# Hebbian Learning in a Random Network Captures Selectivity Properties of Prefrontal Cortex

Grace W. Lindsay[a,e,g], Mattia Rigotti[a,b], Melissa R. Warden[c], Earl K. Miller[d], Stefano Fusi[a,e,f]

[a] Center for Theoretical Neuroscience, College of Physicians and Surgeons, Columbia University, New York, New York, USA

[b] IBM T.J. Watson Research Center, 1101 Kitchawan Rd., Yorktown Heights, NY, USA

[c] Neurobiology and Behavior, College of Agriculture and Life Sciences, Cornell University, Ithaca, NY, USA

[d] The Picower Institute for Learning and Memory & Department of Brain and Cognitive Sciences, Massachusetts Institute of Technology, USA

[e] Mortimer B. Zuckerman Mind Brain Behavior Institute, College of Physicians and Surgeons, Columbia University, New York, New York, USA

[f] Kavli Institute for Brain Sciences, Columbia University, New York, New York, USA

[g] Corresponding Author: gracewlindsay@gmail.com

## Abstract

Complex cognitive behaviors, such as context-switching and rule-following, are thought to be supported by prefrontal cortex (PFC). Neural activity in PFC must thus be specialized to specific tasks while retaining flexibility. Nonlinear 'mixed' selectivity is an important neurophysiological trait for enabling complex and context-dependent behaviors. Here we investigate (1) the extent to which PFC exhibits computationally-relevant properties such as mixed selectivity and (2) how such properties could arise via circuit mechanisms. We show that PFC cells recorded during a complex task show a moderate level of specialization and structure that is not replicated by a model wherein cells receive random feedforward inputs. While random connectivity can be effective at generating mixed selectivity, the data shows significantly more mixed selectivity than predicted by a model with otherwise matched parameters. A simple Hebbian learning rule applied to the random connectivity, however, increases mixed selectivity and allows the model to match the data more accurately. To explain how learning achieves this, we provide analysis along with a clear geometric interpretation of the impact of learning on selectivity. After learning, the model also matches the data on measures of noise, response density, clustering, and the distribution of selectivities. Of two styles of Hebbian learning tested, the simpler and more biologically plausible option better matches the data. These modeling results give intuition about how neural properties important for cognition can arise in a circuit and make clear experimental predictions regarding how various measures of selectivity would evolve during animal training.

**Significance Statement**: Prefrontal cortex (PFC) is a brain region believed to support the ability of animals to engage in complex behavior. How neurons in this area respond to stimuli—and in particular, to combinations of stimuli ("mixed selectivity")—is a topic of interest. Despite the fact that models with random feedforward connectivity are capable of creating computationally-relevant mixed selectivity, such a model does not match the levels of mixed selectivity seen in the data analyzed in this study. Adding simple Hebbian learning to the model increases mixed selectivity

to the correct level and makes the model match the data on several other relevant measures. This study thus offers predictions on how mixed selectivity and other properties evolve with training.

---

## 1. Introduction

1 The ability to execute complex, context-dependent behavior is evolutionarily valu-
2 able and ethologically observed [36, 16]. How the brain carries out complex behaviors
3 is thus the topic of many neuroscientific studies. A region of focus is the prefrontal
4 cortex (PFC), [4, 44, 29, 9], as lesion [42] and imaging [28, 6] studies have implied
5 its role in complex cognitive tasks. As a result, several theories have been put forth
6 to explain how PFC can support complexity on the computational and neural levels
7 [29, 46, 11].
8 Observing the selectivity profiles of its constituent cells is a common way to inves-
9 tigate a neural population's role in a computation. In its simplest form, this involves
10 modeling a neuron's firing rate as a function of a single stimulus, or, perhaps, an addi-
11 tive function of multiple stimuli [39, 8, 30]. More recently, however, the role of neurons
12 that combine inputs in a nonlinear way has been investigated [38, 23, 41, 32, 25, 35, 11],
13 often in PFC. Rather than responding only to changes in one input, or to changes in
14 multiple inputs in a linear way, neurons with nonlinear mixed selectivity have firing
15 rate responses that are a nonlinear function of two or more inputs (Figure 1B). Cells
16 with this selectivity (which we just call "mixed") are important for population coding
17 because of their effect on the dimensionality of the representation: they increase the
18 dimensionality of the population response, which increases the number of patterns that
19 a linear classifier can read out. This means that arbitrary combinations of inputs can
20 be mapped to arbitrary outputs. In relation to complex behaviors, mixed selectivity
21 allows for a change in context, for example, to lead to different behavioral outputs,
22 even if stimulus inputs are the same. For more on the benefits of mixed selectivity, see
23 [11].
24 Theoretical work on how these properties can arise on a circuit level shows that
25 random connectivity is surprisingly efficient at increasing the dimensionality of the
26 neural representation [15, 22, 7, 37, 2, 1, 20]. This means that mixed selectivity can be
27 observed even without learning. However, learning can greatly improve the ability of
28 a linear readout to generalize and hence to make the readout response more robust to
29 noise and variations in the sensory inputs (see e.g. [11]). The ideal situation would be
30 one in which a neural population represents only the task relevant variables and the
31 representation has the maximal dimensionality. In brain areas like PFC, where there
32 is a huge convergence of inputs from many other brain areas, it might be important
33 to bias the mixed selectivity representations toward the task relevant variables, which
34 can be achieved only with learning.
35 In this study, we characterize the response of a population of PFC cells in terms of
36 the distribution of linear and nonlinear selectivity, the response density, and the clus-
37 tering of selectivities. All these properties characterize the dimensionality of neural
38 representations and are important for the readout performance. As described above,
39 nonlinear mixed selectivity is important for increasing dimensionality. High dimension-
40 ality, however, also requires a diversity of responses. We studied this by determining
41 how the preference to different stimuli are distributed across the population. In some
42 lower sensory areas, cells tend to be categorizable—that is, there are groups of cells

that display similar preference profiles [14]. More associative areas tend to lose this clustering of cell types. Such categories may be useful when an area is specialized for a given task, but diversity is needed for flexibility [35].

After characterizing the PFC response, we show that a model with random connectivity can only partially explain the PFC representation. However, with a relatively small deviation from random connectivity—obtained with a simple form of Hebbian learning that is characterized by only two parameters—the model describes the data significantly better.

## 2. Methods

### 2.1. Task Design

The data used in this study comes from previously published work [43]. In brief, two monkeys performed two variants of a delayed match-to-sample task (Figure 1A). In both task types, after initial fixation, two image cues (chosen from four possible) were presented in sequence for 500ms each with a 1000ms delay period in between the first and second cue. After a second delay period also lasting 1000ms, one of two events occurred, depending on the task type. In the recognition task, another sequence of two images were shown and the monkey was instructed to release a bar if this test sequence matched the initial sample sequence. In the recall task, an array of three images appeared on the screen, and the monkey had to saccade to the two images from the sample sequence in the correct order. Blocks of recall and recognition tasks were interleaved during each recording session. Given that each sequence had two different image cues chosen from the four total image identity options and that there were two task types, the total number of conditions was 4 x 3 x 2 = 24.

### 2.2. Neural Data

Recordings were made using grids with 1 mm spacing (Crist Instrument) and custom-made independently moveable microdrives to lower eight dura-puncturing Epoxylite-coated tungsten microelectrodes (FHC) until single neurons were isolated. Cells were recorded from two adult rhesus monkeys (Macaca mulatta), one female and one male, and combined for analysis. No attempt was made to pre-screen neurons, and a total of 248 neurons were recorded (with each neuron observed under both task types).

For the purposes of this study, firing rates for each neuron were calculated as the total number of spikes during the later 900ms of the second delay period, as it was at this point that the identities of all task variables were known. Any cells that did not have at least 10 trials for each condition or did not have a mean firing rate of at least 1 spike/sec as averaged over all trials and conditions were discarded. This left 90 cells.

### 2.3. Fano Factor Measurements

Noise is an important variable when measuring selectivity. High noise levels require stronger tuning signals in order to be useful for downstream areas, and to reach significance in statistical testing. Thus, any model attempting to match the selectivity profile of a population must be constrained to have the same level of noise. Here, we measure noise as the Fano Factor (variance divided by mean) of each cell's activity across trials for each condition (spike count taken from later 900ms of the two-object delay). This gives 24 values per cell. This is the trial Fano Factor. Averaging over conditions gives one trial Fano Factor value per cell, and averaging over cells gives a

3

single number representing the average noise level of the network. Unless otherwise stated, $FF_T$ refers to this network averaged measure.

Another measure of interest is how a neuron's response is distributed across conditions. Do neurons respond differentially to a small number of conditions (i.e., a sparse response), or is the distribution more flat? To measure this, the firing rate for each condition (averaged across trials) was calculated for each neuron and the Fano Factor was calculated across conditions. In this case, a large Fano Factor means that some conditions elicit a very different response than others, while a small Fano Factor suggests the responses across conditions are more similar. Averaging across all cells gives the condition Fano Factor of the network, or $FF_C$.

See Figure 1C for a visualization of these measures in an example neuron.

## 2.4. Selectivity Measurements

A neuron is selective to a task variable if its firing rate is significantly affected by that the identity of that task variable. In this task, each condition contains three task variables: task type (recall or recognition), the identity of the first cue, and the identity of the second cue. Therefore, we used a 3-way ANOVA to determine if a given neuron's firing rate was significantly (p<.05) affected by a task variable or combination of task variables. Selectivity can be of two types: pure or nonlinearly mixed (referred to as just "mixed"), based on which terms in the ANOVA are significant. If a neuron has a significant effect from one of the task variables, for example, it would have pure selectivity to that variable. Interaction terms in the ANOVA represent nonlinear effects from combinations of variables. Therefore, any neurons that have significant contributions from interaction terms as determined by the ANOVA have nonlinear mixed selectivity. So, for example, if a neuron's firing rate can be written as $FR = f(X_{TT}, X_{C2}, X_{TTC1}, b)$, that neuron has pure selectivity to task type (TT), pure selectivity to cue 2 (C2) and mixed selectivity to the combination of task type and cue 1 (TTC1), with $b$ as a bias term and $f$ a linear function of its arguments. Note that having pure selectivity to two or more task variables is not the same as having nonlinear mixed selectivity to a combination of those task variables.

## 2.5. Clustering Measurement

Beyond the numbers of neurons selective to different task variables, an understanding of how preferences to task variable identities cluster can inform network models. For this, we use a method that is inspired by the Projection Angle Index of Response Similarity (PAIRS) measurement as described in [35]. For this measure each neuron is treated as a vector in selectivity space, where the dimensions represent preference to a given task variable identity (Figure 1D). To get these values, neuronal responses are fit with a general linear model (GLM) to find which task variable identities significantly contribute to the firing rate. Note that this gives a beta coefficient value for task variable identities, such as cue 1=A, rather than just each task variable, such as cue 1. It does not include interaction terms. The reason for this is that, given the relatively low number of trials, the high dimensional full GLM model would be difficult to confidently fit. Furthermore, analysis of clustering in a high-dimensional space with a relatively small number of neurons would be difficult to interpret. The beta values found for each cell via this method are shown in Figure 3C (non-significant coefficients—those with p>.05—are set to 0).
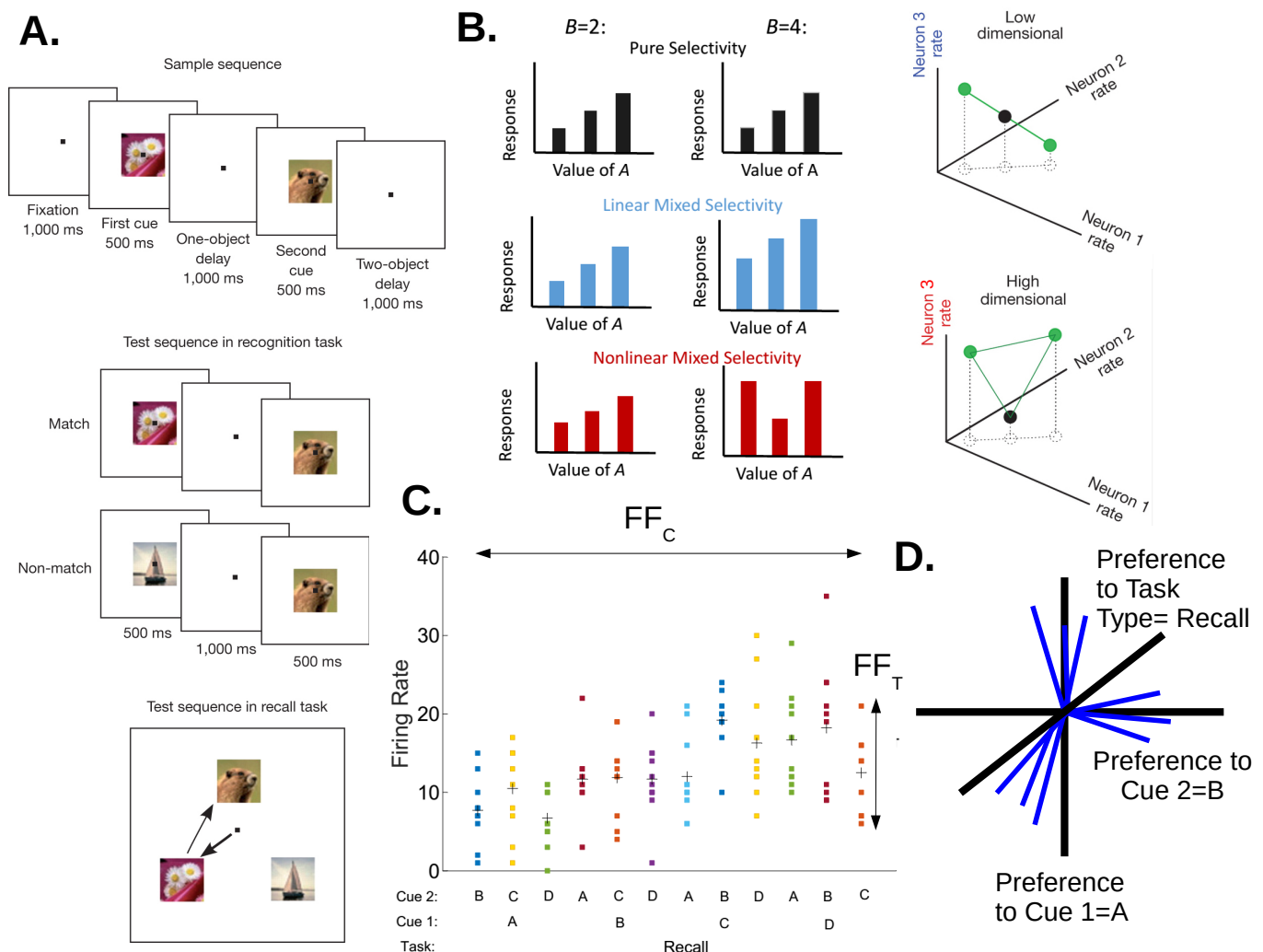
4

Figure 1: Description of prefrontal cortex data and relevant measures of selectivity A.) Task Design. In both task types, the animal fixated as two image cues were shown in sequence. After a delay the animal had to either indicate that a second presented sequence matched the first or not ("recognition") or saccade to the two images in correct order from a selection of three images ("recall"). B.) What nonlinear mixed selectivity can look like in neural responses and its impact on computation. The bar graphs on the left depict three different imagined neurons and their responses to combinations of two task variables A and B. The black neuron has selectivity only to A, as its responses are invariant to changes in B. The blue neuron has linear mixed selectivity to A and B: its responses to different values of A are affected by the value of B, but in a purely additive way. The red neuron has nonlinear mixed selectivity: its responses to A are impacted nonlinearly by a change in the value of B. The figures on the right show how including a cell with nonlinear mixed selectivity in a population increases the dimensionality of the representation. With the nonlinearly-selective cell (bottom), the black dot can be separated with a line from the green dots. Without it (top), it cannot. C.) A depiction of measures of trial-to-trial noise ($FF_T$) and the distribution of responses across conditions ($FF_C$). The x-axis labels the condition, each dot is the firing rate for an individual trial and the crosses are condition means used for calculating $FF_C$ (data from a real neuron; recognition task not shown). D.) Conceptual depiction of the clustering measure. Each cell was represented as a vector (blue) in a space wherein the axes (black) represent preference for task variable identities, as determined by the coefficients from a GLM (only three are shown here). The clustering measure determines if these vectors are uniformly distributed.

The coefficients derived from the GLM define a vector in a 7-D vector space for each neuron (see Figure 1D for a schematic). This clustering method compares the distribution of vectors generated by the data to a uniform distribution on the hypersphere in order to determine if certain combinations of selectivities are more common than expected by chance. In [35] this comparison is done by first computing the average angle between a given vector and its k nearest neighbors and seeing if the distribution of those values differs between the data and a random population.

That approach is less reliable in higher dimensions, therefore we use the Bingham test instead [24]. The Bingham test calculates a test statistic: $S = \frac{p(p+2)}{2}n(Tr(\mathbf{T}^2) - \frac{1}{p})$. This statistic, which we refer to as the clustering value, measures the extent to which the scatter matrix, $\mathbf{T}$, (an approximation of the covariance matrix) differs from the identity matrix (scaled by $1/p$), where $p$ and $n$ are the dimensions of the selectivity space (7) and the number of cells (90), respectively. The higher this value is, the more the data deviates from a random population of vectors wherein selectivity values are IID. Thus, a high value suggests that neurons in the population cluster according to task variable identity preferences. In order to put this clustering value into context we compared the value found from the data to two distributions: one generated by shuffled data and one generated from data designed to be highly clustered. For the shuffled data, we created "fake" cell vectors by shuffling the selectivity values across all cells. For the clustered data, we created 3 categories of fake cells, each defined by pure selectivity to two specific task variable identities. A population of 90 cells was created by combining 30 cells from each category (the population was also designed to have the same average firing rate and $FF_T$ of the data). This results in a population that has 3 clear clusters of cell types in selectivity space. 100 populations based on each type of fake data were created in order to generate distributions that represent random and clustered data.

Using the Gine-Ajne test of uniformity on the hypersphere ([13]) gives very similar results to the Bingham test results.

## 2.6. Circuit Model

To explore the circuit mechanisms behind PFC selectivity, we built a simple two-layer neural model, modeled off of previous work [2] (see Figure 4A for a diagram). The first layer consists of populations of binary neurons, with each population representing a task variable identity. To replicate a given condition, the populations associated with the task variable identities of that condition are turned on (set to 1) and all other populations are off (set to 0). Each population has a baseline of 50 neurons. To capture the biases in selectivities found in this dataset (particularly the fact that, in the 900ms period we used for this analysis, many more cells show selectivity to task type than cue 2 and to cue 2 than cue 1), the number of neurons in the task type and cue 2 populations are scaled by factors that reflect these biases (80 cells in each task type population and 60 in each cue 2 population). The exact values of these weightings do not have a significant impact on properties of interest in the model.

The second layer represents PFC cells. These cells get weighted input from a subset of the first layer cells. Cells from the input layer to the PFC layer are connected with probability .25 (unless otherwise stated), and weights for the existing connections are drawn from a Gaussian distribution ($\mu_W = .207$, and $\sigma_W = \mu_W$ unless otherwise stated. Because negative weights are set to 0, the actual connection probability and $\sigma_W$ may be slightly lower than given).

6

The activity of a PFC cell on each trial, $t$, is a sigmoidal function of the sum of its inputs:

$$r_i^t = k\phi(\sum_j w_{ij}x_j^t + \epsilon_A^t - \Theta_i)$$

$$\phi(z) = \frac{1}{1 + e^{-z}} \tag{1}$$

$$\epsilon_A^t \sim \mathcal{N}(0, \sigma_A{}^2) \qquad \sigma_A = a\mu_W$$

where $x_j$ is the activity (0 or 1) of the $j^{th}$ input neuron and $w_{ij}$ is the weight from the $j^{th}$ input neuron to the $i^{th}$ output neuron. $\Theta_i$ is the threshold for the $i^{th}$ output neuron, which is calculated as a percentage of the total input it receives: $\Theta_i = \lambda\Sigma_j w_{ij}$. The $\lambda$ value is constant across all cells, making $\Theta$ cell-dependent. $k$ scales the responses so that the average model firing rate matches that of the data.

Two sources of noise are used to model trial-to-trial variability. $\epsilon_A$ is an additive synaptic noise term drawn independently on each trial for each cell from a Gaussian distribution with mean zero. The standard deviation for this distribution is controlled by the parameter $a$, which defines $\sigma_A$ in units of the mean of the weight distribution, $\mu_W$. The second noise source is multiplicative and depends on the activity of a given cell on each trial:

$$y_i^t \sim \mathcal{N}(r_i^t, \sigma_{M_i}^{t2})$$

$$\sigma_{M_i}^t = mr_i^t \tag{2}$$

Thus, the final activity of an output PFC cell on each trial, $y_i^t$, is drawn from a Gaussian with a standard deviation that is a function of $r_i^t$. This standard deviation is controlled by the parameter $m$. Both $m$ and $a$ are fit to make the model $FF_T$ match that of the data.

To make the model as comparable to the data as possible, ten trials are run for each condition and 90 model PFC cells are used for inclusion in the analysis.

*2.7. Hebbian Learning*

A simplified version of Hebbian learning is implemented in the network in a manner that captures the "rich get richer" nature of Hebbian learning while keeping the overall input to an individual cell constant. In traditional Hebbian learning, weight updates are a function of the activity levels of the pre- and post-synaptic neurons: $\Delta w_{ij} = g(x_j, y_i)$. In this simplified model we use connection strength as a proxy for joint activity levels: $\Delta w_{ij} = g(w_{ij})$. We also implement a weight normalization procedure so that the total input to a cell remains constant as weights change.

To do this, we first calculate the total amount of input each output cell, $i$, receives from each input population, $p$:

$$I_i^p = \sum_{j \in p} w_{ij} \tag{3}$$

The input populations (each corresponding to one task variable identity) are then

7

209 ranked according to this value. The top $N_L$ populations according to this ranking
210 (that is, those with the strongest inputs onto the output cell) have the weights from
211 their constituent cells increased according to:

$$w_{ij} = (1 + \eta)w_{ij}, \quad j \in P_{1:N_L} \tag{4}$$

212 where $\eta$ is the learning rate (set to .2 unless otherwise stated). After this, all weights
213 into the cell are normalized via:

$$\mathbf{w}_i = \mathbf{w}_i \frac{\sum_{p=1}^{P} I_i^p}{\sum_{j=1}^{J} w_{ij}} \tag{5}$$

214 Note, the numerator in the second term is the sum of all weights into the cell before
215 Eqn. 4 is applied and the denominator is the sum after it is applied.

216     In this work, two versions of Hebbian learning are tested. In the unrestricted, or
217 "free", learning condition described above, the top $N_L$ populations are chosen freely
218 from all input populations (equivalently, all task variable identities) based solely on
219 the total input coming from each population after the random weights are assigned.
220 The alternative, "constrained" learning, is largely the same, but with a constraint
221 on how these top $N_L$ populations are chosen: all task variables must be represented
222 before any can be repeated. So, two populations representing different identities of
223 the same task variable (e.g., cue 1 A and cue 1 B) will not both be included in the
224 $N_L$ populations unless both other task variables already have a population included
225 (which would require that $N_L > 3$). So, with $N_L = 3$, exactly one population from
226 each task variable (task type, cue 1, cue 2) will have weights increased. This variant
227 of the learning procedure was designed to ensure that inputs could be mixed from
228 different task variables, to increase the likelihood that mixed selectivity would arise.
229 Both forms of learning are demonstrated for an example cell in Figure 4B.

230     In both forms of learning, the combination of weight updating and normalization
231 is applied to each cell once per learning step.

232 *2.8. Toy Model Calculations*

233     To make calculations and visualizations of the impacts of learning easier, we use a
234 further simplified toy model (see Figure 8A (left) for a schematic). A cell in this toy
235 model is similar to that in the full model, but instead of a sigmoidal nonlinearity, the
236 heaviside function is used. The toy model has two task variables (T1 and T2) and
237 each task variable has two possible identities (A or B). Four random weights connect
238 these input populations to the output cell: $W_{1A}, W_{1B}, W_{2A}, W_{2B}$. Just as in the full
239 model, on each condition, exactly one task variable identity from each task variable
240 is active (set to 1). This gives four possible conditions, each of which is plotted as a
241 point in the input space in Figure 2. The threshold is denoted by the dotted lines. If
242 the weighted sum of the inputs on a given condition is above the threshold, the cell is
243 active (green), otherwise it is not.

244     The toy model follows the same learning rules defined for the full model. Examples
245 of the impacts of learning on the representation of the 4 conditions are seen in Figure
246 2A and B. In A (top), random weights cause the cell to have pure selectivity to T2.
247 After a learning step that consists of increasing the weights from the two strongest
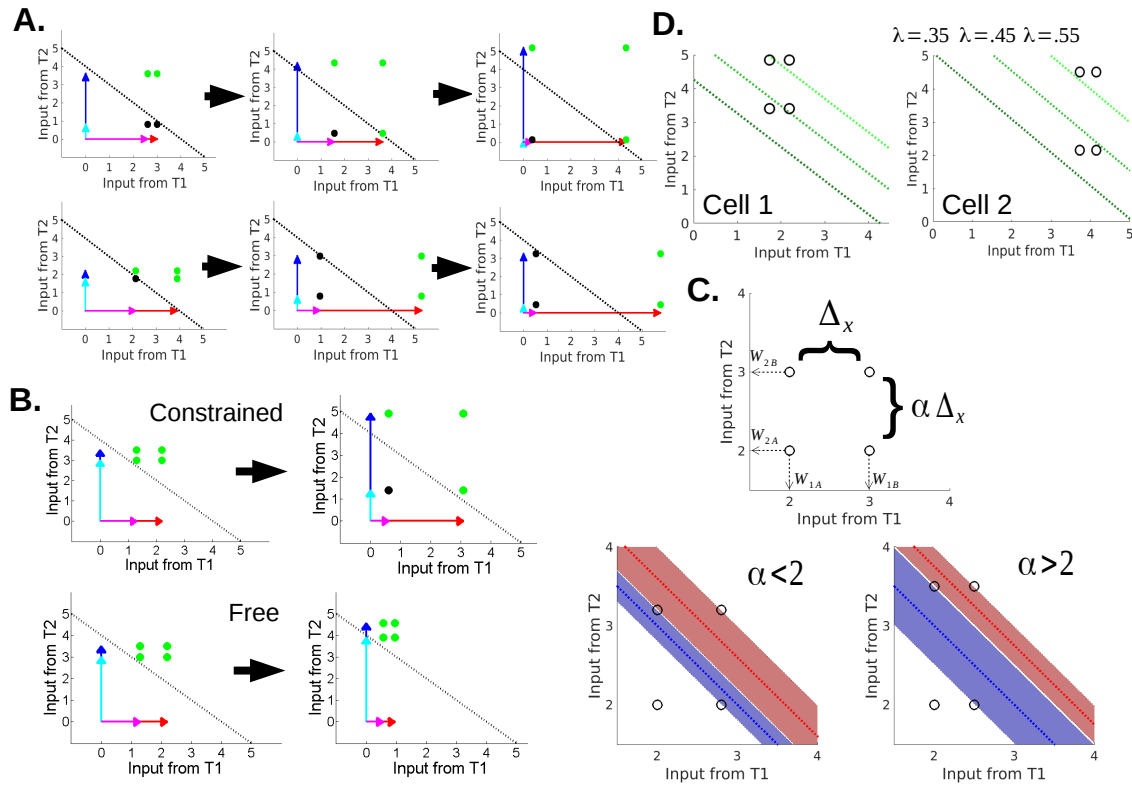248 input populations, T2B and T1B, and then normalizing all weights ($N_L = 2$, learning

8

Figure 2: Signal and noise representation for the toy model shown in Figure 8A. Strength of weights from the 4 input populations are given as arrows in (A and B) and the threshold for the heaviside function is shown as a dotted line. The cell is active for conditions above the threshold (green). Weight arrows omitted for visibility in (C and D). A.) Learning causes the representation of conditions to change. This can change selectivity in multiple ways. Shown here: pure selectivity turns into mixed selectivity (top) and mixed selectivity turns into pure (bottom). B.) Constrained and free learning can lead to different signal changes. Constrained learning (top) guarantees that one population from each task variable is increased. This ensures that the representation spreads out. In this case, the cell goes from no selectivity to mixed selectivity. With these starting weights, free learning increases both populations from T2, and the cell does not gain selectivity. C.) Noise robustness can be thought of as the range of thresholds that can sustain a particular type of selectivity. Relative noise robustness of mixed and pure selectivity depends on the shape of the representation. $\alpha$ is the ratio of the differences between the weights from each task variable (top). In the two figures on the bottom, blue (red) dotted lines show optimal threshold for pure (mixed) selectivity and shaded areas show the range of thresholds created by trialwise additive noise that can exist without altering the selectivity. When $\alpha < 2$, mixed selectivity is robust to larger noise ranges (bottom left). When $\alpha > 2$, pure selectivity is more robust (bottom right). Given normally-distributed weights, $\alpha > 2$ is more common. D. Two example cells showing how selectivity changes with changing $\lambda$. Sets of weights for both cells are drawn from the same distribution. The resulting thresholds at 3 different $\lambda$ values (labeled on the right cell but identical for each) are shown for each cell.

249 rate is 1, weights sum to 10), the cell has lost its pure selectivity and now has nonlinear
250 mixed selectivity. This happens because the T1B-T2A condition was pulled over the
251 threshold by the increase in T1B weight. In another circumstance (bottom), the cell
252 starts with nonlinear mixed selectivity. But the decrease in the weight from T1A
253 with learning pulls the T1A-T2B condition beneath the threshold, resulting in pure
254 selectivity. As the learning process continues until the weights plateau (right column),
255 the new selectivities persist.

256     The changes in selectivity with learning are the result of the representation of the
257 four conditions being expanded. Constrained learning is better able to achieve this
258 expansion. The reason for this is shown in Figure 2B. Unlike Figure 2A, this cell
259 starts off with its two strongest inputs coming from the same task variable (T2). In
260 free learning (bottom), these inputs get increased while the two from T1 get decreased.
261 This shrinks the representation along the T1 dimension and only increases it slightly
262 along the T2 direction. Thus, the selectivity of this cell (no selectivity) doesn't change.
263 With constrained learning (top), the representation is expanded in both directions (as
264 one input from each task variable is increased and the other decreased), and the cell
265 gains mixed selectivity.

266     While some cells will show changes in selectivity, changes in the representation also
267 strongly impact noise robustness. Because additive noise functions like a change in
268 threshold, it can cause a cell's response to flip. Trialwise additive noise drawn from a
269 mean-zero distribution creates a range of effective thresholds centered on the original
270 threshold value, and a cell's selectivity will only remain intact if the range of thresholds
271 that support its selectivity is larger than the noise range. Therefore, a cell's selectivity
272 is more noise robust if there is a larger range of threshold values for which its selectivity
273 doesn't change. To explore noise robustness in this model, we will define:

$$\Delta_x \equiv W_{1B} - W_{1A} \quad \Delta_y \equiv W_{2B} - W_{2A} \qquad \alpha \equiv \Delta_y/\Delta_x \geq 1 \tag{6}$$

274 Thus, $\alpha$ is the ratio of the side lengths of the rectangle formed by the four conditions
275 (see Figure 2C, top). Without loss of generality, we define the larger of the two sides
276 as associated with T2, $W_{2B} > W_{2A}$, and $W_{1B} > W_{1A}$.

277     For the cell to display pure selectivity to T2, the following inequality must hold:

$$W_{1B} + W_{2A} \leq \Theta < W_{1A} + W_{2B} \tag{7}$$

278 Therefore the range of thresholds that give rise to pure selectivity is:

$$(W_{1A} + W_{2B}) - (W_{1B} + W_{2A}) = (W_{2B} - W_{2A}) + (W_{1A} - W_{1B})$$
$$= \Delta_y - \Delta_x = \Delta_x(\alpha - 1) \tag{8}$$

279 The analogous calculations for mixed selectivity (assuming the T1B-T2B condition is
280 active only, but results are identical for T1A-T2A being the only inactive condition)
281 are:

$$W_{1A} + W_{2B} \leq \Theta < W_{1B} + W_{2B}$$
$$W_{1B} + W_{2B} - (W_{1A} + W_{2B}) = (W_{1B} - W_{1A}) = \Delta_x \tag{9}$$

Thus, pure selectivity is more noise robust than mixed selectivity when $\alpha > 2$. This imbalance can be seen in Figure 2C, where the bottom left panel shows that the range of thresholds that support mixed selectivity (red shaded area) is larger than that of pure selectivity (blue shaded area) when $\alpha < 2$. The right panel shows the reverse pattern, when $\alpha > 2$. Here, the dotted colored lines show the optimal (most noise robust) threshold for each selectivity type.

Now we show that, given weights drawn at random from a Gaussian distribution, $\alpha > 2$ is more common than $\alpha < 2$. The argument goes as follows: because $\Delta_x$ and $\Delta_y$ are differences of normally distributed variables, they are themselves normally distributed (with $\mu = 0$, $\sigma = 2\sigma_w$). The ratio of these differences is thus given by a Cauchy distribution. However, because $\alpha$ represents a ratio of lengths, we are only interested in the magnitude of this ratio, which follows a standard half-Cauchy distribution. Furthermore, $\alpha$ is defined such that the larger difference should always be in the numerator. Thus,

$$P(\alpha > 2) = 1 - \int_{1/2}^{2} \frac{2}{\pi(1+u^2)} = .5903 \tag{10}$$

Therefore, the majority of cells can be expected to have $\alpha > 2$ with random weights. This means that most cells have a representation that leads to higher noise robustness for pure selectivity than for mixed.

This comparison of noise robustness, however, assumes an optimal threshold for each type of selectivity. But selectivity (in the absence of noise) and noise robustness change as the threshold varies. Here, the threshold is defined as a fraction of the total weight going into the cell: $\Theta = \lambda \Sigma W$. As we increase $\lambda$ then, the threshold is a line with slope of -1 that moves from the bottom left corner up to the top right. Examples of this are shown in Figure 2D. With the smallest $\lambda$, neither example cell has selectivity. With the middle $\lambda$ value Cell 1 gains mixed. Cell 2 gains pure selectivity, which it retains at the higher $\lambda$, while Cell 1 switches to the other type of mixed. A low $\lambda$ is thus conducive to the type of mixed selectivity where the cell is active in all but one condition, while a high $\lambda$ can create the opposite type of mixed selectivity. Pure selectivity can come from a range of $\lambda$ in the middle.

If $\lambda$ is low, for example, a cell may still achieve pure selectivity, but it will likely do so with low noise robustness, as the threshold will be very near to the condition for mixed selectivity.

To investigate how noise robustness changes with $\lambda$, we generate a large (10000) population of cells, each with four random input weights (drawn from a Gaussian with positive mean. Qualitative results hold for many weight/variance pairs. Weights are strictly non-negative), and calculate the size of the additive noise shift needed to cause each cell to lose its selectivity (whichever it has). For each type of selectivity, we plot these noise values in the form of a cumulative distribution function: Figure 7B plots the fraction of cells that will lose their selectivity at a noise value less than or equal to that given on the x-axis. This function depends on the threshold, and so is plotted for different $\lambda$ values.

To synthesize this, we plot the noise value at which 50% of cells have lost selectivity, as a function of $\lambda$ (Figure 7C, noise values are normalized by the maximum value). On the same plot we show the percent of cells that have mixed and pure selectivity in the absence of noise. The percent of cells that ultimately demonstrate selectivity will

11

depend on the percent present without noise and the noise robustness. For example, starting at $\lambda = .25$ and going to $\lambda = .35$, the percent of cells with mixed selectivity grows, while its noise robustness decreases. So, depending on the noise level, the amount of cells with mixed selectivity may grow or shrink as $\lambda$ changes this way. This plot is used to understand the choice of threshold in the model.

Assuming a fixed threshold, we then explore how noise robustness varies with learning. In doing so, it is important to note the effect of starting from a $\lambda$ value that has unequal noise robustness for pure and mixed selectivities. Given a fixed noise value, if most cells with pure selectivity are already robust to it, an increase in noise robustness for pure will only have a moderate effect on the population levels of pure selectivity. Conversely, if most mixed cells have noise robustness less than the current noise value, an increase in that robustness could strongly impact the population. In the same vein, a decrease in robustness will impact the pure population more than the mixed.

In the case of constrained learning with $N_L = 2$, $\Delta_x$ and $\Delta_y$ both increase. According to Eqn. 7 and Eqn. 9, robustness to both selectivities increases with $\Delta_x$, which is why constrained learning causes increases in both mixed and pure selectivity (Figure 6A).

The relative increase in robustness will depend on how $\alpha$ changes. It can be shown that if $\frac{W_{1B}}{W_{1A}} < \frac{W_{2B}}{W_{2A}}$ then $\Delta_x$ will expand more than $\Delta_y$ and $\alpha$ will decrease, meaning the increase in noise robustness favors mixed selectivity. If $\frac{W_{1B}}{W_{1A}} > \frac{W_{2B}}{W_{2A}}$, then $\alpha$ will grow, and the increase in noise robustness will be larger for pure than mixed. Because the latter condition is less common, pure noise robustness doesn't increase as much as mixed (see Figure 8C, where constrained learning with $N_L = 2$ is used.)

When $N_L = 1$, only one side length will increase and the other decrease, leading ultimately to lower length of the shortest side but a larger ratio between the sides (so more robustness to noise for pure selectivity and less for mixed). This is straightforward for $W_{2B} > W_{1B}$ ($\Delta_y$ grows and $\Delta_x$ shrinks) and contributes to the increase in pure selectivity with $N_L = 1$ in Figure 6A. However, if $W_{1B} > W_{2B}$, $\alpha$ will first decrease as $\Delta_x$ grows and $\Delta_y$ shrinks. This is good for mixed noise robustness. The ratio then flips ($\Delta_x > \Delta_y$), and $\Delta_y$ (the side that is now shorter) is still shrinking and $\Delta_x$ is growing. In this circumstance, if $\Delta_y/\Delta_x$ becomes less than $\frac{1}{2}$, the representation will favor pure noise robustness over mixed. This pattern is reflected in the shape of the mixed selectivity changes seen with $N_L = 1$ in Figure 6A (mixed selectivity increases then decreases). This flipping of $\alpha$ is possible for some cells when $N_L = 2$ if $\frac{W_{1B}}{W_{1A}} < \frac{W_{2B}}{W_{2A}}$, but the weights would likely plateau before $\alpha$ became less than $\frac{1}{2}$, and so the drop in mixed selectivity does not occur.

In free learning with $N_L = 2$, cells that have $W_{1A} > W_{2B}$, will see both weights from T1 increase and (due to the weight normalization) both weights from T2 decrease. Because the weights change in proportion to their value, $\Delta_x$ increases, $\Delta_y$ decreases and so $\alpha$ goes down. This leads to more noise robustness for mixed and less for pure. If $W_{2A} > W_{1B}$, these trends are reversed and the cell has more noise robustness for pure and less for mixed.

## 3. Results

In this study, we analyzed various measures of selectivity of a population of PFC cells recorded as an animal carried out a complex delayed match-to-sample task.
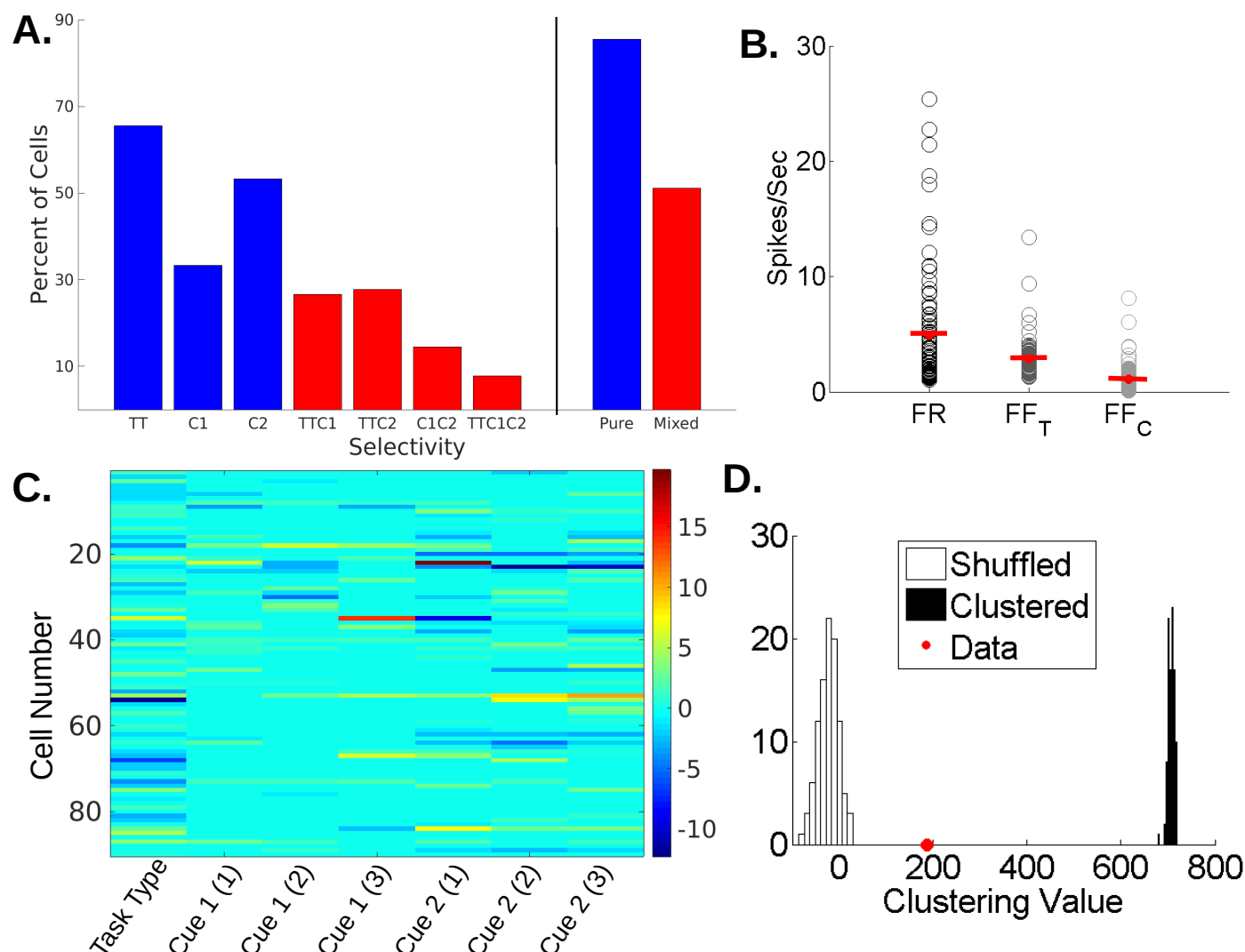
Figure 3: Results from the experimental data. A.) Selectivity profile of the 90 cells analyzed. A cell had pure selectivity to a given task variable if the term in the ANOVA associated with that task variable was significant (p<.05). A cell had nonlinear mixed selectivity to a combination of task variables if the interaction term for that combination was significant. On the right of the vertical bar are the percent of cells that had at least one type of pure selectivity (blue) and percent of cells that had at least one type of mixed selectivity (red). B.) Values of firing rate, $FF_T$, and $FF_C$ for this data. Each open circle is a neuron and the red markers are the population means. C.) Beta coefficients from GLM fits for each cell. The first regressor corresponds to task type, regressors 2-4 correspond to cue 1 and 4-7 to cue 2. These values were used to determine the clustering value D.) Histograms of clustering values generated for different distributions. The shuffled data comes from shuffling the selectivity coefficients across cells. The clustered data is designed to have 3 different categories of cell types defined according to selectivity. The red dot shows the data value.

13

Through this process, several properties of the representation in PFC were discovered and a simple circuit model that included Hebbian learning was able to replicate them. These properties, combined with the modeling results, provide strong support for the notion that PFC selectivities are the result of Hebbian learning in a random network.

### 3.1. PFC Population is Moderately Specialized and Selective

The average firing rate of of cells in this population was 4.90+/-5.14 spikes/s. Fano Factor analyses provided measurements of the noise and density of response in the data (Figure 3B). The average value of the across-trial Fano Factor ($FF_T = 2.86 +/-1.68$), shows that the data has elevated levels of noise compared to a Poisson assumption. Looking at $FF_C$—a measure of how a cell's response is distributed across conditions—suggests that PFC cells are responding densely across the 24 conditions ($FF_C = 1.11 +/-1.19$, for comparison, at the observed average firing rates, a cell that responded only to a single condition would have $FF_C \approx 120$, one that responded to two conditions would have $FF_C \approx 57$). This finding suggests that these cells are not responding sparsely and are not very specialized for the individual conditions of this task.

Each condition is defined by a unique combination of 3 task variables: task type, identity of image cue 1 and identity of image cue 2 (Figure 1A). Selectivity to task variables was determined via a 3-way ANOVA. The results of this analysis are shown in Figure 3A. This figure shows the percentage of cells with selectivity to each task variable and combination of task variables (as determined by a significant (p<.05) term in the ANOVA). A cell that has selectivity to any of the regular task variables (task type, cue 1, cue 2) has pure selectivity, while a cell that has selectivity to any of the interaction terms (combination of task variables such as task type-cue1, task type-cue 2, etc) has nonlinear mixed selectivity. The final two bars in Figure 3A show the number of cells with pure and mixed selectivity defined this way. Note that a cell can have both pure and mixed selectivity, thus the two values sum to more than 100%.

The majority of cells (77/90) showed pure selectivity to at least one task variable. But the population shows clear biases in the distribution of these pure selectivities: task type selectivity is the most common (59 cells) and cue 2 is represented more than cue 1 (48 vs. 30 cells) (these biases are observable in the GLM fits as well, see Figure 3C). This latter effect may be due to the time at which these rates were collected: these rates were taken during the second delay, which comes directly after the presentation of the second cue. The former effect is perhaps more surprising. While the task type is changed in blocks and thus knowable to the animal on each trial (with the exclusion of block changes), there is no explicit need for the animal to store this information: the presence of a second sequence or an array of images will signal the task type without the need for prior knowledge. However, regardless of its functional role in this task, contextual encoding is a common occurrence ([10, 19]). Furthermore, the fact that the recall task is more challenging than the recognition task may contribute to clear representation of task type. That is, it is possible that the animals keep track of the task type in order to know how much effort to exert during the task.

Approximately half of the cells (46) had some form of mixed selectivity, mostly to combinations of two task variables. The small number of cells with selectivity to the 3-way interaction term (TT-C1-C2) is consistent with the relatively low value of $FF_C$ in this population, as a strong preference for an individual condition would lead to a

14

high $FF_C$. The number of cells with only mixed selectivity was low (only 1 out of 90 cells), 32 cells had only pure selectivity, and 12 cells had no selectivity.

We use a population-level analysis inspired by [35] to measure the extent to which cell types are clustered into categories. Here, we used this analysis to determine if cells cluster according to their responsiveness to different task variable identities (i.e., recognition vs recall). That is, are there groups of neurons which all prefer the same task type and image identities, beyond what would be expected by chance? In order to explore this, we first use a GLM, with task variable identities as regressors, to fit each neuron individually. The beta coefficients from these fits define a neuron's position in selectivity space (these beta coefficient values are shown in Figure 3C, and a schematic of how the clustering measure works is shown in Figure 1D). The clustering measure then determines the extent to which the population of neurons deviates from a uniform distribution in this space. The data had a clustering value of 186.22. Comparing this to the mean values of two distributions of artificially generated populations suggests the data has a mild but significant deviation from random: the average clustering value for populations generated by randomly shuffling the coefficient values is -22.59+/-21.75, and the average value of populations that have 3 distinct clusters of selectivity is 706.68+/-6.84. As the data clustering value sits in between these values and closer to the shuffled data, we conclude that some structure does exist in the data, yet the cells in this population do not appear to form strongly separable categories as defined by task variable identity preference (Figure 3D).

### 3.2. Circuit Model without Hebbian Learning Cannot Replicate Mix of Density and Specialization

A simple circuit model was made to replicate the selectivity properties found in the data. The model contains two layers: an input layer consisting of binary neurons that represent task variable identities and an output layer consisting of "PFC" neurons which get randomly-weighted input from the first layer and whose activity is a nonlinear function of the sum of that input. The model also has two forms of noise: an additive term applied before the nonlinearity (which replicates input/background noise, and implicitly shifts the threshold of the cell), and a multiplicative term applied after (which enforces the observed relationship between firing rate and variance) (see Methods and Figure 4A).

The output of the initial circuit model, prior to any Hebbian learning, was analyzed in the same way as the data to determine if it matched the properties found in PFC. The results of this can be found in Figure 5. First, in Figure 5A, we demonstrate the impact of the noise parameters on $FF_T$, pure and mixed selectivity, and the clustering value. As expected, increasing the additive and/or multiplicative noise terms increases the $FF_T$, as this is a measure of trial variability. Increasing noise also makes it harder for cells to reach significance, and thus the percentage of cells with pure and mixed selectivity are inversely related to the noise parameters, (the relative sensitivities of mixed and pure selectivity to noise will be discussed in depth later). For similar reasons, clustering value also decreases with noise (cells need to display significant preferences to task variable identities in order to form clusters based on that).

To determine the impact other properties of the model had on our measures of interest, we varied several other parameters. Figure 5B shows what happens at different values of the threshold parameter. Here, the threshold is given as the amount of input the cell needs to reach half its maximal activity, expressed as a fraction of its
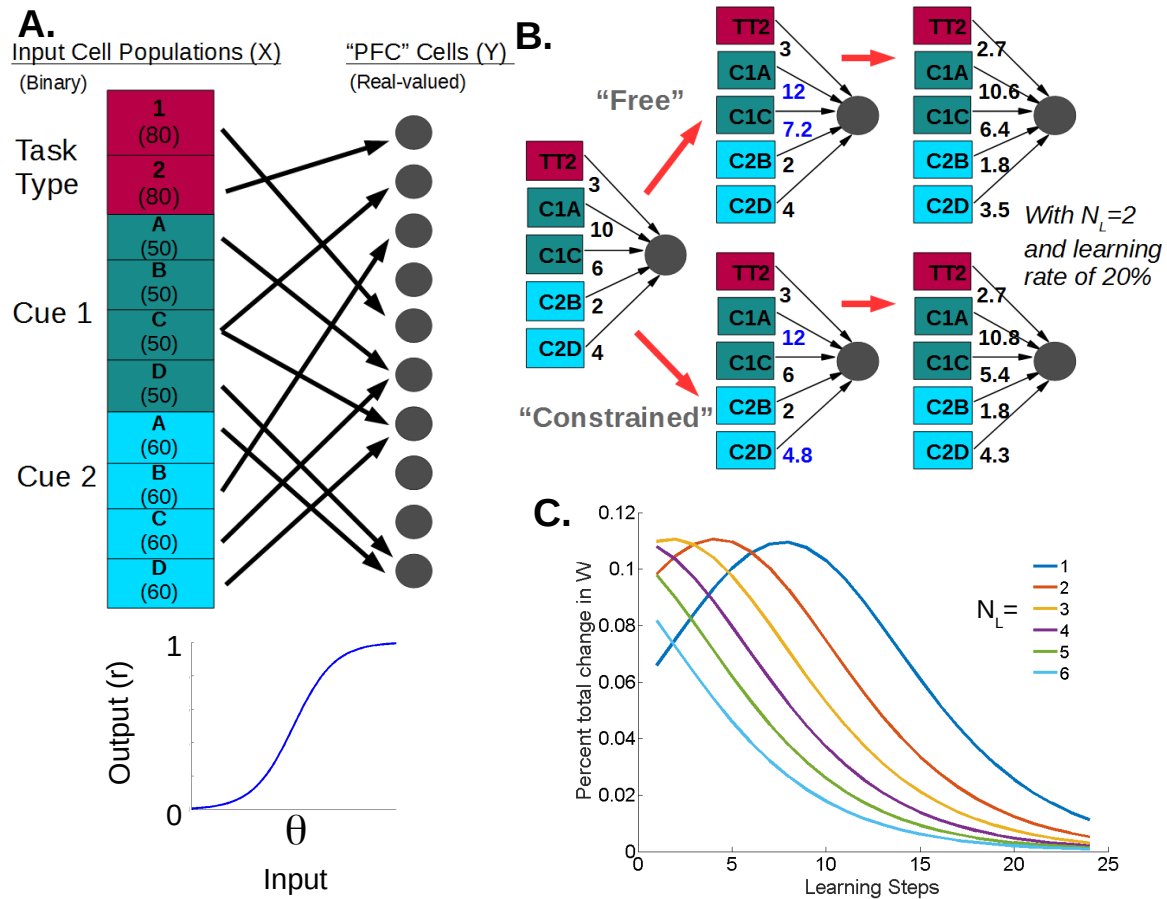
15

Figure 4: The full model and how learning occurs in it. A.) The model consists of groups of binary input neurons (colored blocks) that each represent a task variable identity. The number of neurons per group is given in parenthesis. Each PFC cell (gray circles) receives random input from the binary cells. Connection probability is 25% and weights are Gaussian-distributed and non-negative. The sum of inputs from the binary population and an additive noise term are combined as input to a sigmoidal function (bottom). The output of the PFC cell on a given trial is a function of the output of the sigmoidal function, $r$ and a multiplicative noise term (see Methods). The threshold, $\Theta$, is given as percentage of total input to each cell B.) Two styles of learning in the network, both of which are based on the idea that the input groups that initially give strong input to a PFC cell have their weights increased with learning (sum of weights from each population are given next to each block). In free learning, the top $N_L$ input populations are chosen freely. In this example, that means two groups from the cue 1 task variable have their weights increased (marked in blue). In constrained learning, the top $N_L$ populations are chosen with the constraint that they cannot come from the same task variable. In this case, that means that cue 2D is chosen over cue 1C despite the latter having a larger summed weight. In both cases, all weights are then normalized. C.) Learning curves as a function of learning steps for different values of $N_L$. Strength of changes in the weight matrix expressed as a percent of the sum total of the weight matrix are plotted for each learning step (a learning step consists of both the weight increase and normalization steps). Different colors represent different $N_L$s.

total input (keep in mind that, given the number of input cells in each population and the task structure, roughly one-third of input cells are on per trial). The colored lines are, for each measure, the extent to which the model differs from the data, expressed in units of the model's standard deviation (calculated over 100 instantiations of the model). Due to the impact of noise parameters discussed above, at each point in this graph the noise parameters were fit to ensure the model was within $+/-$ 1.5 standard deviations of the data $FF_T$ (this generally meant that it varied from $\sim$ 2.8 to 2.9).

With an increasing threshold, the $FF_C$ (green line in Figure 5B) increases. This is because higher thresholds mean cells respond to only a few combinations of input, rather than responding similarly to many, and the $FF_C$ is a measure of variability in response across conditions (note that while $FF_C$ appears to peak at $\approx$ .35 and decrease, this particular trend is driven by an increase in $FF_C$ standard deviation; the mean continues to increase). The percentage of cells with mixed selectivity (red line) also increases with threshold. With a higher threshold, the majority of conditions give input to the cell that lies in the lower portion of the sigmoidal function (bottom of Figure 4A). The nonlinearity is strong here—with some input producing little to no response—thus, more cells can attain nonlinear mixed selectivity. Pure selectivity also increases with threshold, and the percent of cells with pure selectivity goes quickly to 100 (and the standard deviation of the model gets increasingly small). We go into more detail about the reliance of selectivity on threshold later.

The clustering value relies on cells having preference for task variable identities and so increases as selectivity increases initially. However, just having selectivity is not enough to form clusters, and so the clustering value in the model levels off below the data value even as the number of cells with pure selectivity reaches full capacity. Thus, with the exception of the clustering value, the model can reach the values found in the data by using different thresholds. As Figure 5B shows, however, at no value of the threshold are all measures of PFC response in the model simultaneously aligned with those in the data.

Figure 5C shows how the same measures change when the width of the weight distribution from input to PFC cells is varied. Here, the standard deviation of the distribution from which connection strengths are drawn ($\sigma_W$) is given as a factor of the mean weight, $\mu_W$. Increasing this value increases pure and mixed selectivity as well as $FF_C$. Because a wider weight distribution increases the chances of a very strong weight existing from an input cell to an output cell, it makes it easier for selectivity to emerge (that is, the output cell's response will be strongly impacted by the task variable identity the input cell represents). The $FF_C$ increase occurs for similar reasons: a cell may have uneven responses across conditions due to strong inputs from single input cells. Clustering values, however, are unaffected by this parameter. At no point, then, can the model recreate all aspects of the data by varying the weight distribution. Furthermore, while values of mixed selectivity and $FF_C$ approach the data values with large $\sigma_W/\mu_W$, such large values are likely unrealistic. Data show that a $\sigma_W/\mu_W$ ratio of around 1 is consistent with observations of synaptic strengths from several brain areas [3].

Varying other parameters such as the mean weight, number of cells per population, and connection probability similarly doesn't allow the model to capture all properties of the data (not shown).

Figure 5D shows the values of the model as compared to the data for the set of parameters marked with arrows in Figure 5B and 5C. For reasons that will be discussed
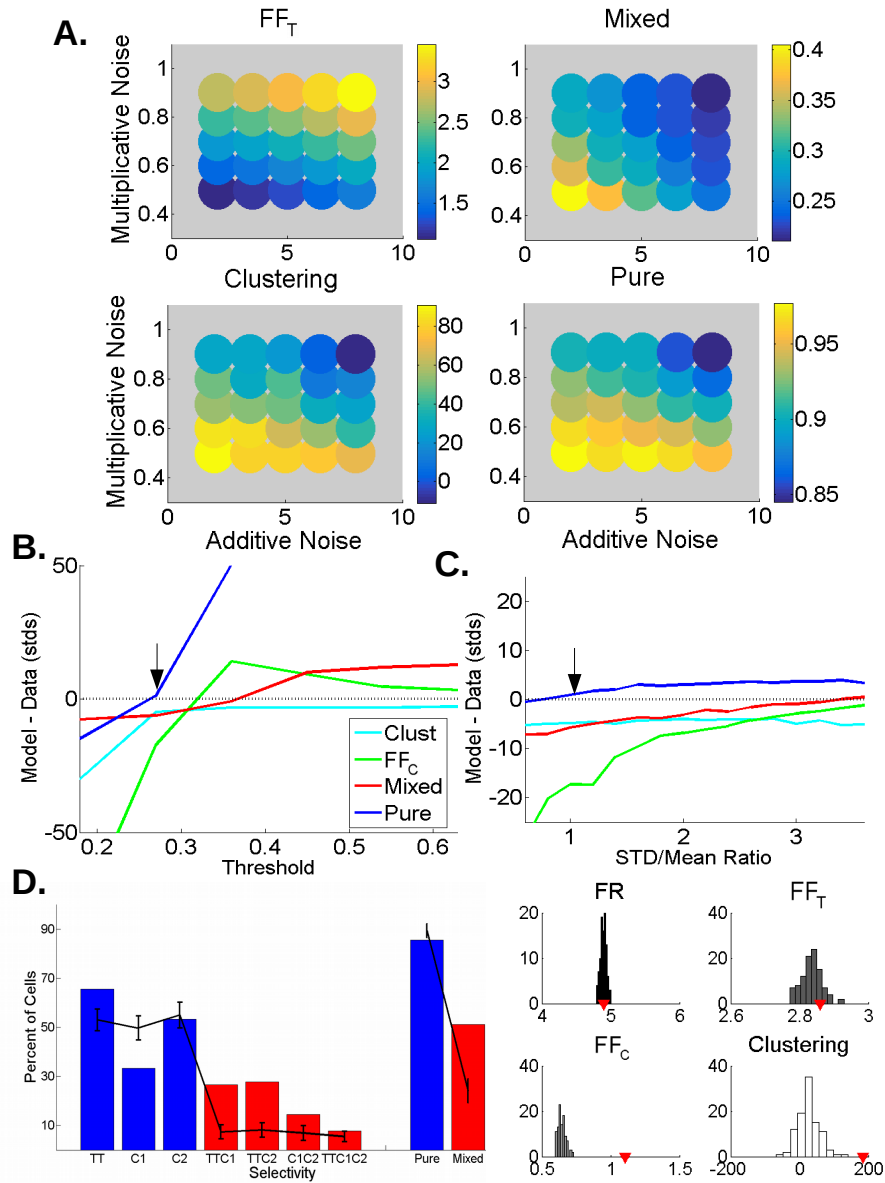
17

Figure 5: Results from the model without learning. A.) $FF_T$ and other measures can be controlled by the additive and multiplicative noise parameters. Each circle's color shows the value for the given measure averaged over 25 networks, for a set of $a$ and $m$ values (see Methods). $FF_T$ scales predictably with both noise parameters. Mixed selectivity, pure selectivity, and clustering scale inversely with the noise parameters. Other model parameters are taken from the arrow locations in (B) and (C). B.) How the threshold parameter, $\lambda$, affects measures of selectivity. Lines show how the average value of the given measure in the model (in units of standard deviations away from the data value) varies as a function of the threshold parameter $\lambda$, where $\Theta_i = \lambda \Sigma_j w_{ij}$ At each point noise parameters are fit to keep $FF_T$ close to the data value. C.) Same as (B), but varying the width of the weight distribution rather than the threshold parameter. D.) Example of the model results at the points given by the black arrows in (B) and (C). On the left, blue and red bars are the data values as in Fig 2. The lines are model values (averaged over 100 networks, errorbars +/-1 std). On the right, histograms of model values over 100 networks. The red markers are data values. This model has no learning.

18

more later, these parameters were chosen because they were capable of capturing the amount of pure selectivity in the model (any higher value of the threshold would lead to too many cells with pure selectivity, for example). On the left are the percentage of cells with different selectivities as in Figure 3C. The bars are the data and the lines are the model. On the right, are histograms of model values from 100 instantiations, with the red markers showing the data values. The model matches the average firing rate and $FF_T$ of the model, as it was fit to do so. Clustering, $FF_C$, and the amount of mixed selectivity are too low in the model. We use these parameters as the starting point for learning in this model.

### 3.3. Circuit Model with Hebbian Learning Captures PFC Responses

As described above, responses of PFC cells have a set of qualities that cannot be explained by random connectivity. In particular, the inability of the random network to simultaneously capture the values of $FF_C$, clustering, pure, and mixed selectivity shows that PFC cells have a balance of specialization that may require learning to achieve. Here, we tested two variants of Hebbian learning to determine if a network endowed with synaptic plasticity can capture the elements of the data that the random network could not. The simple form of Hebbian learning that we use is based on the idea that the input populations that randomly start out giving strong inputs to a cell would likely make that cell fire and thus have their weights increased. In both variants of learning tested, each cell has the weights from a subset ($N_L$) of its input populations increased while the rest are decreased to keep overall input constant (this is done via a weight increase step and a normalization step). Mechanisms for such balancing of Hebbian and homeostatic plasticity have been observed experimentally ([17]), particularly via the type of synaptic up and down regulation used here ([5, 40, 21]).

The difference between the two variants of learning comes from which input populations are increased. In general, the top $N_L$ input populations from which the cell already receives the most input have their weights increased (to capture the "rich get richer" nature of Hebbian learning). In the "constrained" variant, however, weight increases onto a PFC cell are restricted to populations of input cells that come from different task variables (e.g., cue 1 and cue 2. For a detailed explanation see Methods). This was done to ensure that cells had enough variety of inputs to create mixed selectivity. In the free variant, the populations from which a cell receives increased input due to learning are unrestricted. That is, they are determined only by the amount of input that the cell originally received from each population as a result of the random connectivity. This unrestricted form of learning is more biologically plausible as it can be implemented locally, without knowledge of other inputs. A toy example of each variant can be found in Figure 4B. Given random weights, free and constrained learning will select the same input populations in some cells.

Figure 4C shows how the weight matrix changes with different $N_L$ values (the number of populations from which weights are increased during learning). The higher the $N_L$ the faster the matrix converges to its final state. When $N_L$ is low, convergence takes longer as all the weight is transferred to a small number of cells. This plot is shown with a learning rate of .2.

The results of both forms of learning are shown in Figure 6A. The effects of learning are dependent on $N_L$, and different $N_L$ values are in different colors ($N_L = 1, 2, 3$ are tested here). Free learning is shown with solid lines, and constrained with dotted lines, except for the case of $N_L = 1$, where free and constrained learning do not differ and
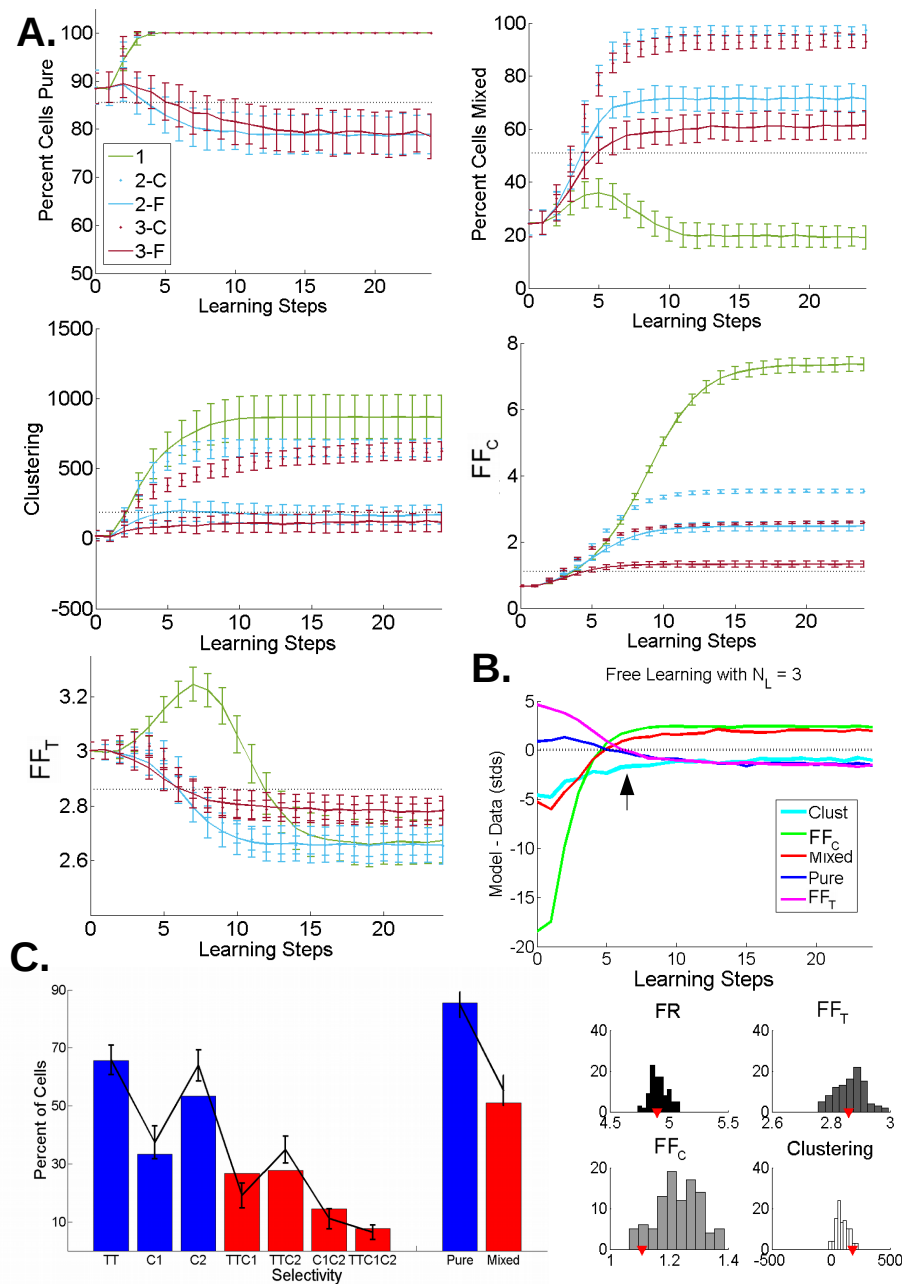
Figure 6: The model with learning. A.) How selectivity measures change with learning. In each plot, color represents $N_L$ value, solid lines are free learning, and dotted lines are constrained learning (only one line is shown for $N_L = 1$ as the free and constrained learning collapse to the same model in this circumstance). Step 0 is the random network. Black dotted lines are data values and errorbars are $+/-1$ std over 100 networks. In the pure selectivity plot, with constrained learning and when $N_L = 1$, the value maxes out at 100% in essentially all networks, leading to vanishing errorbars. B.) All measures as a function of learning for the $N_L = 3$ free learning case. Values are given in units of model standard deviation away from the data value as in Figure 5B and C. C.) The model results at the learning step indicated with the black arrow in (B), same as in Figure 54D. Here, the model provides a much better match to the data.

561 only one line is shown. In each plot, the data value is shown as a small black dotted
562 line.

563     Clustering, mixed selectivity, and $FF_C$ all increase with learning, for any value of
564 $N_L$ and both learning variants. When $N_L = 1$ (green line), mixed selectivity peaks and
565 then plateaus at a lower value (as connections to all but one population are pruned),
566 while other values of $N_L$ plateau at their highest values. As it was designed to do so,
567 constrained learning is very effective at increasing mixed selectivity, eventually getting
568 to nearly 100 percent of cells. Free learning produces more modest increases in mixed
569 selectivity, with $N_L = 2$ leading to slightly larger increases than $N_L = 3$.

570     A factor impacting selectivity in this model—and especially with this task structure—
571 is that cells that receive inputs from multiple populations from a single task variable
572 may not end up having significant selectivity to that variable. This is especially true
573 for the 'task type' variable, as cells can easily end up with input from both 'recall' and
574 'recognition' populations. If the inputs from these populations are somewhat similar in
575 strength, the cell does not respond preferentially to either. This can help understand
576 the discrepancy in how pure selectivity changes with free and constrained learning.
577 In constrained learning, pure selectivity necessarily increases with learning (to the
578 point where nearly all networks have 100% pure selectivity), whereas free learning can
579 have inputs that effectively cancel each other out. A more direct investigation of how
580 selectivity changes with learning occurs in the next section.

581     In these plots, both noise parameters are fixed, which allows us to see how $FF_T$
582 varies with learning (this is also why the values at step 0 in Figure 6A do not always
583 match those shown in Figure 5, as that model has noise parameters fit to match the
584 data). The changes in $FF_T$ stem from both changes in robustness to the additive noise
585 and from changes in the mean responses, which impacts $FF_T$ via the multiplicative
586 noise term. Figure 6A shows that the variant of learning has less of an impact on $FF_T$
587 than $N_L$ does. In all cases, however, learning ultimately leads to lower trial variability
588 in the model. This is consistent with observation made in PFC during training [34].

589     Overall, low $N_L$ leads to more acutely distributed weights and stronger structure
590 and selectivity in the model. Constrained learning, with its guarantee of enhancing
591 weights from different task variables, is also more efficient at enhancing structure
592 and selectivity. The prefrontal cortex data shows a moderate level of structure and
593 selectivity, therefore the approach that is best able to capture it is free learning with
594 $N_L = 3$. In Figure 6B, we show how all of the model values compare to the data as
595 this form of learning progresses. These plots, similar to Figure 5B and C, show values
596 in units of standard deviations away from the model. It is clear from these plots that
597 this form of learning leads all values in the model closer to those of the data, and all
598 values eventually plateau within +/- 2.5 model standard deviations of the data. The
599 best fit to the data comes after 6 learning steps with a learning rate of .2 (marked
600 with a black arrow). At this point the ratio of the standard deviation to the mean of
601 the distribution has only slightly increased, remaining within a biologically plausible
602 range. We plot the values of the data in comparison to model in Figure 6C, similarly
603 to Figure 5D. At this point, the average percent of cells with only pure selectivity is
604 25.40+/-4.16, with only mixed 4.42+/-2.15, and with no selectivity 15.9+/-4.08 (the
605 comparable data values are ≈ 36%, 1%, and 13%, respectively). Thus, the model with
606 learning is a much better fit to the data than the purely random network.

21

*3.4. Understanding Properties of Selectivity Before Learning*

We have shown that Hebbian learning can impact selectivity properties in a model
of PFC. Some of these impacts, particularly the increase in mixed selectivity, may seem
counterintuitive. Here we use a further simplified toy neuron model to understand the
properties of the network before learning and then demonstrate how learning causes
these changes.

A schematic of this toy model is in Figure 7A and 8A, and it is fully described in
the Methods. Briefly, the cell gets four total inputs–two (A and B) from each of two
task variables (T1 and T2). The output of the cell is binary: if the weighted sum of
the inputs is above the threshold, $\Theta$, the cell is active and otherwise it is not. As in
the full model, $\Theta$ is defined as a fraction, $\lambda$, of the sum of the input weights.

This format makes it easy to spot nonlinear mixed selectivity: if the cell is active
(or inactive) for exactly one of the four conditions, it has nonlinear mixed selectivity
to the combination of T1-T2. If the cell's output can be determined by the identity of
only one task variable, it has pure selectivity (and would be active for two of the four
conditions). Otherwise it has no selectivity (active or inactive for all conditions) (see
examples in Figure 2A and B).

Learning impacts selectivity by altering the way a cell represents these four condi-
tions. To say more about how this occurs, we must first describe the properties of the
representation in the random network before learning.

To be robust to noise, the cell's response should be constant across conditions.
Additive noise can be thought of as a shift in the threshold, which may lead to a
change in the cell's response. Thus, trialwise additive noise drawn from a distribution
centered on zero can be thought of as a range of effective thresholds centered on the
original one (gray shaded area in Figure 8A, black dotted line is the threshold without
noise). If the inputs for a given condition fall in this range, the response of the cell
will be noisy, i.e. flipping from trial to trial, and selectivity will be lost. Robustness to
noise, then, can be measured as the range of thresholds a representation can sustain
without any responses flipped, with a larger range implying higher noise robustness.

Assuming optimal threshold values for each, the relative noise robustness of mixed
and pure selectivity can be calculated (see Methods). We find that, thinking of the
four conditions as the corners of a rectangle (as visualized in Figure 2C), mixed se-
lectivity robustness depends on the length of the shorter side, while pure selectivity
noise robustness depends on the difference between the two side lengths. We also find
that, with random weights, most cells will have a representation that has higher noise
robustness for pure selectivity than for mixed (see Methods).

Noise robustness changes, however, as thresholds deviate from optimal. The type
of selectivity cells have in the absence of noise also varies with threshold (see Figure
2D for examples). To quantify these trends, we varied the threshold parameter $\lambda$ and
determined both the probability of different types of selectivity as well as the noise
robustness for each type (see Methods for details). In Figure 7B, we show the fraction
of cells that lose selectivity at a given noise level, for three different values of $\lambda$. Noise
robustness (plotted as a function of $\lambda$ in Figure 7C) is defined then as a normalized
measure of the noise value that causes 50% of cells to lose selectivity.

Figure 7C demonstrates why the random network from which we start learning is
necessarily in a condition of low mixed selectivity. The value of $\lambda$ we choose to start
from is constrained by the fact that the data shows high levels of pure selectivity.
Therefore, we need a value that has high probability of pure selectivity and high
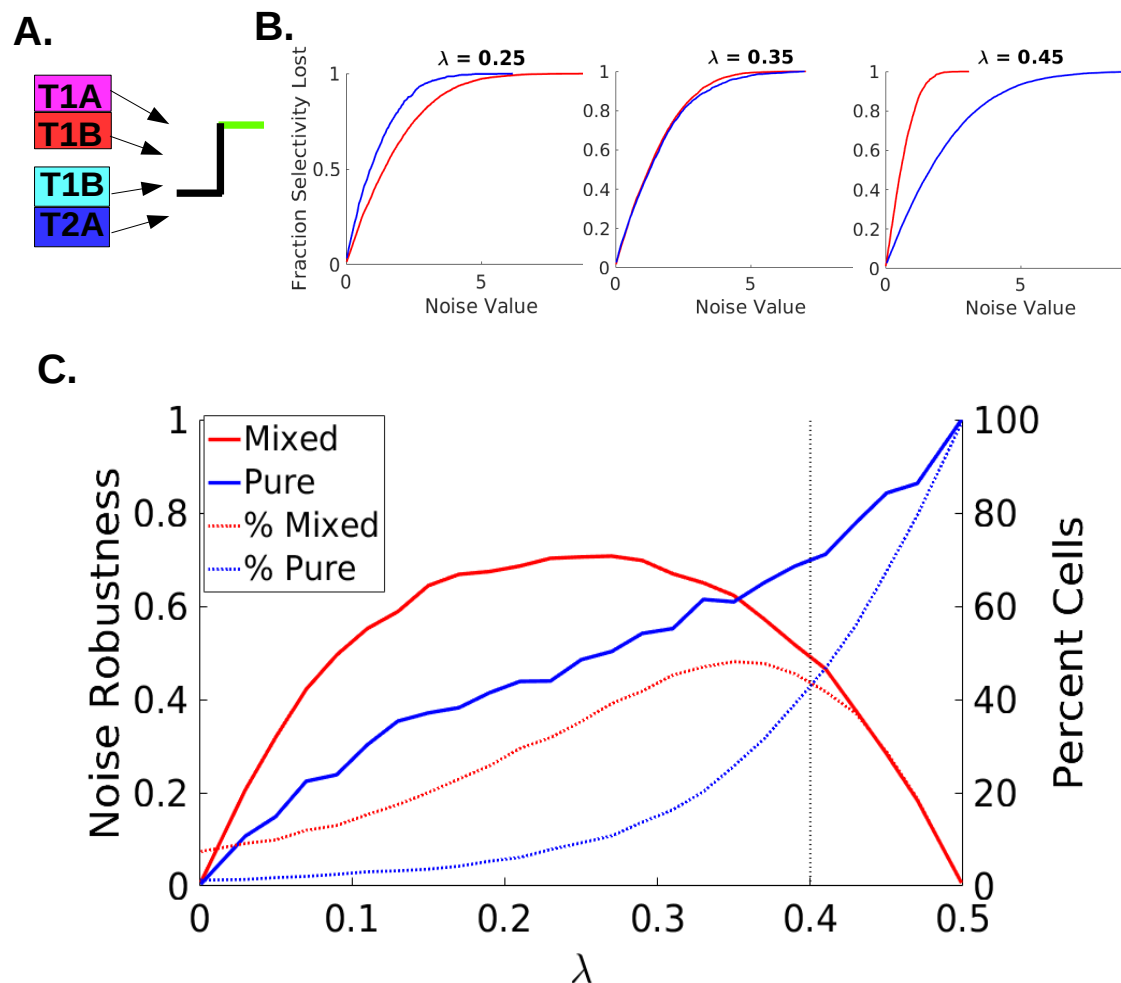
22

Figure 7: How noise robustness varies with threshold in a random network using the toy model A.) Schematic of the toy model: four input populations (two from each task variable) send weighted inputs to a cell with a threshold ($\Theta$) nonlinearity B.) For a given noise value, the fraction of cells that would lose selectivity if that noise value were used. Values are separated for cells with pure (blue) and mixed (red) selectivity. Three $\lambda$ values shown, where $\Theta = \lambda \Sigma W$. C.) Based on plots like those in (B), the noise value at which 50% of cells have lost selectivity is calculated ("Noise Robustness" refers to these values normalized by the peak value. Higher values are better) and plotted as a function of $\lambda$ (solid lines). On the same plot, the percent of cells with each type of selectivity in the absence of noise is shown (dotted lines). The black doted line marks a $\lambda$ value at which the probability of mixed and pure selective cells is equal, but their noise robustness is unequal. This plot is mirror-symmetric around $\lambda = .5$

noise robustness for it. Values of $\lambda$ that meet this condition are not favorable for mixed selectivity. Therefore, the best we can do is choose a value of, for example, .4, where probabilities of pure and mixed are even, but pure has higher noise robustness (therefore effective rates of pure selectivity are higher). The fact that mixed selectivity is less noise robust than pure in the full model can be seen in Figure 5A.

Note that while the $\lambda$ used for the random version of the full model shown in Figure 5D was around .27, that value is not directly comparable to the $\lambda$ values in these plots for many reasons. First, the full model has 3 task variables, compared to the 2 used in the toy model. This means that, from the perspective of mixed selectivity for 2 task variables, a given $\lambda$ value will create a higher $\Theta$ in the full model with 3 task variables than in the toy one that has only 2 (because $\Theta$ is a function of the sum total

23

of all weights, not just those relevant for the 2-way selectivity). In addition, in the toy model, 50% of the inputs are on for any given condition, whereas the nature of the task in the full model means that only 25% of inputs are on when looking at C1-C2 mixed selectivity, while one-third are on for TT-C1, TT-C2, and TT-C1-C2 mixed selectivity. The percentage of cells are also not directly comparable, as cells in the full model are labeled as pure if they have any of 3 different types of pure selectivity, and mixed if they have any of 4 different types of mixed. This toy model is thus meant to provide intuition only.

### 3.5. How Learning Impacts Selectivity

For the reasons just discussed, the random model starts in a regime where pure selectivity has high noise robustness and mixed does not. In order to match the amount of mixed selectivity seen in the data, we must then rely on learning to increase noise robustness for mixed selectivity, allowing more mixed cells to reach significance.

Learning impacts noise robustness by expanding the representation of the different conditions. An example of this is in Figure 8A, where the gray shaded area represents the noise-induced range of the threshold. Before learning, the cell's response is impacted by the noise. With learning, different conditions get pulled away from each other and the threshold, creating a much more favorable condition for mixed selectivity to be robust to noise. As can be seen, the responses are now outside the noise range.

For the same reason that learning increases noise robustness (because the expansion increases the range of thresholds that support mixed selectivity), it can also increase the probability of a cell having mixed selectivity in the absence of noise. This can be seen in Figure 8C (left), where learning steps are indicated by increasing color brightness (constrained learning with rate of .25). At lower $\lambda$ values, cells that are initially above threshold for all conditions (no selectivity) gain mixed selectivity with learning. But for $\lambda$ values that support higher levels of pure selectivity (e.g., $\lambda = .4$, marked with a black dotted line), the percent of cells with mixed is not as impacted by learning. The percent of cells with pure selectivity increases only slightly at most $\lambda$ values.

Noise robustness has a different pattern of changes with learning (Figure 8C, right). In particular, at $\lambda = .4$, the noise robustness still increases with learning even when the percent of cells with mixed doesn't change. Thus, changes in noise robustness are more relevant for the increase in mixed selectivity observed in the full model.

In particular, constrained learning with $N_L = 2$ always increases the lengths of both sides of the rectangle (as one weight from each task variable increases and the other decreases). As mentioned above, noise robustness for mixed selectivity scales with the length of the shorter side and so it necessarily increases with learning in this condition. Under certain weight conditions, noise robustness will also increase for cells with pure selectivity (this can be seen in Figure 8C, see Methods for details).

If $N_L = 1$, only one side length will increase and the other decrease. If the shorter side decreases, mixed selectivity noise robustness decreases. If the shorter side increases, mixed noise robustness increases, up until the point at which side lengths are equal. At that point the shorter side is now the decreasing side and mixed noise robustness goes down. This trend is reflected in the shape of the mixed selectivity changes seen with $N_L = 1$ in Figure 6A (mixed selectivity increases then decreases).

When using free learning (with $N_L = 2$), a portion of the cells will by chance have the same changes as with constrained learning. The remaining cells cause the

24

differences observed between the two versions of learning, and can be of two types. In the first type, the larger side length increases and the smaller shrinks, causing a decrease in mixed noise robustness. Free learning doesn't achieve the same levels of mixed selectivity as constrained because these cells continue to be too noisy. In the other type, the shorter side increases and the larger decreases, reducing the difference between the two side lengths and thus reducing pure noise robustness. Free learning loses pure selectivity as these cells become too noisy (as seen in 6A). More detailed descriptions of changes with learning can be found in the Methods.

Inputs from additional task variables can be thought of as a source of noise as well. In Figure 8B, we add a third task variable to the toy model. Now, in the case of the T1B-T2A condition, the identity of T3 determines if the cell is active or not. From the perspective of T1-T2 mixed selectivity, this has the same impact as shifting the threshold, and thus creates noise. If both T3 inputs are weaker than the strongest two inputs from T1 and T2 (as they are here), they will decrease with learning. This means that not only do different T1-T2 conditions get pulled apart with learning, but the same T1-T2 conditions become closer. This reduces the impact of "noise" from other task variables, and explains why mixed increases more with $N_L = 2$ than with $N_L = 3$ (Figure 6A).

In sum, learning changes a cell's representation of the task conditions. Depending on the threshold value, this can create changes in the probability of mixed and pure selectivity and the relative noise robustness for each. Here, in order to match the high levels of pure selectivity seen in the data, we use a threshold regime where mixed selectivity noise robustness increases with learning. This causes a gain in the number of cells with mixed selectivity, such that it reaches the level seen in the data.

### 3.6. How Learning Impacts Other Properties

The visualization of this toy model gives intuition for why other properties change with learning as well. $FF_C$, for example, increases with learning (Figure 6A). The expansion that comes with learning places different conditions at different distances from the threshold. With a sigmoidal nonlinearity, this would translate to more variance in the responses across conditions, increasing $FF_C$. Because constrained learning ensures the most expansion, it increases $FF_C$ more. These increases depend on $N_L$ because lower $N_L$ allows for a more extreme skewing of weights, and thus a subset of conditions will be far above threshold while the rest are below (leading to a high $FF_C$). $FF_C$ has a limit, however, because even with $N_L = 1$, the cell would still respond equally to a quarter of the conditions (assuming an input from a cue variable)

Clustering values are also impacted by how selectivity changes. Clustering in the data appears to be driven by task type selectivity (Figure 3C), and as task type preferences develop in the model the clustering value increases. Here, the relative sizes of the the input populations play a role. Because the input populations that represent task type contain more cells (Figure 4A), these populations are more likely to be among the strongest inputs to a cell, and thus have their weights increased (Note that this bias in favor of task type could also arise from the fact that only two task types are possible, and thus these inputs are on twice as often as cue inputs. Such a mechanism cannot be implemented in this model, however, so we use uneven numbers of input cells). Therefore, task type selectivity becomes common and clusters form around the axis representing the first regressor (which captures task type preference). This effect is weaker with free learning because both task type populations may have
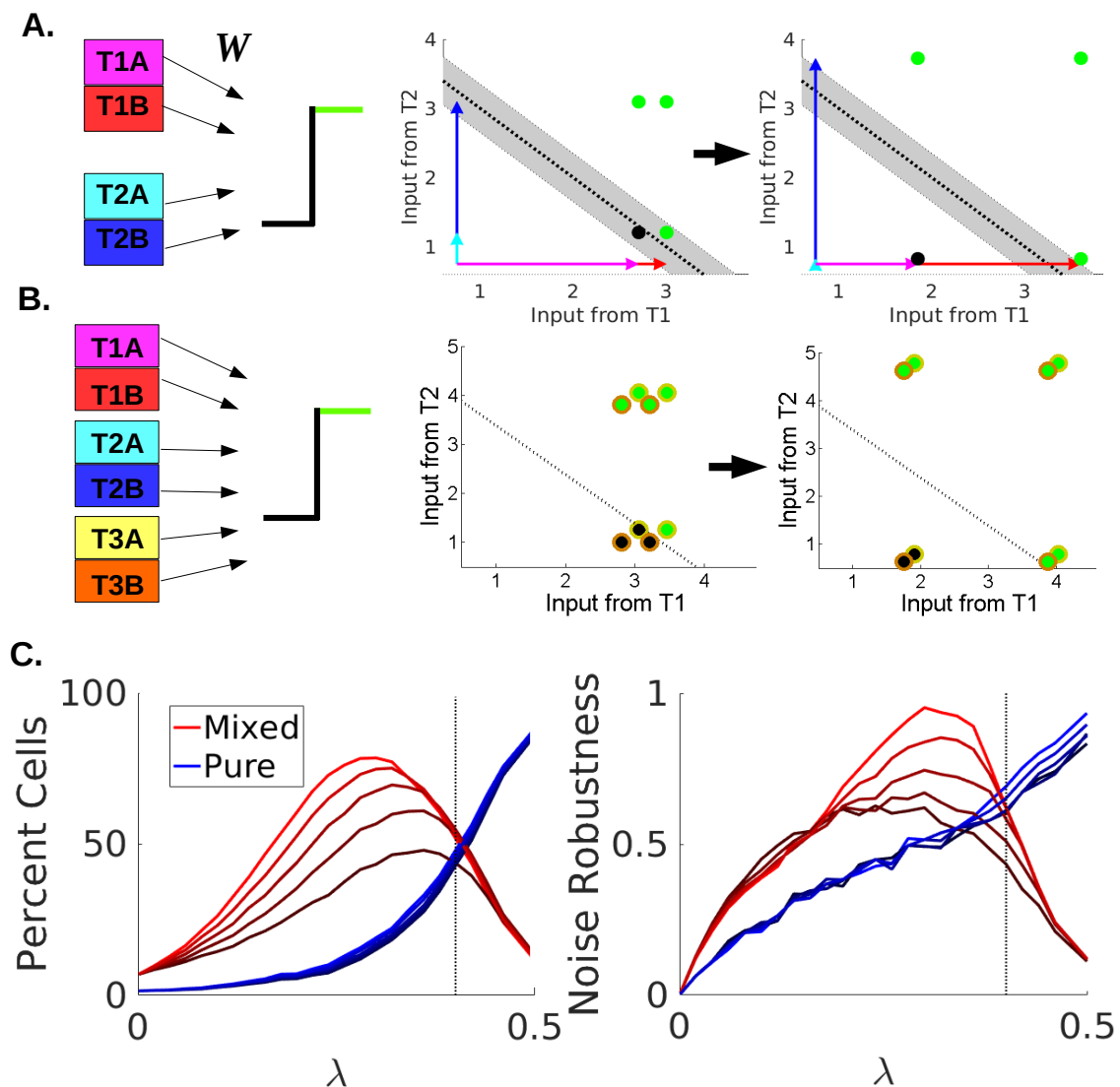
25

Figure 8: How learning impacts noise robustness A.) A simple toy cell (left) with 2 task variables is used to show the effects of learning. The 4 possible conditions are plotted as dots (green if above threshold, black if not), with the threshold as a dotted black line. Colored arrows represent the weights from each population. Before learning (middle), the cell's input on two of the conditions falls within the range of the shifting threshold created by additive noise (gray area). After learning, all conditions are outside the noise range. B.) A third task variable is added to the model and is another source of additive noise from the perspective of T1-T2 selectivity. The model's outputs are color-coded according to which T3 population is active. Weight arrows are omitted for visibility. After learning with $N_L = 2$, input strength from T3 populations are decreased and the points from the same T1-T2 condition are closer together (less noisy). C.) How the percent of cells with a given selectivity (left) and their noise robustness (right) change with constrained learning as a function of the threshold parameter $\lambda$. Learning steps are symbolized by increasing color brightness (the darkest line is the random model as displayed in Figure 7C, and the dashed line shows where the percent of mixed and pure are the same in the random model)

760 their weights increased, which diminishes the strength of task type preference. Lower
761 $N_L$, which minimizes preferences to other task variable identities, allows these clusters
762 to be tighter.

763     Finally, it is important to note that the strength of inputs shown in Figures 2
764 and 8 (the colored arrows) correspond to, in the full model, the summed input from
765 all cells representing a given task variable identity (i.e., $I_i^p$), not just to weights from
766 individual cells. These summed values are what need to change in order to expand the
767 representation and see the observed changes. This is important for why the Hebbian
768 procedure described here is effective at changing selectivity, as it assumes that many
769 cells, acting in unison to cause post-synaptic activity, would lead to the increase of their
770 individual synaptic weights, and thus an increase in the sum of those weights. Merely
771 increasing the variance of the individual weights does not cause such a coordinated
772 effect and would be less effective at driving these changes (as was shown in Figure 5C),
773 especially with larger input population size.

## 4. Discussion

775     Here, motivated by several theoretical proposals about properties that would ben-
776 efit encoding, we explored how prefrontal cortex represents task variables during a
777 complex task. In particular we were interested in measures of selectivity (particularly
778 nonlinear mixed selectivity), response density, and clustering of cell types according
779 to selectivity. By quantifying and measuring these properties in a PFC dataset, this
780 work connects theoretical literature with experimental data to give insight into how
781 PFC is able to support complex and flexible behavior. Furthermore, we explored how
782 these response properties could be generated by a simple network model. Through
783 this, we find evidence that the particular level of specialization and structure in the
784 PFC response is not achievable in a random network without Hebbian learning. After
785 Hebbian learning, the model—despite its relative simplicity—is able to capture many
786 response properties of PFC. The changes that come with learning act via an expansion
787 of the way cells represent conditions, and corresponding changes in noise robustness.

788     Interestingly, the variant of Hebbian learning that best matches the data is not the
789 most effective at increasing mixed selectivity. It may be that the more effective method
790 ("constrained" learning) would be too difficult to implement biologically, but perhaps
791 there is also a computational benefit to the balance of mixed and pure selectivity
792 found in the data. Particularly, in order to read out the task variable identity inputs
793 themselves, pure selectivity may be of more use. Retaining pure selectivity could be a
794 tool then for staying flexible.

795     In addition to retrospectively matching experimental results, this model also makes
796 predictions regarding how certain values should change with training. In particular,
797 clusters of cells defined by selectivity are expected to emerge with training and cell
798 responses should become less dense across conditions. Previous work [38] has shown
799 the value of mixed selectivity for the ability of a population to perform complex tasks.
800 This work shows that mixed selectivity increases with learning, and these changes
801 in PFC may correspond to increases in performance [33]. Perhaps surprisingly, this
802 model also predicts a concurrent, though small, decrease in pure selectivity. However,
803 studies that have tracked PFC responses during training show signs of these changes.
804 For example, in [27], the ability to decode the identity of the stimuli (in the comparable
805 portion of the trial) decreases slightly after training, suggesting a possible decrease in

pure selectivity. The ability to readout match/nonmatch of the two stimuli, however, increases dramatically, suggesting an increase in mixed selectivity. In [26], the amount of pure selectivity was measured directly pre- and post-training, and a significant drop in the percent of cells with pure selectivity was indeed observed. In hippocampus, an increase in mixed selectivity and slight decrease in pure was also observed with learning ([18]).

Our model makes many simplifying assumptions. The inputs, for instance, are binary cells that encode only the identity of different task variables. While this implies that the cells representing cue identities already have mixed selectivity (responding to the combination of the image and its place as either cue 1 or cue 2), it is still an assumption that the cells providing input to PFC are otherwise unmixed. This is something that, given current experimental evidence seems plausible [32], but would benefit from further experimental exploration.

Another valuable endeavor would be to expand this model in the temporal domain. Currently in the model, all the task variable inputs are given to the network simultaneously. In the experiment, of course, there is a delay between cue 1 and cue 2. Delay activity is known to exist in areas like IT [45, 12], and so this information could be being feed into PFC at the same time. But presumably, recurrent connections in PFC, and even possibly between PFC and its input areas, can enhance or alter selectivity. A recurrent model could also explore how PFC responses and representation vary over the time course of the trial, as recent experimental work has provided insight on this [31].

## 5. Acknowledgements

## 6. References

[1] Baktash Babadi and Haim Sompolinsky. Sparseness and expansion in sensory representations. *Neuron*, 83(5):1213–1226, 2014.

[2] Omri Barak, Mattia Rigotti, and Stefano Fusi. The sparseness of mixed selectivity neurons controls the generalization–discrimination trade-off. *The Journal of Neuroscience*, 33(9):3844–3856, 2013.

[3] Boris Barbour, Nicolas Brunel, Vincent Hakim, and Jean-Pierre Nadal. What can we learn from synaptic weight distributions? *TRENDS in Neurosciences*, 30 (12):622–629, 2007.

[4] Matthew M Botvinick. Hierarchical models of behavior and prefrontal function. *Trends in cognitive sciences*, 12(5):201–208, 2008.

[5] Jennifer N Bourne and Kristen M Harris. Coordination of size and number of excitatory and inhibitory synapses results in a balanced structural plasticity along mature hippocampal ca1 dendrites during ltp. *Hippocampus*, 21(4):354–373, 2011.

[6] Lior Bugatus, Kevin S Weiner, and Kalanit Grill-Spector. Task alters category representations in prefrontal but not high-level visual cortex. *NeuroImage*, 2017.

[7] Dean V Buonomano and Wolfgang Maass. State-dependent computations: spatiotemporal processing in cortical networks. *Nature Reviews Neuroscience*, 10(2): 113–125, 2009.

[8] Jean-René Duhamel, Carol L Colby, and Michael E Goldberg. Ventral intraparietal area of the macaque: congruent visual and somatic response properties. *Journal of neurophysiology*, 79(1):126–136, 1998.

[9] John Duncan. An adaptive coding model of neural function in prefrontal cortex. *Nature Reviews Neuroscience*, 2(11):820–829, 2001.

[10] Howard Eichenbaum, Andrew P Yonelinas, and Charan Ranganath. The medial temporal lobe and recognition memory. *Annu. Rev. Neurosci.*, 30:123–152, 2007.

[11] Stefano Fusi, Earl K Miller, and Mattia Rigotti. Why neurons mix: high dimensionality for higher cognition. *Current opinion in neurobiology*, 37:66–74, 2016.

[12] Joaquin M Fuster and John P Jervey. Neuronal firing in the inferotemporal cortex of the monkey in a visual memory task. *Journal of Neuroscience*, 2(3):361–375, 1982.

[13] Evarist Giné. Invariant tests for uniformity on compact riemannian manifolds based on sobolev norms. *The Annals of statistics*, pages 1243–1266, 1975.

[14] Michael J Goard, Gerald N Pho, Jonathan Woodson, and Mriganka Sur. Distinct roles of visual, parietal, and frontal motor cortices in memory-guided sensorimotor decisions. *Elife*, 5:e13764, 2016.

[15] Herbert Jaeger and Harald Haas. Harnessing nonlinearity: Predicting chaotic systems and saving energy in wireless communication. *science*, 304(5667):78–80, 2004.

[16] Ned H Kalin, Steven E Shelton, and Lorey K Takahashi. Defensive behaviors in infant rhesus monkeys: ontogeny and context-dependent selective expression. *Child development*, 62(5):1175–1183, 1991.

[17] Tara Keck, Taro Toyoizumi, Lu Chen, Brent Doiron, Daniel E Feldman, Kevin Fox, Wulfram Gerstner, Philip G Haydon, Mark Hübener, Hey-Kyoung Lee, et al. Integrating hebbian and homeostatic plasticity: the current state of the field and future research directions. *Phil. Trans. R. Soc. B*, 372(1715):20160158, 2017.

[18] Robert W Komorowski, Joseph R Manns, and Howard Eichenbaum. Robust conjunctive item–place coding by hippocampal neurons parallels learning what happens where. *Journal of Neuroscience*, 29(31):9918–9929, 2009.

[19] Robert W Komorowski, Carolyn G Garcia, Alix Wilson, Shoai Hattori, Marc W Howard, and Howard Eichenbaum. Ventral hippocampal neurons are shaped by experience to represent behaviorally relevant contexts. *Journal of Neuroscience*, 33(18):8079–8087, 2013.

[20] Ashok Litwin-Kumar, Kameron Decker Harris, Richard Axel, Haim Sompolinsky, and LF Abbott. Optimal degrees of synaptic connectivity. *Neuron*, 2017.

[21] Yi-Jiuan Lo and Mu-ming Poo. Activity-dependent synaptic competition in vitro: heterosynaptic suppression of developing synapses. *Science*, 254(5034):1019, 1991.

[22] Wolfgang Maass, Thomas Natschläger, and Henry Markram. Real-time computing without stable states: A new framework for neural computation based on perturbations. *Neural computation*, 14(11):2531–2560, 2002.

[23] Valerio Mante, David Sussillo, Krishna V Shenoy, and William T Newsome. Context-dependent computation by recurrent dynamics in prefrontal cortex. *Nature*, 503(7474):78–84, 2013.

[24] Kanti V Mardia and Peter E Jupp. Distributions on spheres. *Directional Statistics*, pages 159–192, 2000.

[25] Miriam LR Meister, Jay A Hennig, and Alexander C Huk. Signal multiplexing and single-neuron computations in lateral intraparietal area during decision-making. *Journal of Neuroscience*, 33(6):2254–2267, 2013.

[26] Travis Meyer, Xue-Lian Qi, Terrence R Stanford, and Christos Constantinidis. Stimulus selectivity in dorsal and ventral prefrontal cortex after training in working memory tasks. *Journal of Neuroscience*, 31(17):6266–6276, 2011.

[27] Ethan M Meyers, Xue-Lian Qi, and Christos Constantinidis. Incorporation of new information into prefrontal cortical activity after learning working memory tasks. *Proceedings of the National Academy of Sciences*, 109(12):4651–4656, 2012.

[28] Brian T Miller and Mark D'Esposito. Searching for the top in top-down control. *Neuron*, 48(4):535–538, 2005.

[29] Earl K Miller and Jonathan D Cohen. An integrative theory of prefrontal cortex function. *Annual review of neuroscience*, 24(1):167–202, 2001.

[30] Edvard I Moser, Emilio Kropff, and May-Britt Moser. Place cells, grid cells, and the brain's spatial representation system. *Annu. Rev. Neurosci.*, 31:69–89, 2008.

[31] John D Murray, Alberto Bernacchia, Nicholas A Roy, Christos Constantinidis, Ranulfo Romo, and Xiao-Jing Wang. Stable population coding for working memory coexists with heterogeneous neural dynamics in prefrontal cortex. *Proceedings of the National Academy of Sciences*, page 201619449, 2016.

[32] Marino Pagan, Luke S Urban, Margot P Wohl, and Nicole C Rust. Signals in inferotemporal and perirhinal cortex suggest an untangling of visual target information. *Nature neuroscience*, 16(8):1132–1139, 2013.

[33] Anitha Pasupathy and Earl K Miller. Different time courses of learning-related activity in the prefrontal cortex and striatum. *Nature*, 433(7028):873–876, 2005.

[34] Xue-Lian Qi and Christos Constantinidis. Variability of prefrontal neuronal discharges before and after training in a working memory task. *PLoS One*, 7(7): e41053, 2012.

[35] David Raposo, Matthew T Kaufman, and Anne K Churchland. A category-free neural population supports evolving demands during decision-making. *Nature neuroscience*, 17(12):1784–1792, 2014.

[36] Drew Rendall, Robert M Seyfarth, Dorothy L Cheney, and Michael J Owren. The meaning and function of grunt variants in baboons. *Animal Behaviour*, 57 (3):583–592, 1999.

[37] Mattia Rigotti, Daniel D Ben Dayan Rubin, Xiao-Jing Wang, and Stefano Fusi. Internal representation of task rules by recurrent dynamics: the importance of the diversity of neural responses. *Frontiers in computational neuroscience*, 4:24, 2010.

[38] Mattia Rigotti, Omri Barak, Melissa R Warden, Xiao-Jing Wang, Nathaniel D Daw, Earl K Miller, and Stefano Fusi. The importance of mixed selectivity in complex cognitive tasks. *Nature*, 497(7451):585–590, 2013.

[39] Maneesh Sahani and Jennifer F Linden. How linear are auditory cortical responses? *Advances in neural information processing systems*, pages 125–132, 2003.

[40] Massimo Scanziani, Robert C Malenka, and Roger A Nicoll. Role of intercellular interactions in heterosynaptic long-term depression. *Nature*, 380(6573):446, 1996.

[41] Mark G Stokes, Makoto Kusunoki, Natasha Sigala, Hamed Nili, David Gaffan, and John Duncan. Dynamic coding for cognitive control in prefrontal cortex. *Neuron*, 78(2):364–375, 2013.

[42] Sara M Szczepanski and Robert T Knight. Insights into human behavior from lesions to the prefrontal cortex. *Neuron*, 83(5):1002–1018, 2014.

[43] Melissa R Warden and Earl K Miller. Task-dependent changes in short-term memory in the prefrontal cortex. *The Journal of Neuroscience*, 30(47):15801–15810, 2010.

[44] Michael L Waskom, Dharshan Kumaran, Alan M Gordon, Jesse Rissman, and Anthony D Wagner. Frontoparietal representations of task context support the flexible control of goal-directed cognition. *The Journal of Neuroscience*, 34(32): 10743–10755, 2014.

[45] Luke Woloszyn and David L Sheinberg. Neural dynamics in inferior temporal cortex during a visual working memory task. *Journal of Neuroscience*, 29(17): 5494–5507, 2009.

[46] Jacqueline N Wood and Jordan Grafman. Human prefrontal cortex: processing and representational perspectives. *Nature Reviews Neuroscience*, 4(2):139–147, 2003.

31