29)

See Figure S5I for exmplar raw turning curves. In practice, the longest delay time between images is at most 2 s, so those neurons tuned to higher times are almost always silent in our setup. The resulting neural population responses are then each individual neuron's response to the corresponding s (Figure S5J). Due to the higher density of neurons centered to small s, the relative magnitude of the population response vectors decreases gradually for larger times (Figure S5K). Notably the trend in magnitude is the opposite to that of the ramping response, i.e., smaller s have the largest magnitude. A verification of the decrease in similarity for times further apart is shown in Figure S5L. Since we assume this history stimulus represents some unknown subset of the total top-down input into the VIP population, we simply set the magnitude of the cell activity to be comparable to the present stimulus input so that the omission ramping has a similar response to image presentations.

Similar to the present stimulus input, noise is added to the stimulus history stimuli and thresholded to be positive definite,

$$\mathbf{x}_t^{\text{hist}} = \text{ReLU}(\mathbf{r}(\mathbf{s}) + \epsilon).$$
 (Equation 30)

We note that this encoding of the history of what the mouse viewed is purposely simplistic and likely misses other effects that could be observed experimentally. For instance, an image that lasts longer or a shorter delay would not elicit a large response from the VIP cells despite these being outside of the the normal rhythm of the task. A more thorough encoding of the history of the task would allow the model to react to additional disruptions to the regular task flow, but we leave such exploration for future work.



#### **Time-correlated noise**

To account for contributions from neurons not included in the circuit, as well as contributions from other task-relevant effects (e.g., behavior), the excitatory, SST, and VIP populations are injected with additional time-correlated noise (in addition to the noise added to the present stimulus and stimulus history inputs described above).

Specifically, the noise is generated by convolving white noise with a Gaussian smoothing kernel over time, and then weighting the noise injected into each neuron to account for any variance the population may have from such effects. The continuous Gaussian smoothing kernel is given by  $k(t) = 1/(\pi\sigma^2)^{1/4} \exp(-t^2/2\sigma^2)$  with  $\sigma = 125$  ms. The discrete smoothing kernel,  $k_t$ , is found by evaluating the above at  $k_t = k(t\Delta t)$  for  $t = -(\sigma/\Delta t)^2, ..., (\sigma/\Delta t)^2$ . Note the normalization of k(t) is chosen such that the convolution does not change the variance of the uncorrelated noise, e.g.,  $\sum_{k} k_t^2 \approx 1$ . The weight accounting for how much noise is injected into neuron i is

$$w_i \sim \mathcal{U}(0,1)$$
. Thus we have

$$\mathbf{n}_t^{\rho} = \mathbf{w} \odot (k * \tilde{\mathbf{n}}^{\rho})_t,$$
 (Equation 31)

where  $\tilde{\mathbf{n}}_t^{\rho}$  is the uncorrelated noise with  $\tilde{\mathbf{n}}_t^{\rho} \sim \mathcal{N}(\mathbf{0}, \sigma^{\rho})$ . As shown above, this noise is added to the preactivations of all neurons for each population. The variance of the noise for each population,  $\sigma^{\rho} = \sigma^{\rho} \mathbf{1}$ , is adjusted to match experimental baselines, see STAR Methods for details.

#### **Training schedule details**

As discussed in the main text, the mices' training schedule consists of many sessions that each last on the order of one to two hours.<sup>24</sup> Several sessions are required for the mice to learn the task completely, meaning often they have been exposed to on the order of 10 hours of the image change task to achieve the task performance threshold needed to progress onto imaging. Since neuronal responses are only collected after this performance threshold is achieved, it is not yet known how many sessions are required for the neuronal responses to the familiar image change task to stabilize.

For numerical tractability, we do not explicitly simulate the full tens of hours of the training sequences for the microcircuit model. Instead, we expose the model cortical circuit to a shorter version of the task and increase its learning rates so that it achieves stabilized responses to the familiar data over a shorter simulated time. As we saw in the FMSN, higher learning rates are capable of becoming familiarized with responses at a quicker rate, at the cost of fitting the noise to a greater degree. Thus, explicitly simulating full training/imaging times at equivalent lower learning rates should only improve the results we have shown throughout this work. Additionally, as we mentioned above, we did not find that using that exact distribution of image change times affected any results outside of the cell subpopulation analysis, and thus to further expedite training we reduced the number of repeated presentations that are between each image change to between 4 and 9.

An explicit demonstration of the training time equivalence is explored in the FMSN in Figures S2O and S2P. There, the FMSN is trained on sequences that vary in length over two orders of magnitude and it is shown that, by correctly adjusting  $\lambda$  and  $\eta$ , the networks essentially develop almost identical responses despite the large differences in training time. For two training sequence lengths T and T', the equivalence is given by

$$T \to T', \eta \to \frac{T}{T'} \eta' \quad \lambda \to (\lambda')^{T/T'}.$$
 (Equation 32)

Thus, if T' > T, this corresponds to a reduction in the learning rate and decay rate for the longer train sequence. Note the last relation is equivalent to  $\tau_{\text{decay}} \rightarrow \frac{T'}{T} \tau'_{\text{decay}}$ .

Specifically, we train our cortical circuit models on 2000 s of change detection, which consists of approximately 2660 total image presentations or 330 presentations of each familiar image. With the expedited image change time, the training session consists of roughly 400 image change events.

To gather responses in cortical circuit model's 'imaging' sessions, we monitor the network's responses while it continues to train. We measure the network's responses on 250 s of change detection, but we take advantage of batch training to allow us to gather responses to several distinct input streams simultaneously. However, since the cell coding analysis is dependent on statistics over an entire session, we use the actual image change distribution when collecting data for said analysis.

Since we would like to simulate the gradual familiarization of the novel set during the novel imaging session, but we have used a larger learning rate to expedite the training procedure, we reduce the learning rate of both  $FMS_A$  and  $FMS_{AH}$  during imaging sessions. There is evidence the mices' response changes between sessions even without explicit exposure to the stimuli.<sup>24</sup> This may be due to replay. To simulate this additional familiarization that occurs between sessions, as well as additional stimulus exposure during the passive session, we train the networks on the novel images for an additional session equal in length to the imaging sessions, but at a higher learning rate, similar to training.





#### Analytical intuition and additive modulations Additive modulations

An alternative form of modulations that has also recently been considered has the modulations directly added to the fixed weights rather than the multiplicative form we consider throughout this work. Explicitly, rather than the transformation of the form of Equation 1, these modulations affect the fixed weights via

$$\mathbf{W} \to \mathbf{W} + \mathbf{M}_t$$
. (Equation 33)

The same modulation update expressions continue to be used: for the associative case, Equation 2a, and pre-only, Equation 2b. Equivalent modulation bounds to those in Equation (7) are enforced such that the associative and pre-only modulations do not exceed biological bounds of LTP/D and STSP, respectively. Many equivalent FMSN results to those in the main text and Methods using this additive modulation are shown in Figure S11.

#### **Analytical intuition**

Here we discuss some analytic properties of FMSs that help us understand them better. Throughout this section, we will investigate both the multiplicative modulations considered in the main text as well as the additive modulations that were introduced n above. The analytical approximations for the additive modulations and their connection to Hopfield networks is more straightforward to understand, but similar qualitative results hold for the multiplicative modulations. Like the analysis here, feedforward version of Hopfield networks as familiarity discriminators were explored in ref. <sup>32</sup>, though their setup requires specialized weights that take on the value of the input neurons to compute the energy function used for familiarity discrimination.

Like the main text, let  $\mathbf{M}$  be the modulation matrix and  $\mathbf{W}$  be the fixed matrix of synaptic connecting the input and the output. Let  $\mathbf{x}_1,...,\mathbf{x}_m$  be the set of familiar inputs we would like the network to memorize. Additionally, let  $\tilde{\mathbf{x}}_{\alpha}$  for  $\alpha=1,...,m$  be novel inputs, that also obey  $\tilde{\mathbf{x}}_{\alpha} \cdot \tilde{\mathbf{x}}_{\beta} = \delta_{\alpha\beta}$  and  $\mathbb{E}\tilde{\mathbf{x}}_{\alpha} = 1$  for all  $\alpha,\beta=1,...,m$ , but also  $\mathbf{x}_{\alpha} \cdot \tilde{\mathbf{x}}_{\beta} = 0$  for all  $\alpha$  and  $\beta$ . Here we assume the novel and familiar sets are of the same size, but it is straightforward to generalize what we show here for different size sets. Note since we consider inputs that are positive definite, in practice the dot product between any two inputs is finite, but it can approach zero as the size of the input space gets large and the sparsity is small.

To begin with, we establish some properties of the element-wise product that will be useful for the multiplicative form of the modulations. We will often make use of the identity that shows how an outer product of vectors (e.g., the modulations) act through matrix multiplication.

$$[(\mathbf{y}\mathbf{x}^T)\odot\mathbf{W}]\mathbf{x}'=\mathbf{y}\odot[\mathbf{W}(\mathbf{x}\odot\mathbf{x}')], \qquad (Equation 34)$$

from ref. <sup>36</sup>. Additionally, it will be useful to compare the (L2) magnitudes of the element-wise product between two stimulus vectors. Let nonzero values of the vectors be A with probability p. If the two vectors are different we have

$$\mathbb{E}\|\mathbf{x}\odot\mathbf{x}'\|_{2}^{2} = \sum_{l=1}^{d} \mathbb{E}(x_{l}x'_{l})^{2} = \sum_{l=1}^{d} p^{2}(A)^{4} = dp^{2}A^{4},$$
 (Equation 35)

where we have used the fact that the only nonzero element of the element-wise product occurs when the Meanwhile, if the two stimulus vectors are the same, instead we have

$$\mathbb{E}\|\mathbf{x}\odot\mathbf{x}\|_{2}^{2} = \sum_{I=1}^{d} \mathbb{E}(x_{I}x_{I})^{2} = \sum_{I=1}^{d} p(A)^{4} = dpA^{4},$$
 (Equation 36)

which is larger by a factor of p. Note a similar property holds for the dot product between two vectors, where  $\mathbb{E}\|\mathbf{x}^T\mathbf{x}'\|_2^2 = d^2p^4A^4$  and  $\mathbb{E}\|\mathbf{x}^T\mathbf{x}'\|_2^2 = d^2p^2A^4$ , but now the same vector result is larger by a factor of  $p^2$ . Thus, in the limit that  $p \to 0$  and proper normalization of A, we make the analogous approximations

$$\mathbf{x}_{\alpha} \odot \mathbf{x}_{\beta} \approx \delta_{\alpha\beta} \mathbf{x}_{\alpha}^{2}, \mathbf{x}_{\alpha} \odot \tilde{\mathbf{x}}_{\beta} \approx \mathbf{0},$$
 (multiplicative) (Equation 37a)

$$\mathbf{x}_{\alpha} \cdot \mathbf{x}_{\beta} \approx \delta_{\alpha\beta}, \mathbf{x}_{\alpha} \odot \tilde{\mathbf{x}}_{\beta} \approx 0, \text{ (additive)}$$
 (Equation 37b)

where we use the shorthand  $\mathbf{x}_{\alpha}^2 = \mathbf{x}_{\alpha} \odot \mathbf{x}_{\alpha}$ .

Let us start with an approximate setting that will serve as a rough representation for the function of the FMSN. We assume that the modulations are *not* involved in the the feedforward pass of our network but are still updated as we have described in the main text. That is, the output activity is given by  $\mathbf{y} = \varphi(\mathbf{W}\mathbf{x} + \mathbf{b})$  but we still update  $\mathbf{M}_t$  via Equation 2a even though it has no effect on the network's behavior. We also presume the modulations do not decay, i.e.,  $\lambda = 1$ , and that modulation are small enough such that they do not encounter the biological bounds or violate the bounds of Equation (5).



Over a training time where the familiar inputs are each shown N times, the associative update will lead to an expected  $\mathbf{M}$  given by a sum of m outer products

$$\mathbf{M} = \eta N \sum_{\alpha=1}^{m} \mathbf{y}_{\alpha} \mathbf{x}_{\alpha}^{T}.$$
 (Equation 38)

Notably, if  $\mathbf{W} = \mathbf{I}$  and  $\mathbf{y}_{\alpha} = \mathbf{x}_{\alpha}$ , i.e., if  $\varphi(x) = x$ , this would be the same form of the updates to the lateral connections in a Hopfield network with an associative learning rule. Now consider how this modulation matrix acts on a given familiar input (for next three equations, top: multiplicative, bottom: additive),

$$(\mathbf{W} \odot \mathbf{M})\mathbf{x}_{\alpha} = \eta N \sum_{\beta=1}^{m} \mathbf{y}_{\beta} \odot \left[ \mathbf{W} \left( \mathbf{x}_{\beta} \odot \mathbf{x}_{\alpha} \right) \right] \approx \eta N \sum_{\beta=1}^{m} \delta_{\alpha\beta} \mathbf{y}_{\beta} \odot \mathbf{W} \mathbf{x}_{\beta}^{2} = \eta N \mathbf{y}_{\alpha} \odot \mathbf{W} \mathbf{x}_{\alpha}^{2},$$
 (Equation 39a)

$$\mathbf{M}\mathbf{x}_{\alpha} = \eta N \sum_{\beta=1}^{m} \mathbf{y}_{\beta} \mathbf{x}_{\beta}^{T} \mathbf{x}_{\alpha} \approx \eta N \sum_{\beta=1}^{m} \mathbf{y}_{\beta} \delta_{\alpha\beta} = \eta N \mathbf{y}_{\alpha},$$
 (Equation 39b)

where for the multiplicative case (top) we have used Equation 34 for the first equality and in both lines we have used the approximation of Equation 37. Now compare this to a novel input

$$(\mathbf{W} \odot \mathbf{M})\tilde{\mathbf{x}}_{\alpha} = \eta N \sum_{\beta=1}^{m} \mathbf{y}_{\beta} \odot \left[ \mathbf{W} \left( \mathbf{x}_{\beta} \odot \tilde{\mathbf{x}}_{\alpha} \right) \right] \approx \eta N \sum_{\beta=1}^{m} \delta_{\alpha\beta} \mathbf{y}_{\beta} \odot \mathbf{W} \mathbf{0} = \mathbf{0},$$
 (Equation 40a)

$$\mathbf{M}\tilde{\mathbf{x}}_{\alpha} = \eta N \sum_{\beta=1}^{m} \mathbf{y}_{\beta} \mathbf{x}_{\beta}^{T} \tilde{\mathbf{x}}_{\alpha} \approx \eta N \sum_{\beta=1}^{m} \mathbf{x}_{\beta}(0) = \mathbf{0}.$$
 (Equation 40b)

Thus a familiar input yields a non-zero modulation but a novel input simply yields zero. From the above results, we have

$$(\mathbf{W} + \mathbf{W} \odot \mathbf{M})\mathbf{x}_{\alpha} \approx \tilde{\mathbf{y}}_{\alpha} + \eta N \mathbf{y}_{\alpha} \odot \mathbf{W} \mathbf{x}_{\alpha}^{2}, (\mathbf{W} + \mathbf{W} \odot \mathbf{M})\tilde{\mathbf{x}}_{\alpha} \approx \tilde{\mathbf{y}}_{\alpha},$$
 (Equation 41a)

$$(\mathbf{W} + \mathbf{M})\mathbf{x}_{\alpha} \approx \tilde{\mathbf{y}}_{\alpha} + \eta N \mathbf{y}_{\alpha}, (\mathbf{W} + \mathbf{M})\tilde{\mathbf{x}}_{\alpha} \approx \tilde{\mathbf{y}}_{\alpha},$$
 (Equation 41b)

Thus we see that if  $\eta > 0$  (or  $\eta < 0$ ) the familiar preactivations grow (shrink) in size from the effects of the modulations, while the novel preactivations are left approximately unchanged.

Let the *familiar subspace* be the subspace of the input stimulus space spanned by the familiar stimuli. Then, since any stimulus can be decomposed into parts that lie within the familiar subspace and perpendicular to it, a generalization of the above arguments shows that the modulation matrix  $\mathbf{M}$  will yield a nonzero result for any vector that has components in the familiar subspace. Since for a large input stimulus space the novel inputs are close to perpendicular to the familiar subspace (so long as  $m \ll d$ ), they yield approximately zero output.

For the additive case, we can see Equation (41) is similar to checking the energy function of a Hopfield network (up to the vector  $\tilde{\mathbf{y}}_{\alpha}$ ), which has been used previously as a method of familiarity detection. <sup>32</sup> Indeed, since the activity in our network is positive definite, taking the L1 normalization of this output is equivalent to taking the dot product with **1**, so it is similar to a Hopfield energy measurement with one occurrence of the stimulus replaced by **1**.

Now in practice, the modulations are involved in the forward pass, so as the modulations get updated during training they affect the output. For the FMSN,  $\mathbf{y} = \varphi[(\mathbf{M} \odot \mathbf{W})\mathbf{x} + \mathbf{W}\mathbf{x} + \mathbf{b}]$  and thus the modulations affect its own update. Notably, it is only the *output* activity that is affected by our approximation above, and so what will change are the  $\mathbf{y}_{\alpha}$  dependence of Equations 38–41. However, what causes the significantly different behavior between Equations 39 and 40 is the *input* activity dependence of  $\mathbf{M}$ , and this is unchanged when we include modulations in the forward pass.

#### **QUANTIFICATION AND STATISTICAL ANALYSIS**

#### **Fitting FMS learning rates**

To determine the learning rates of the three FMS mechanisms we add to the microcircuit network, we perform a brute-force scan over grids of learning rates and evaluate how well the resulting modulation effects fit the novelty data. Note that procedure still amounts to fitting only three numbers:  $\eta^{(A)}$ ,  $\eta^{(P)}$ , and  $\eta^{(AH)}$ . We evaluate the parameters by seeing how well their mean image change and omission responses match the experimental data. Specifically, we compute the MSE loss of the mean fit to the experimental data over all three cell populations,

$$L_{\text{MSE}}(\mathbf{z}_{\text{c}}, \mathbf{z}_{\text{o}}; \widehat{\mathbf{z}}_{\text{c}}, \widehat{\mathbf{z}}_{\text{o}}) = \sum_{s = \text{F,Np}} \sum_{s = \text{E,V,s}} \sum_{t} \ell^{p} \left[ \left( z_{t,\text{c}}^{p,s} - \widehat{z}_{t,\text{c}}^{p,s} \right)^{2} + \delta_{s,\text{F}} \left( z_{t,\text{o}}^{p,s} - \widehat{z}_{t,\text{o}}^{p,s} \right)^{2} \right]$$
 (Equation 42)



where  $z_{t,c}^{\rho,S}$  is the experimental mean image change response of population  $\rho$  in session s at time t relative to the image change,  $\widehat{z}_{t,c}^{\rho,S}$  is the cortical circuit model equivalent,  $z_{t,o}^{\rho,S}$  is the omission equivalent,  $\ell^p$  is a population-dependent weight, and  $\delta_{s,F}$  is a Kronecker delta function ensuring omission loss is only computed during the familiar session (since we do not try to model the suppressed omission of novel sessions). We take  $\ell^V = 5$  and  $\ell^E = \ell^S = 1$ , so that fitting the VIP response is more important than the excitatory or SST populations. The sum over s represents a sum over the familiar and novel sessions. The sum over t represents a sum over the relative time to the change/omission for a given mean response. We take the relative time window to be 25 time steps before and after the corresponding image change/omission event, which is roughly  $\pm 800$  ms.

#### **Cell subpopulation analysis**

We reproduce the functional cell subtype analysis pipeline of ref. <sup>24</sup> to compare our model to experimental results on equal footing. Here we give a summary of said pipeline for completeness, additional details and justification for certain parameters we match to the experimental analysis can be found in ref. <sup>24</sup>. Throughout this section, we suppress indices that indicate the population and session of a given cell unless needed.

#### **Experimental data**

We take the computed coding scores directly from ref. <sup>24</sup>. The codings scores are for cells collected across several different brain areas and layers. Although our cortical circuit model specifically takes cell counts and connection data for L2/3, we note that the vast majority of VIP cells were found in upper cortical layers and there does not appear to be a significant difference in coding scores with brain area. <sup>24</sup>

The experimental coding score analysis focuses on four primary input feature categories (also called 'components'): images, omissions, behavioral, and task. These feature categories are further subdivided into various features that each have their own kernel and input data. For example, the image feature category contains one feature for each of the eight possible images in the corresponding image set. When a feature category is removed to compute its coding score (see details below), all feature kernels within that category are removed. See ref. <sup>24</sup> for additional details.

#### **Model fitting**

To understand how the various features coded in the task explain individual cell activity across the VIP, SST, and excitatory populations, we fit each cell's activity using a linear regression model with time-dependent kernels. The feature categories we consider are image presentations, omitted images, and image changes. With the exception of behavioral feature category, these are the same categories considered in ref. <sup>24</sup>. Also note that the image change feature category can no longer be divided into behavior-dependent features representing hits and misses.

We thus define the ten time-dependent features vectors;  $f_t^{\gamma}$  for  $\gamma = \text{image1}$ , image2, ..., image8, omission, change; to have value 1 at the onset of a given feature and to be 0 otherwise (Figure 6A, top). These features each belong to one of three feature categories;  $\alpha = \text{image,omission,change}$ ; with the eight image features belonging to the 'image' feature category, and the omission and change features belonging to the category of the same name.

For each cell *i* in each 'imaging' session, we fit its full session response (post smoothing, see above),  $y_{i,t}$ , using time-dependent feature kernels,  $k_{i,t}^{\gamma}$ , such that an estimate of its response is given by the convolution

$$\widehat{y}_{i,t} = \sum_{\gamma} (f^{\gamma} * k_i^{\gamma})_t + c_i,$$
 (Equation 43)

where  $c_i$  is bias term. Each kernel's width in time is matched to that used in ref. <sup>24</sup>: the image, omission, and change kernels persist for 0.75, 3.0, and 2.25 s after the corresponding feature onset, respectively.

The kernels  $k_{i,t}^{\gamma}$  and bias terms  $c_i$  are fit using ordinary least squares regression with an L2 penalty (i.e., ridge regression, see ref. <sup>24</sup> for additional details). We evaluate the fit of the models by computing their *variance explained* on a test set,

$$VE_{i} = 1 - \frac{Var_{\mathcal{T}}(y_{i} - \widehat{y}_{i})}{Var_{\mathcal{T}}(y_{i})}, Var_{\mathcal{T}}(y_{i}) = \frac{1}{|\mathcal{T}|} \sum_{t \in \mathcal{T}} (y_{i,t} - \overline{y}_{i})^{2},$$
 (Equation 44)

where  $\overline{y}_i$  is the cell's mean activity over the entire imaging session. Here, the  $\mathcal{T}$  subscript on  $Var_{\mathcal{T}}$  indicates the subset of sequence times over which the variance is computed and  $|\mathcal{T}|$  represents the number of time steps (see below). To find the optimal L2 regularization, we scan over regularization coefficients, evaluate said fits, and choose the regularization that yielded the highest mean variance explained across the entire cell population. Train/test splits are computed over distinct batches.

Since certain feature categories are quite sparse across the full input sequence (e.g., omissions and changes), their corresponding feature kernels influence only a small subset of sequence time steps. To account for the different possible kernel coverage over the entire sequence, below it will be useful to compute the variance explained over only the subset of sequence time steps where a given feature category's kernel(s) could have possibly had an influence. Let  $\mathcal{T}^{\alpha}$  be the set of time steps a feature category's kernel(s) could



have possibly influenced the response given the sequence's feature vectors  $f_t^{\gamma}$  and the kernel widths in time. We define the *adjusted* variance explained as

$$\overline{\mathsf{VE}}_{i}^{\alpha} = 1 - \frac{\mathsf{Var}_{\mathcal{T}^{\alpha}}(y_{i,t} - \widehat{y}_{i,t})}{\mathsf{Var}_{\mathcal{T}^{\alpha}}(y_{i,t})}, \tag{Equation 45}$$

where the variance is now only computed on the subset of sequence times  $\mathcal{T}^{\alpha}$ . Since  $|\mathcal{T}^{\alpha}|$  is the number of time steps in  $\mathcal{T}^{\alpha}$  and  $|\mathcal{T}|$  is the total number of sequence time steps in the session,  $|\mathcal{T}^{\alpha}| < |\mathcal{T}|$ , for all three feature categories we consider. Specifically,  $|\mathcal{T}^{\alpha}| < |\mathcal{T}| \approx 0.95, 0.19$ , and 0.16 for the image, omission, and change categories, respectively. Lastly, note that from our definition in Equation 44, the adjusted variance is always computed relative to the mean cell activity over the entire session.

#### **Coding scores**

For each cell in each session, we compute its *coding score* with respect to each of the three feature categories we introduced above. Intuitively, a category's coding score represents how important its feature(s) are for fitting the cell's response. To compute coding score, we compare the cell's adjusted variance explained of a model fit *without* a given feature category's kernel(s) to the model fit with all kernels. Explicitly, the raw coding score is defined as

$$c_i^{\alpha} = \frac{\overline{VE}_{i,\text{full}}^{\alpha} - \overline{VE}_{i,\text{sans }\alpha}^{\alpha}}{\overline{VE}_{i,\text{full}}^{\alpha}},$$
 (Equation 46)

where  $\overline{\mathsf{VE}}_{i,\mathrm{full}}^{\alpha}$  and  $\overline{\mathsf{VE}}_{i,\mathrm{sans}\ \alpha}^{\alpha}$  are the adjusted variance explained of the models fit with all kernels and all kernels except those belonging to category  $\alpha$ , respectively.

Finally, it will sometimes be useful to compare coding scores across sessions, in which case we want to normalize all coding scores on equal footing. The across-session coding score of session s is defined as

$$\overline{c}_{i}^{\alpha,s} = \left(\frac{\left|\mathcal{T}^{\alpha,S}\right|}{\left|\mathcal{T}^{\alpha,S}\right|}\right) \frac{\overline{V}\overline{\mathsf{E}}_{i,\mathsf{full}}^{\alpha,s} - \overline{V}\overline{\mathsf{E}}_{i,\mathsf{sans}\ \alpha}^{\alpha,s}}{\overline{V}\overline{\mathsf{E}}_{i,\mathsf{full}}^{\alpha,S}}, \text{where} \quad S = \underset{s}{\operatorname{argmax}} \overline{V}\overline{\mathsf{E}}_{i,\mathsf{full}}^{\alpha,s}. \tag{Equation 47}$$

Since we have three feature categories  $\alpha$  and three sessions s, each cell will have a 9-dimensional across-session coding score vector.

In the experiment, a minimum variance explained is required for a nonzero coding score. Specifically,  $\overline{VE}_{i,\text{full}}^{\alpha,s} > 0.005$ . Relative to the full fit variance explained, approximately 54.5%, 34.7%, and 52.1% of VIP cell fits fall under this threshold in the familiar, novel, and novel-plus sessions, respectively. Since the cortical microcircuit model has overall higher variance explained for all cell populations, we adjust this minimum coding score threshold to compensate for the different distribution. Setting a threshold of  $\overline{VE}_{i,\text{full}}^{\alpha,s} > 0.075$  results in similar rates as the experiment, namely 56%, 20%, and 54% of VIP cells across initializations fall under this value.

#### **Cell clustering**

We use spectral clustering to cluster the set of  $\overline{c}_i^{\alpha,s}$  for each population. In this subsection, we use  $\mathbf{c}_i$  to denote the 9-dimensional coding score vector of cell i.

To compute the ideal number of clusters for cell population, we use two measures: the gap statistic and the eigengap. For the gap statistic, we scan over cluster sizes from k=2 to 15. We use the SpectralClustering method from scikit-learn with default parameters and a given k to fit the data and compute the pairwise Euclidian distance within each cluster. Let the n th cluster contain the set of cells indexed by  $i_n$ . Then,

$$\overline{D}(k) = \frac{1}{k} \sum_{n=1}^{k} \sum_{i_n \neq j_n} d_{i_n j_n} \quad d_{ij} = \|\mathbf{c}_i - \mathbf{c}_j\|_2.$$
 (Equation 48)

This metric is computed for the actual clusters and compared to a baseline of shuffled data. The shuffled data is the across-session coding scores shuffled across experience-level and feature categories. For the metric over the shuffled data  $\overline{D}_s(k)$ , the gap statistic is then  $\overline{D}_s(k) - \overline{D}(k)$ , and the optimal k is the one that maximizes this metric (Figure S8E).

To compute the *eigengap*, we compute differences in consecutive (ordered) eigenvalues of the Laplacian of the coding score's affinity matrix. Specifically, the affinity matrix has elements  $e^{-\gamma \sigma_{a,b}^2}$ , with  $d_{a,b}$  the Euclidean distance computed above. The eigengap is then the difference in eigenvalues of the Laplacian, where large gaps are associated with sudden changes of the amount of similarity explained by additional cluster partitions (Figure S8D).

Once the optimal number of clusters is computed, we perform spectral clustering on the set of coding score vectors of a given population for 150 different initial seeds. Across all these fits, we compute the symmetric matrix of co-cluster probabilities for all cell pairs. This co-clustering matrix is then passed through scikit-learn's AgglomerativeClustering method, again with





the optimal number of clusters as determined above, with default parameters except for affinity='euclidean' and linka-ge='average'. Finally, cell clustered are ordered by mean across-session coding scores, with the clusters with the smallest mean being ordered first.

#### **Additional Figure details**

The equivalent of Figure 1 for the pre-only dependent update rules is given in Figure S1A.

The equivalent of Figure 2c for strengthened modulations is shown in Figure S1B.

For Figures 2D–2G], an exemplar FMSN was trained using a set of 8 familiar stimuli. To ensure an equal distribution of the familiar stimuli over the training time, a training schedule where each of the 8 familiar stimuli is shown every 8 inputs is used. Note the order of the 8 examples is shuffled within every 8 inputs. The FMSN has a population size of 300 input and 500 output neurons. As mentioned in the main text, we take all input neurons to be excitatory. The distribution of **W** elements is taken to have w = 1/300 and  $p_W = 0.2$ , see Equation 8. For the stimuli, the sparse random binary vector population has  $p_{\text{stim}} = 0.05$ , with a minimum number of nonzero inputs of 1, and nonzero elements of size 0.15.(Any two vectors drawn from the sparse random binary vector distribution we consider in this example have cosine similarity of 0.14. Cosine similarity of stimuli decreases with increased sparsity and larger input dimension.) The standard deviation of the Gaussian noise added to the inputs is taken to be 0.1 times the size of the nonzero elements. We threshold is adjusted at initialization such that the output population has a firing rate of 30%. The associative modulations obey the bounds discussed below Equation (7). We take  $\eta = -5 \times 10^2$  and  $\tau_{\text{decay}}/\Delta t = 2 \times 10^4$  so that the modulations undergo practically no decay during the stimulus learning period. These parameters are used in FMSNs throughout this work unless otherwise stated.

Note in order to track both the familiar and novel output activity throughout training, we treat them as "test sets" when we pass them to the network, which distinguishes them from the sequential training set we use to change the modulations of a network. For any input that belongs to the test set, we do not update the synapses. In this way, we can understand what the network's response to these various stimuli without would be actually updating the network's modulations as if it truly "saw" the stimuli during training. The full familiar and novel sets were treated as test sets in order to track their output activity over training shown in Figure 2E.

In Figure 2G we introduce the idea of an *important synapse*. An important synapse is stimulus dependent, and as such a given synapse can be important for multiple stimuli. Important synapses are also defined before any modulations in the network occur, and are thus independent of the FMS mechanism. We define important synapses in our model as those synapses that satisfy two requirements: (1) there must be a synapse there and (2) the synapse's pre- and postsynaptic neurons both fire when the stimulus is input into the network (without modulations). Formally, for a given stimulus input  $\mathbf{x}$ , let  $\mathbf{y}$  be the corresponding activity of the output layer *without any synapse modulations from FMSs*. For example, for the FMSN,  $\mathbf{y} = \varphi(\mathbf{W}\mathbf{x} + \mathbf{b})$ , but this generalizes to other possible postsynaptic expressions. The mask  $\mathbf{S}$  that defines the important synapses contained within  $\mathbf{W}$  for stimulus  $\mathbf{x}$  is given by

$$S_{li} = \begin{cases} 1 & W_{li}y_ix_i \neq 0, \\ 0 & \text{otherwise.} \end{cases}$$
 (Equation 49)

Intuitively, whether or not a synapse is important tells us whether or not said synapse would be modulated from an associative FMS mechanism. In practice, since we add noise to all stimuli being passed through the FMSN, many synapses that are not important are also modulated.

Figure 3A quantifies how a vector's distance from the familiar subspace influences its output magnitude. The familiar subspace is defined as the subspace spanned by the familiar set of vectors. To measure the distance of any random vector  $\mathbf{v} \in \mathbb{R}^d$  to this subspace, we orthonormalized the familiar set to obtain the matrix  $\tilde{\mathbf{F}} \in \mathbb{R}^{8 \times d}$  and then formed the projection matrix onto the familiar subspace via  $\mathbf{P}^F = \tilde{\mathbf{F}}(\tilde{\mathbf{F}}^T\tilde{\mathbf{F}})^{-1}\tilde{\mathbf{F}} \in \mathbb{R}^{d \times d}$ . A given vector's distance from the familiar subspace is then measured by calculating the cosine similarity of the vector and its projection onto the familiar subspace,

$$d^{\mathsf{F}}(\mathbf{v}) = \frac{\mathbf{v} \cdot \mathbf{P}^{\mathsf{F}} \mathbf{v}}{\|\mathbf{v}\|_{2} \|\mathbf{P}^{\mathsf{F}} \mathbf{v}\|_{2}}.$$
 (Equation 50)

By this definition, any vector from the familiar set or any linear combination thereof has  $d^F = 1.0$ . Note that noise is added to all vectors before being passed through the network, and this noisy vector was used to compute  $d^F$  in Figure 3A. Hence, even the noisy familiar vectors do not lie exactly in the familiar subspace.

Since simply drawing from the same distribution that generated the familiar and novel sets almost always generates stimuli that are far from the familiar subspace, we generated inputs as follows. We first drew a vector  $\mathbf{v}'$  from said distribution and also for each vector drew  $f \sim \mathcal{U}(0,1)$  and a random vector from the familiar set  $\mathbf{f}$ . Then, each element of the final vector  $\mathbf{v}$  has probability f to be the same as  $\mathbf{f}$ , and is otherwise equal to the corresponding element of  $\mathbf{v}'$ . In this way, as f varies from 0 to 1, we interpolate between vectors that are drawn from the original distribution (i.e., the novel set) to those in the familiar set. Finally, to generate random vectors in the familiar subspace, 8 binary weights were drawn, their sum was normalized to 1, and these weights determined the linear combination of familiar vectors that formed the new vector within its subspace.

#### **Article**



Figures 3B and 3C show how the evolution of the modulation matrix can change with different learning rates,  $\eta$ , and decay rates,  $\lambda$ . Both setups use a more specialized learning schedule than those in Figure 2. Figure 3B consists of a single input vector passed over and over again to observe how quickly modulations can grow and when they saturate. Figure 3C consists of a single input vector passed on the first time step and then only noise afterward. The goal of this plot is to observe how quickly the modulations created by a single input can shrink over time.

Figure 3D shows the average result of several KS tests on a network as we scan over number of exposures and the learning rate,  $\eta$ . With the exception of the parameters scanned over, the network and training parameters used in this setup are identical to that in Figures 2D-2G. KS test results are averaged of  $log_{10}(p)$  over 10 distinct network/stimuli initializations. Figures 3E-3F are the same as Figure 3D, but now scan over the FMS's decay time constant and learning rate.

In Figure 4, additional details of the experimental training procedure can be found in ref. 24. Details of the model network (Figure 2C) and task (Figure 2D) can be found in the STAR Methods.

In Figure 4G, for synaptic matrix  $\mathbf{W}^{p,p'}$  between presynaptic population p and postsynaptic p', the model connection strength values were computed via

$$\frac{1}{N_p} \sum_{i,l} W_{il}^{p,p'}.$$
 (Equation 51)

Note here p and p' are treated on unequal footing to account for the fact that, for a given postsynaptic cell, the relative count of presynaptic inputs contributes to its total activity. The theoretical values were computed via Equation 16.

In Figure 5, unless otherwise stated, we take the cortical circuit model to have 400 excitatory, 200 VIP, and 200 SST neurons, though see STAR Methods for how synaptic strength is adjusted to compensate for deviations from realistic cell counts. Weights are initialized as described in the STAR Methods, with multiplicative constant c = 0.18. The three FMS learning rates are scanned over to determine the best fit to experimental responses, see STAR Methods.

Figures 5A and 5H show comparisons of the mean responses of our cortical microcircuit model and responses measured in experiment. We match event traces that are smoothed by a half-normal filter, see above. See Ref. 24 for significantly more details on the experimental details including details about the event extraction.

To extract the mean responses of the model, let the set of all times of interest (e.g., for image changes or omissions) be denoted by  $\mathcal{T}$ . We denote the mean response as

$$Z_t^p = \sum_{S \in \mathcal{T}} \sum_{i=1}^N Y_{i,t+S}^p$$
 (Equation 52)

The full familiar and novel set responses of Figure 5B are gathered similarly to the test sets of the FMSN. That is, a test set consisting of the stimuli from the familiar and novel sets representing the image changes is passed to the network at particular times during training. No temporal history or time-correlated noise input is passed to the network in these test sets so that the VIP's change in response to the present input from  $FMS_A$  can be isolated (the temporal history response also changes during training from  $FMS_{AH}$ ). Once again, we do not update the network's modulations during these passes, and so the microcircuit has no memory of viewing these stimuli that are solely used to monitor training progress. The test set is evaluated at every image change during training, specifically at the step corresponding to the peak of the smoothing kernel (see above). The modulation magnitudes of Figure 5C are computed analogously to Figure 2F and also measured at each image change during training.

Figures 5D, 5G, and 5K all show the mean responses as a particular FMS learning rate is varied. Networks and tasks are initialized identically to those used to produce the analogous figures in Figures 5A and 5H, the only thing that changes is the corresponding FMS's learning rate and thus its asymptotic modulations.

The responses in Figure 5E show the VIP response of a test set, evaluated in an identical manner to that described for Figure 5B above. The test set now consists of all novel stimuli with the temporal history part of the input corresponding to the zero-time encoding and no noise input (since the FMSAH modulations have stabilized, they no longer significantly affect the change in VIP responses for the relevant timescales shown here). Test responses are again gathered at the step corresponding to the peak of the smoothing kernel. The responses shown are averaged over 50 image changes measured during a novel session. Additionally, the responses are normalized by the largest mean response. Figure 5F shows the modulations of FMS<sub>P</sub> during an exemplar image change of the novel imaging period. The modulation magnitudes shown are gathered at every time step.

Figure 5I shows an exemplar mean VIP response to all the encodings of time-since-last image from the stimulus history input. These responses are once again gathered as a test set (see above), where now the present stimulus and time-correlated noise inputs are taken to be zero so that the VIP's change in response from FMSAH can be isolated. Responses are gathered before and after training on the familiar image set. The modulation magnitudes shown in Figure 5J are gathered analogously to those in Figure 5C.

Figure 6A shows exemplar input features and fits over time. The top subfigure dots correspond to when the corresponding input features is on, see STAR Methods for additional details. The bottom subfigure shows exemplar raw cell data, full kernel fit, and fit without the image kernels. Figure 6B shows the clustered VIP across-session coding scores from experiment.<sup>24</sup> The middle plot shows cluster-average coding scores and the right plot shows average coding scores across all VIP cells. Figure 6C shows the clustered VIP coding score from the model, see STAR Methods for details of how these are computed. Clustered are ordered by smallest





mean coding score to largest. Figure 6D shows image kernel fits from the full kernel regression model. Image kernel fits are averaged across all 8 image kernels, all VIP cells, and all initializations. Figure 6E shows fits and raw data of the cluster-averaged novel image (across-session) coding score as a function of a cluster-averaged network property, see details of Figures S9 and S10 for full details of this network property and all others. Fits are done using linear regression and we use the resulting correlation to measure how much a network property influences the value of a given across-session coding score. This plot shows only one exemplar network property and one coding score, the resulting correlations across all 16 network properties we investigate and all 9 coding scores can be found in Figure S9. In Figure 6F we show the median correlation for the familiar and novel image coding scores as a function of two network properties for both the cluster-averaged values and the raw cell data. Figure 6G shows the amount of familiar and novel input cells belonging to a particular cluster receive. Each point represents a single cluster for a given initialization. Points are colored by whether they are familiar-coded (clusters 3), novel coded (clusters 2, 5, 7), both familiar and novel coded (clusters 6, 8), or not image coded (clusters 1, 4). Figure 6H is generated similarly, with clusters colored by whether they are omission-coded (clusters 4, 7, 8) or not (all other clusters).