

# Noise Correlations for Faster and More Robust Learning

 Matthew R. Nassar,<sup>1,2</sup> Daniel Scott,<sup>1,3</sup> and  Apoorva Bhandari<sup>1,3</sup>

<sup>1</sup>Robert J. and Nancy D. Carney Institute for Brain Science, Brown University, Providence, Rhode Island 02912-1821, <sup>2</sup>Department of Neuroscience, Brown University, Providence, Rhode Island 02912-1821, and <sup>3</sup>Department of Cognitive, Linguistic, and Psychological Sciences, Brown University, Providence, Rhode Island 02912-1821

Distributed population codes are ubiquitous in the brain and pose a challenge to downstream neurons that must learn an appropriate readout. Here we explore the possibility that this learning problem is simplified through inductive biases implemented by stimulus-independent noise correlations that constrain learning to task-relevant dimensions. We test this idea in a set of neural networks that learn to perform a perceptual discrimination task. Correlations among similarly tuned units were manipulated independently of an overall population signal-to-noise ratio to test how the format of stored information affects learning. Higher noise correlations among similarly tuned units led to faster and more robust learning, favoring homogenous weights assigned to neurons within a functionally similar pool, and could emerge through Hebbian learning. When multiple discriminations were learned simultaneously, noise correlations across relevant feature dimensions sped learning, whereas those across irrelevant feature dimensions slowed it. Our results complement the existing theory on noise correlations by demonstrating that when such correlations are produced without significant degradation of the signal-to-noise ratio, they can improve the speed of readout learning by constraining it to appropriate dimensions.

**Key words:** decision making; inductive biases; learning; noise correlations; perceptual learning; population code

## Significance Statement

Positive noise correlations between similarly tuned neurons theoretically reduce the representational capacity of the brain, yet they are commonly observed, emerge dynamically in complex tasks, and persist even in well-trained animals. Here we show that such correlations, when embedded in a neural population with a fixed signal-to-noise ratio, can improve the speed and robustness with which an appropriate readout is learned. In a simple discrimination task such correlations can emerge naturally through Hebbian learning. In more complex tasks that require multiple discriminations, correlations between neurons that similarly encode the task-relevant feature improve learning by constraining it to the appropriate task dimension.

## Introduction

The brain represents information using distributed population codes in which particular feature values are encoded by large numbers of neurons. One advantage of such codes is that a pooled readout across many neurons can effectively reduce the impact of stimulus-independent variability (noise) in the firing of individual neurons (Pouget et al., 2000). However, the extent to which this benefit can be employed in practice is constrained by noise correlations, or the degree to which stimulus-independent variability is shared across neurons in the population

(Averbeck et al., 2006). In particular, positive noise correlations between neurons that share the same stimulus tuning can reduce the amount of decodable information in the neural population (Averbeck et al., 2006; Hu et al., 2014; Moreno-Bote et al., 2014). Despite their detrimental effect on encoding, noise correlations of this type are reliably observed, even after years of training on perceptual tasks (Cohen and Kohn, 2011). Furthermore, noise correlations between neurons are dynamically enhanced under conditions where two neurons provide evidence for the same response in a perceptual categorization task (Cohen and Newsome, 2008), raising questions about whether they might serve a function rather than simply reflect a suboptimal encoding strategy.

At the same time, learning to effectively read out a distributed code also poses a significant challenge. Learning the appropriate weights for potentially tens of thousands of neurons in a low signal-to-noise regime is a difficult, high-dimensional problem, requiring a very large number of learning trials and entailing considerable risk of overfitting to specific patterns of noise encountered during learning trials. Nonetheless, people and animals can rapidly learn to perform perceptual discrimination

Received Dec. 2, 2020; revised June 8, 2021; accepted June 10, 2021.

Author contributions: M.R.N. and A.B. designed research; M.R.N. and D.S. performed research; M.R.N. and D.S. analyzed data; and M.R.N. wrote the paper.

This work was funded by National Institutes of Health Grant R00AG054732 (M.R.N.) and National Institute of Neurological Disorders and Stroke Grant R21NS108380 (A.B.). The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript. We thank Josh Gold, Rex Liu, Michael Frank, Drew Linsley, Chris Moore, and Jan Drugowitsch for discussion.

The authors declare no competing financial interests.

Correspondence should be addressed to Matthew R. Nassar at [matthew\\_nassar@brown.edu](mailto:matthew_nassar@brown.edu).

<https://doi.org/10.1523/JNEUROSCI.3045-20.2021>

Copyright © 2021 the authors

tasks, albeit with performance that does not approach theoretically achievable levels (Hawkey et al., 2004; Stringer et al., 2019). In comparison, deep neural networks capable of achieving human-level performance typically require a far greater number of learning trials than would be required by humans and other animals (Tsvividis et al., 2017). This raises the question of how brains might implement inductive biases to enable efficient learning in high-dimensional spaces.

Here we address open questions about noise correlations and learning by considering the possibility that noise correlations facilitate faster learning. Specifically, we propose that noise correlations aligned to task-relevant dimensions could reduce the effective dimensionality of learning problems, thereby making them easier to solve. For example, perceptual stimuli often contain a large number of features that may be irrelevant to a given categorization. At the level of a neural population, individual neurons may differ in the degree to which they encode task-irrelevant information, thus making the learning problem more difficult. In principle, noise correlations in the relevant dimension could reduce the effects of this variability on learned readout. Such an explanation would be consistent with computational analyses of Hebbian learning rules (Oja, 1982), which can both facilitate faster and more robust learning (Krotov and Hopfield, 2019) and, in turn, may induce noise correlations. We propose that faster learning of an approximate readout is made possible through low-dimensional representations that share both signal and noise across a large neural population. In particular, we hypothesize that representations characterized by enhanced noise correlations among similarly tuned neurons can improve learning by focusing adjustments of the readout onto task-relevant dimensions.

We explore this possibility using neural network models of a two-alternative forced-choice perceptual discrimination task in which the correlation among similarly tuned neurons can be manipulated independently of the overall population signal-to-noise ratio (SNR). Within this framework, noise correlations, which can be learned through Hebbian mechanisms, speed learning by forcing learned weights to be similar across pools of similarly tuned neurons, thereby ensuring learning occurs over the most task-relevant dimension. We extend our framework to a cued multidimensional discrimination task and show that dynamic noise correlations similar to those observed *in vivo* (Cohen and Newsome, 2008) speed learning by constraining weight updates to the relevant feature space. Our results demonstrate that when information is extrinsically limited, noise correlations can make learning faster and more robust by controlling the dimensions over which learning occurs.

## Materials and Methods

Our goal was to understand the computational principles through which correlations in the activity of similarly tuned neurons affect the speed with which downstream neurons could learn an effective readout. Previous work has demonstrated that manipulating noise correlations while maintaining a fixed variance in the firing rates of individual neurons leads to changes in the theoretical encoding capacity of a neural population (Averbeck et al., 2006; Moreno-Bote et al., 2014). To minimize the potential impact of such encoding differences, we took a different approach; rather than setting the variance of individual neurons in our population to a fixed value, we set the signal-to-noise ratio of our population to a fixed value. Thus, our approach does not ask how maximum information can be packed into the activity of a given neural population but rather how the strategy for packing a fixed amount of information in a population affects the speed with which an appropriate

readout of that information can be learned. We implement this approach in a set of neural networks described in more detail below.

**Learning readout in perceptual learning task.** Simulations and analyses for a simple perceptual discrimination task were performed with a simplified and statistically tractable two-layer feedforward neural network (see Fig. 3A). The input layer consisted of two homogenous pools of 100 units that were each identically tuned to one of two motion directions (left, right). On each trial normalized firing rates for the neural population were drawn from a multivariate normal distribution that was specified by a vector of stimulus-dependent mean firing rates (signal: +1 for preferred stimulus, −1 for nonpreferred stimulus) and a covariance matrix. All elements of the covariance matrix corresponding to covariance between units that were tuned to different stimuli were set to zero. The key manipulation was to systematically vary the magnitude of diagonal covariance components (e.g., noise in the firing of individual units) and the within-pool covariance elements (e.g., shared noise across identically tuned neurons) while maintaining a fixed level of variance in the summed population response for each pool as follows:

$$\sigma_{pool}^2 = n\sigma_{unit}^2 + n(n-1)\text{Cov}(\text{withinpool}), \quad (1)$$

where  $\sigma_{pool}^2$  is the variance on the sum of normalized firing rates from neurons within a given pool,  $n$  is the number of units in the pool and the within-pool covariance,  $\text{Cov}(\text{withinpool})$ , specifies the covariance of pairs of units belonging to the same pool. SNR was defined as the population signal (preferred/antipreferred) divided by the SD of the population response in the signal dimension. SNR was set to be 2 for each individual pool of neurons, leading to a signal-to-noise ratio for the entire population (both pools) equal to  $2\sqrt{2}$ . Given this constraint, the fraction of noise that was shared across neurons within the same pool was manipulated as follows:

$$\sigma_{unit}^2 = \frac{\sigma_{pool}^2}{n + n(n-1)\phi} \quad (2)$$

$$\text{Cov}(\text{withinpool}) = \phi\sigma_{unit}^2, \quad (3)$$

where  $\phi$  reflects the fraction of noise that is correlated across units, which we refer to in the text as noise correlations. Noise correlations ( $\phi$ ) were manipulated across values ranging from 0 to 0.2 for simulations. Note that because  $\phi$  appears in the denominator of Equation 2, adding noise correlations while sustaining a fixed population signal-to-noise ratio leads to lower variance in the firing rates of single neurons, differing from previous theoretical assumptions (see Fig. 2).

The input layer of the neural network was fully connected to an output layer composed of two output units representing left and right responses. Output units were activated on a given trial according to a weighted function of their inputs as follows:

$$F_{output} = wF_{input}, \quad (4)$$

where  $F_{output}$  is a vector of firing rates of output units,  $F_{input}$  is a vector of firing rates of the input units, and  $w$  is the weight matrix. Firing of an individual output unit can also be written as a weighted sum over input unit activity as in the following:

$$F_j = \sum_{i=1}^{200} w_{i,j}F_i, \quad (5)$$

where  $F_j$  reflects the firing of the  $j$ th output unit,  $F_i$  reflects the firing of the  $i$ th input unit, and  $w_{i,j}$  reflects the weight of the connection between the  $i$ th input unit and the  $j$ th output unit. Actions were selected as a softmax function of output firing rates as follows:

$$p(A_j) = \frac{e^{\beta F_j}}{\sum_k e^{\beta F_k}}, \quad (6)$$

where  $\beta$  is an inverse temperature, which was set to a relatively deterministic value (10,000). Learning was implemented through reinforcement of weights to the selected output neuron (subscript  $j$  below) as follows:

$$\Delta w_{i,j} = \alpha \delta F_i, \quad (7)$$

where  $F_i$  is the normalized firing rate of the  $i$ th input neuron,  $\delta$  is the reward prediction error experienced on a given trial (+0.5 for correct trials and  $-0.5$  for error trials), and  $\alpha$  is a learning rate (set to 0.0001 for simulations; see Fig. 2). The network was trained to correctly identify two stimuli (each of which was preferred by a single pool of input neurons) over 100 trials (of which the last 20 trials were considered testing). Simulations were repeated 1000 times for each level of  $\phi$ , and performance measures were averaged across all repetitions. Mean accuracy per trial across all simulations was convolved with a Gaussian kernel (SD = 0.5 trials) for plotting (see Fig. 2B). Mean accuracy across the final 20 trials was used as a measure of final accuracy (see Fig. 2E). Statistics on model performance were computed as Pearson correlations between noise correlations  $\phi$  and performance measures across all simulations and repetitions.

**Analytical learning trajectories.** One advantage of our simple network architecture is its mathematical tractability. To complement the simulations described above, we also explored learning in the network analytically. Specifically, we decomposed weight updates into two categories: weight updates in the signal dimension and weight updates perpendicular to the signal dimension. Weight updates in the signal dimension improved performance through alignment with the signal itself, whereas weight updates in the perpendicular dimension limited performance through chance alignment with trial-to-trial noise. An intuition for our approach and derivation are provided below.

The two-alternative discrimination task is a one-dimensional signal detection problem because it depends only on the difference between two scalars. In particular, if  $y = [y_1, y_2]$  denotes the readout activity in the pair of pools, and  $r$  denotes the response (e.g.,  $r = -1$  is respond left, and  $r = 1$  is respond right), then  $r = r(y_1 - y_2) = r(\Delta y)$ . In addition,  $\Delta y = w_1 x - w_2 x \equiv \Delta w x$ , where  $x$  reflects the firing rates of the input units and  $w_1$  reflects the vector of weights mapping input activation onto output unit 1 ( $y_1$ ). To determine how accuracy is affected by noise correlations, we ask how Mahalanobis distance ( $d'$ ), mean separation ( $d$ ), and signal variance ( $\sigma_{s^*}^2$ ) diverge over training time for the different noise correlation conditions. The effective variance,  $\sigma_{s^*}^2$ , differs from the true noise variance in the signal dimension because of the fact that out-of-signal-dimension noise is transferred into the signal dimension by imperfect readout weights. Intuitively, learning speed may be improved by noise correlations because less out-of-dimension noise is learned into the weights, thereby reducing the transfer of out-of-dimension noise into the signal space on any given trial.

The logic of training is as follows. On a correct trial, the weights to the chosen unit are incremented by a multiple of the input vector  $x$ , as in the following:

$$w_i \rightarrow w_i + \alpha \delta x. \quad (8)$$

Here  $\alpha$  reflects a positive learning rate,  $x$  reflects the activity of the input units, and  $\delta$  is the reward prediction error, which we use as the absolute reward prediction error instead of the signed one in this section for convenience.

Now the input is a sum of signal and zero mean noise as follows:

$$x = \mu + \xi. \quad (9)$$

The expectation of noise is zero ( $E(\xi) = 0$ ), and the signal  $\mu$  can take only two values  $\mu \in \{\pm \mu_0\}$ . Therefore, if the weights start from some value  $\Delta w(0)$ , we will find the following:

$$E[\Delta w(t)] = t \alpha \delta \mu_0 + \Delta w(0), \quad (10)$$

where  $t$  reflects the current timestep of learning. In words, we expect the amount of signal in the weights to increase linearly over time. This means that we expect the response to a noise-free signal ( $\mu_0$ ) after  $t$  time-steps to be the following:

$$\Delta y(\mu_0, t) = \Delta w(t) \mu_0 + \Delta w_0 \mu_0 = t \alpha \delta \|\mu_0\|^2 + \Delta w_0 \mu_0. \quad (11)$$

This is the measure  $d$  between the two Gaussian peaks in the one-dimensional signal detection problem described above. Below, we ignore the initial weight term as it does not change over time. To compute accuracy and  $d'$  over training time, we also need to compute the effective variance along the signal dimension. First we note that the noise can be decomposed as follows:

$$\xi = \xi_s + \xi_{\perp}, \quad (12)$$

where  $\xi_s$  and  $\xi_{\perp}$  are orthogonal components of the noise in the signal dimension ( $\xi_s$ ) and perpendicular to the signal dimension ( $\xi_{\perp}$ ). Here we consider cases where the noise along the signal dimension ( $\xi_s$ ) has constant variance, following the assumption that SNR is set to a constant value and that the mean signal is the same for all noise correlation conditions.

The difference  $\Delta y$  on any given trial decomposes into a sum of terms, one reflecting a weight-based transfer of signal and one reflecting the transfer of orthogonal noise. This latter term arises because the weights are not, at any finite time, a perfect matched filter for the signal. Letting subscripts  $s$  and  $\perp$  continue to denote signal and perpendicular dimensions, we have the following:

$$\Delta y = \Delta w x \quad (13)$$

$$\Delta y = (\Delta w_s + \Delta w_{\perp})(\mu + \xi_s + \xi_{\perp}) \quad (14)$$

$$\Delta y = \Delta w_s(\mu + \xi_s) + \Delta w_{\perp} \xi_{\perp}, \quad (15)$$

where the final equation reflects the absence of terms that have zero products by definition of the perpendicular subspaces. The variance of  $\Delta y$  can be computed using independence and orthogonality properties as follows:

$$\text{Var}(\Delta y) = \text{Var}(\Delta w_s(\mu + \xi_s) + \Delta w_{\perp} \xi_{\perp}) \quad (16)$$

$$\text{Var}(\Delta y) = \Delta w_s^2 E[\xi_s^2] + \Delta w_{\perp}^2 E[\xi_{\perp}^2]. \quad (17)$$

For any given network, the term  $\Delta w(t)_{\perp}^2$  is a mean-zero diffusion process arising from the fact that noise is added to the weights at every timestep. For the Gaussian white noise case,  $\Delta w(t)_{\perp}^2$  is equivalent to Brownian motion in the  $(n - 1)$  dimensions perpendicular to the signal. Because  $(n - 1)$  is not small, the summed empirical variance of these processes, operative on each component, is likely to be close to the theoretical total variance. If we split the term  $\Delta w(t)_{\perp}^2$  into the  $(n - 1)$  components and index them with  $i$ , this gives the following:

$$\Delta w_{\perp}^i = \alpha \delta \sqrt{\frac{t}{n-1}} \sigma_{\perp}. \quad (18)$$

The denominator of  $(n - 1)$  appears here because Brownian motion determines growth in the variance of each of the  $(n - 1)$  perpendicular noise directions among which the total variance  $\sigma_{\perp}$  is distributed. Technically, our manipulation of the noise covariance fixes the variance in a second direction of the space as well, so that noise variance is actually evenly distributed over only  $(n - 2)$  of the  $(n - 1)$  perpendicular dimensions, but this inhomogeneity is inconsequential if  $n$  is not small.

In effect, we are ignoring an order 1 term relative to an order  $n$  term for simplicity. To understand how perpendicular weights grow with time, we need only to determine  $\sigma_{\perp}(\phi)$ , where  $\phi$  is the parameter controlling the noise covariance matrix in our simulations. Specifically, the first row of the covariance matrix takes the following form:

$$\sum (\xi)_1 = [b, \phi b, \phi b, \dots, \phi b, 0, \dots, 0]. \quad (19)$$

Using the additional fact that row sums are set to  $\sigma_s^2$  to control the signal variance, we find the following:

$$b + \left(\frac{n}{2} - 1\right) \phi b = \sigma_s^2 \quad (20)$$

$$b = \frac{2\sigma_s^2}{2 + (n - 2)\phi}. \quad (21)$$

Because the eigenvalues of  $\sum(\xi)$  are the variances in different dimensions of the space, we can find the total variance perpendicular to the signal by subtracting the known signal variance from the trace of  $\sum(\xi)$  as follows:

$$\sigma_{\perp}^2 = \text{Var}(\xi_{\perp}) = \text{Tr}\left(\sum(\xi)\right) - \text{Var}(\xi_s) \quad (22)$$

$$\sigma_{\perp}^2 = \text{Var}(\xi_{\perp}) = nb - \sigma_s^2. \quad (23)$$

Putting this together with previous results, we have the following:

$$\text{Var}(\Delta y) = (t\alpha\delta\mu\sigma_s)^2 + \frac{t(\alpha\delta)^2\sigma_{\perp}^4}{n-1}. \quad (24)$$

This provides analytic prediction for the variance of our readout decision variable  $\Delta y$  after learning for  $t$  trials, using a learning rate  $\alpha$  to learn from prediction errors of magnitude  $\delta$ . Note that  $\sigma_s$  was fixed in our simulations but that  $\sigma_{\perp}^4$  depends on  $\phi$  through  $b$ , so that larger values of  $\phi$  lead to smaller values of  $b$ , and thus a smaller  $\sigma_{\perp}^2$ , reducing the second term in Equation 24. Furthermore, as the first term in Equation 24 scales with  $t^2$ , its contributions dominate as more trials are observed. This leads to identical asymptotic variance in the limit of large  $t$  because the first term does not depend on  $\phi$ .

By combining the mean and variance information in Equations 11 and 24 we computed accuracy as one minus the cumulative probability density of the Gaussian distribution as follows:  $N(t\alpha\delta\|\mu_0\|^2 + (t\alpha\delta\mu\sigma_s)^2 + \frac{t(\alpha\delta)^2\sigma_{\perp}^4}{n-1})$ , evaluated from negative infinity to zero.

*Noise correlations with fixed signal-to-noise ratio and single-unit variance.* Noise correlations produced by the simulations above lead to reductions in the overall variance of single-unit firing rates. To validate that our results depend on maintaining signal-to-noise, rather than depending on single-unit variance, we also consider the case where noise correlations are introduced with a fixed level of single-unit variance. In this case, signal-to-noise ratio was maintained by scaling the amount of signal according to the level of noise correlations (see <https://github.com/NassarLab/NoiseCorrelation> for the full derivation) as follows:

$$S_{neuron} = \sqrt{\frac{\sigma_{unit}^2(1 + (n-1)\phi)}{n}}, \quad (25)$$

where  $S_{neuron}$  reflects the amount of signal provided by each unit,  $\sigma_{unit}^2$  reflects a fixed variance assigned to each unit,  $n$  reflects the number of units in the pool, and  $\phi$  reflects the level of noise correlations. Thus, when we simulated correlated noise using this equation, neurons maintained the same variance ( $\sigma_{unit}^2$ ) but increased

the signal of the neurons relative to the zero noise correlation condition ( $\phi = 0$ ).

*Noise correlations that are bounded to a maximum signal-to-noise ratio.* To examine the importance of our assumption regarding fixed signal-to-noise ratio, we also considered a parameterized model, where signal ( $S_{neuron}$ ) was set according to a linear mixture as follows:

$$S_{neuron} = m\sqrt{\frac{\sigma_{unit}^2(1 + (n-1)\phi)}{n}} + (1-m)\sqrt{\frac{\sigma_{unit}^2}{n}}, \quad (26)$$

where  $m$  is a mixing parameter that combines the signal producing a fixed signal-to-noise ratio (first term) with a fixed signal that does not depend on the level of noise correlations (second term). When  $m$  is set to 1, this parameterized model obeys our assumptions regarding fixed signal-to-noise ratio, but when  $m$  is set to 0, the model conforms to more standard assumptions regarding fixed single-unit variance and signal.

*Hebbian learning of noise correlations in three-layer network.* We extended the two-layer feedforward architecture described above to include a third hidden layer to test whether Hebbian learning could facilitate the production of noise correlations among similarly tuned neurons (see Fig. 5A). The input layer was fully connected to the hidden layer, and each layer contained 200 neurons. In the input layer, neurons were homogeneously tuned (100 leftward, 100 rightward) as described above, with  $\phi$  set to zero (e.g., no noise correlations). Weights to the hidden layer were initialized to favor one-to-one connections between input layer units and hidden layer units by adding a small normal random weight perturbation (mean = 0, SD = 0.01) to an identity matrix (although an alternate initialization produced qualitatively similar results). During learning, weights between the input and hidden layer were adjusted according to a normalized Hebbian learning rule as follows:

$$\Delta W = \alpha_{hebb} F_1 F_2, \quad (27)$$

where  $F_1$  is a normalized vector of firing rates corresponding to the input layer, and  $F_2$  is a normalized vector of firing rates corresponding to the hidden layer units. The learning rate for Hebbian plasticity ( $\alpha_{hebb}$ ) was set to 0.00005 for simulations (see Fig. 4). Weights were normalized after Hebbian learning to ensure that the Euclidean norm of the incoming weights to each unit in layer two was equal to one. The model was trained over 100 trials in the same perceptual discrimination task described above, and an additional 100 trials of the task were completed to measure emergent noise correlations in the hidden layer. Noise correlations were measured by regressing out variance attributable to the stimulus on each trial, and then computing the Pearson correlation of residual firing rate across each pair of neurons for the 100 testing trials (see Fig. 4B,C).

*Learning readout in multiple discrimination task.* To test the impact of contextual noise correlations on learning (Cohen and Newsome, 2008), the perceptual discrimination task was extended to include two dimensions and two interleaved trial types, one in which an up/down discrimination was performed (vertical), and one in which a right/left discrimination was performed (horizontal). Each trial contained motion on the vertical axis (up or down) and on the horizontal axis (left or right), but only one of these motion axes was relevant on each trial as indicated by a cue.

To model this task, we extended our two-layer feedforward network to include four pools of input units, four output units, and two task units (see Fig. 5A). Each homogenous pool of 100 input units encoded a conjunction of the movement directions (up-right, up-left, down-right, down-left). On each trial, the mean firing rate of each input unit population was determined according to the tuning preferences of each unit population as follows:

$$\mu = V + H, \quad (28)$$

where  $V$  was  $+1/-1$  for trials with the preferred/antipreferred vertical motion direction,  $H$  was  $+1/-1$  for trials with the preferred/antipreferred horizontal motion direction. Firing rates for individual neurons were

sampled from a multivariate Gaussian distribution with mean  $\mu$  and a covariance matrix that depended on trial type (vertical vs horizontal) and the level of same-pool, relevant-pool, and irrelevant-pool correlations.

To create a covariance matrix, we stipulated a desired SEM for summed population activity ( $SEM = 20$  for simulations; see Fig. 7) and determined the summed population variance that would correspond to that value ( $\sigma_{pool}^2$ ). We then determined the variance on individual neurons that would yield this population response under a given noise correlation profile as follows:

$$\sigma_{unit}^2 = \frac{\sigma_{pool}^2}{n + n(n-1)\phi_{same} + n^2\phi_{relevant} - n^2\phi_{irrelevant}}, \quad (29)$$

where  $\phi_{same}$  is the level of same-pool correlations (range, 0–0.2 in our simulations),  $\phi_{relevant}$  is the level of relevant-pool correlations (range, 0–0.2 in our simulations),  $\phi_{irrelevant}$  is the level of irrelevant-pool correlations (range, 0–0.2 in our simulations). Note that increasing the same pool or in-pool correlations reduces the overall variance to preserve the same level of variance on the task-relevant dimension in the population response, but increasing irrelevant-pool correlations has the opposite effect. Covariance elements of the covariance matrix were determined as follows:

$$Cov(samepool) = \phi_{same} \sigma_{unit}^2 \quad (30)$$

$$Cov(relevantpool) = \phi_{relevant} \sigma_{unit}^2 \quad (31)$$

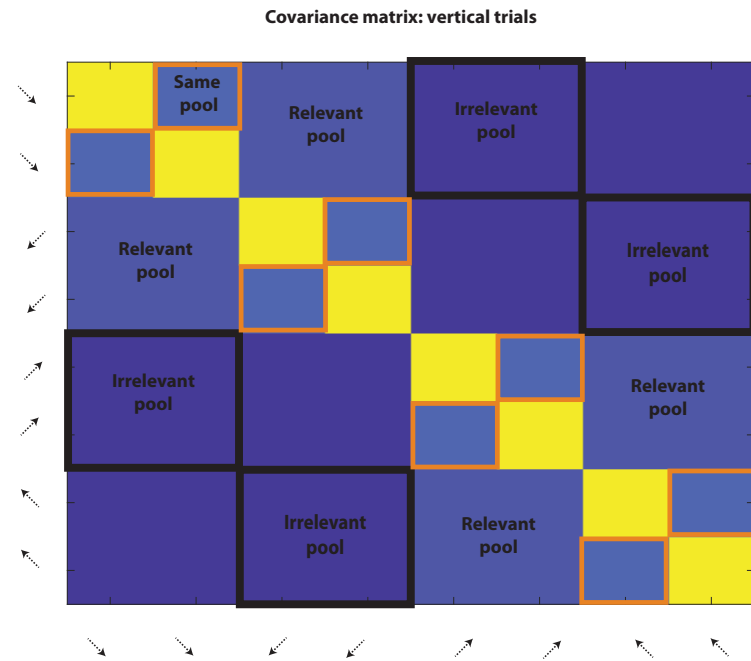
$$Cov(irrelevantpool) = \phi_{irrelevant} \sigma_{unit}^2 \quad (32)$$

The variance and covariance values above were used to construct a covariance matrix for each trial type (vertical/horizontal) as depicted in Figure 1

Output units corresponded to the four possible task responses (up, down, left, right) and were activated according to a weighted sum of their inputs as described previously. Task units were modeled as containing perfect information about the task cue (vertical vs horizontal), and each task unit projected with strong fixed weights (1000) to both responses that were appropriate for that task. Decisions were made on each trial by selecting the output unit with the highest activity level. Weights to a chosen output unit were updated using the same reinforcement learning procedure described in the two-alternative perceptual learning task.

## Results

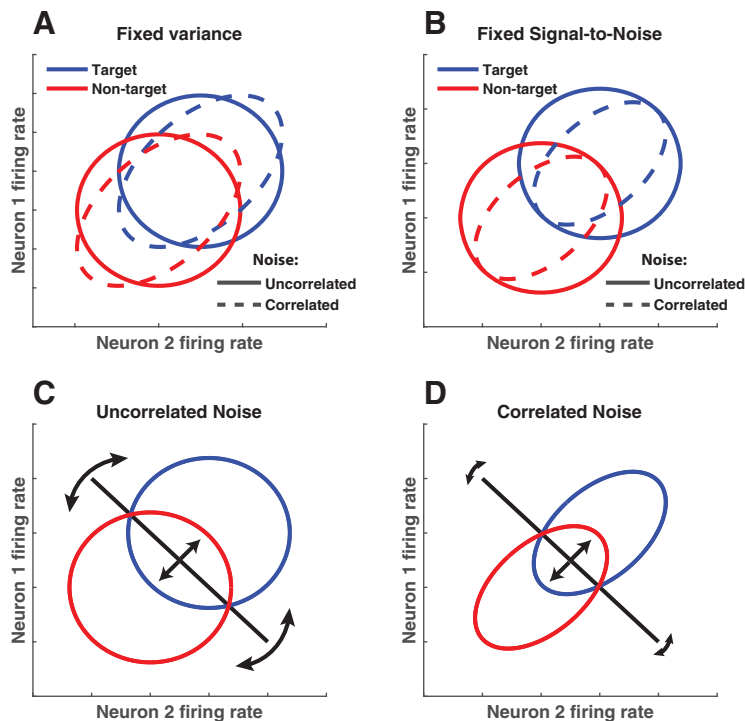
We examine how noise correlations affect learning in a simplified neural network where the appropriate readout of hundreds of weakly tuned units is learned over time through reinforcement. To isolate the effects of noise correlations on learning, rather than their effects on other factors such as representational capacity, we consider population encoding schemes at the input layer that can be constrained to a fixed signal-to-noise ratio. This assumption differs from previous work on noise correlations where the variance of the neural population is assumed to be fixed, and covariance is changed to produce noise correlations, thereby affecting the representational capacity of the population (Fig. 2A; Averbeck et al., 2006; Moreno-Bote et al., 2014). Under our assumptions, a fixed signal-to-noise ratio can be achieved for any level of noise correlations by scaling the variance (Fig. 2B; Eqs. 1–3), or



**Figure 1.** Schematic of covariance matrix for two-dimensional motion discrimination task. The covariance between units with different motion tuning (reflected by the arrows labeling columns and rows) is schematically represented for a simplified input layer, where only two identically tuned neurons are in each pool (In actual simulations there were 100 units per pool). Same pool correlations are controlled by covariance elements between neurons with identical tuning (orange boxes). Relevant pool correlations are controlled by covariance elements between neurons that are similarly tuned to the task-relevant feature. Task-irrelevant correlations are controlled by covariance elements between neurons that are similarly tuned to the task-irrelevant feature. The covariance matrix shown here is for a vertical trial; on a horizontal trial the irrelevant pool and relevant pool locations would be reversed. Covariance elements for pairs of neurons that differed in tuning on both dimensions were set to zero. Each input population has been depicted as two units here for presentation purposes. Background color reflects the case where same pool correlations = 0.2 and relevant pool correlations = 0.1.

alternately scaling the magnitude of the signal (Eq. 25). Although we do not discount the degree to which noise correlations affect the encoding potential of neural populations, we believe that in many cases the relevant information is limited by extrinsic factors (e.g., the stimulus itself or upstream neural populations providing input; Ecker et al., 2011; Beck et al., 2012; Kanitscheider et al., 2015). Under such conditions, reducing noise correlations can increase information only until it saturates because all the available incoming information is encoded. Beyond that, increasing encoding potential is not possible as it would be tantamount to the population creating new information that was not communicated by inputs to the population. Therefore, our framework can be thought of as testing how best to format limited available information in a neural population to ensure that an acceptable readout can be rapidly and robustly learned.

We propose that within this framework, noise correlations of the form that have previously been shown to limit encoding are beneficial because they constrain learning to occur over the most relevant dimensions. In general, a linear readout can be thought of as a hyperplane serving as a classification boundary in an  $N$  dimensional space, where  $N$  reflects the number of neurons in a population. Learning in such a framework involves adjustments of the hyperplane to minimize classification errors. The most useful adjustments are in the dimension that best discriminates signal from noise (Fig. 2C,D, central arrows), but adjustments may also occur in dimensions orthogonal to the relevant one (such as twisting of the hyperplane, depicted by curved arrows in Fig. 2C,D) that could potentially impair performance or slow



**Figure 2.** Modeling noise correlations with extrinsic constraint on signal-to-noise ratio. **A**, Previous work has modeled noise correlations by assuming that population variance is fixed and that covariance is manipulated to produce noise correlations. Under such assumptions, the firing rate of two similarly tuned neurons is plotted in the absence (solid line) or presence (dotted line) of information-limiting noise correlations. **B**, Here we assume that the signal-to-noise ratio of the neural population is limited to a fixed value so that noise correlations between similarly tuned neurons do not affect theoretical performance. Thus, the percentage overlap of blue (target) and red (nontarget) activity profiles does not differ in the presence (dotted line) or absence (solid line) of noise correlations. **C**, **D**, Under this assumption, noise correlations among similarly tuned neurons could compress the population activity to a plane orthogonal to the optimal decision boundary, thereby minimizing boundary adjustments in irrelevant dimensions (**C**) and maximizing boundary adjustments on relevant ones (**D**).

down learning. Our motivating hypothesis is that by focusing population activity into the task-relevant dimension, noise correlations can increase the fraction of hyperplane adjustments that occur in the task-relevant dimension (Figure 2D), thus reducing the effective dimensionality of readout learning.

### Noise correlations enable faster learning in a fixed signal-to-noise regime

To test our overarching hypothesis, we constructed a fully connected two-layer feedforward neural network in which input layer units responded to one of two stimulus categories (pool 1 and pool 2), and each output unit produced a response consistent with a category perception (Figure 3A, left/right units). On each trial, the network was presented with one stimulus at random, and input firing for each pool was drawn from a multivariate Gaussian with a covariance that was manipulated while preserving the population signal-to-noise ratio. Output units were activated according to a weighted average of inputs, and a response was selected according to output unit activations. On each trial, weights to the selected action were adjusted according to a reinforcement learning rule that strengthened connections that facilitated a rewarded action and weakened connections that facilitated an unrewarded action (Law and Gold, 2009).

Noise correlations led to faster and more robust learning of the appropriate stimulus-response mapping. All neural networks learned to perform the requisite discrimination, but neural networks that employed correlations among similarly tuned

neurons learned more rapidly (Fig. 3B). After learning, networks that employed such noise correlations assigned more homogeneous weights to input units of a given pool than did networks that lacked noise correlations (compare Fig. 3C,D). This led to better trained task performance (Fig. 3E; Pearson correlation between noise correlations and test performance,  $R = 0.29$ ,  $p < 10e-50$ ) and greater robustness to adversarial noise profiles (Fig. 3F;  $R = 0.81$ ,  $p < 10e-50$ ) in the networks that employed noise correlations. Critically, these learning advantages emerged despite the fact that optimal readout of all networks achieved similar levels of performance and robustness (Fig. 3E,F; compare optimal readout across conditions).

### Learning benefits from noise correlations are greatest for large, low SNR populations

To better understand how noise correlations promoted faster learning, we developed an analytical method for describing learning trajectories (see above, Materials and Methods). Our method considered the impacts of the following two influences on weight updates over time: (1) weight updates in the signal dimension that tend to align with the signal and improve performance and (2) weight updates perpendicular to the signal dimension, which through chance alignment with trial-to-trial firing rate variability allow noise to have an impact on decisions and therefore hinder performance (Fig. 4A). Noise correlations implemented using our methods decreased

the latter form of weight updates (Fig. 4B), leading updates in the signal dimension to more quickly dominate performance (Fig. 4C), thereby speeding analytical predictions for learning (Fig. 4D,E). The analytically derived learning advantage for fixed-SNR noise correlations was greatest for situations in which SNR was relatively low and neural populations were large (Fig. 4F).

The advantage of noise correlations for learning speed did not depend on specific assumptions about whether SNR was balanced by adjusting signal or noise. We employed an alternate method for creating fixed-SNR noise correlations that amplified signal, rather than reducing variance, to maintain SNR for higher levels of noise correlation (Eq. 25). Such noise correlations could be thought of as reflecting amplification of both signal and shared noise that would result from top-down recurrent feedback (Haefner et al., 2016). Under such assumptions, noise correlations sped learning and led to more robust weight profiles, similar to our previous simulations (Fig. 5A).

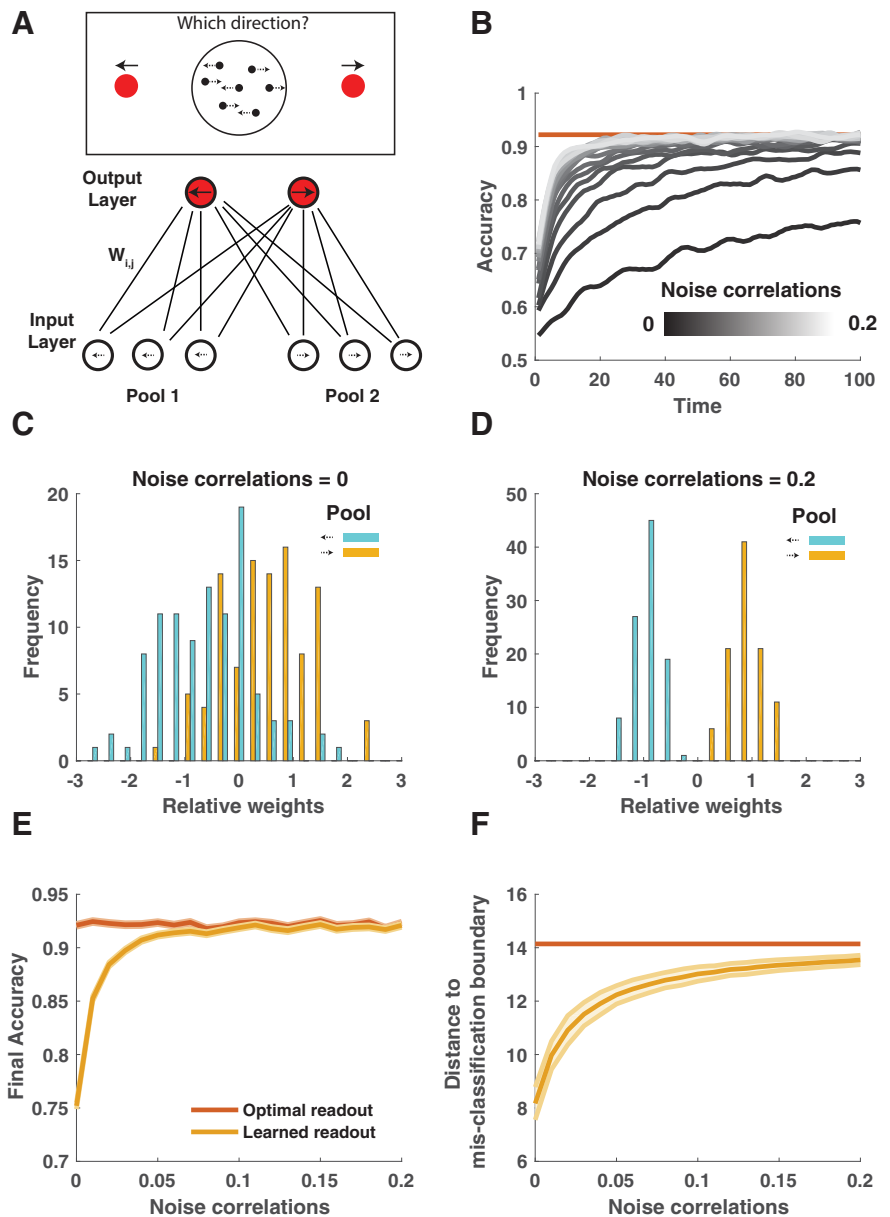
### Noise correlations that do not maintain signal-to-noise ratio can introduce a trade-off between learning speed and asymptotic performance

In contrast, our learning speed results depended critically on the assumption that signal-to-noise ratio is maintained across different levels of noise correlation. To test this dependency, we

examined performance of a family of models that contained a single parameter, allowing them to range in assumptions from fixed SNR ( $m = 1$ ) to fixed single-unit signal and variance ( $m = 0$ ), analogous to assumptions of Averbeck et al. (2006). Consistent with our previous results, noise correlations improve learning in the  $m = 1$  case, and consistent with Averbeck et al. (2006), asymptotic performance is reduced by noise correlations in the  $m = 0$  case (Fig. 5B). Interestingly, for intermediate assumptions between these two extremes, noise correlations promote faster learning, improving performance in the short run but at the cost of lower asymptotic accuracy. Thus, under such assumptions, adjusting noise correlations between similarly tuned neurons could potentially optimize a trade-off between short-term gains from rapid learning and long-term gains from higher asymptotic performance.

### Hebbian learning can produce useful noise correlation structure

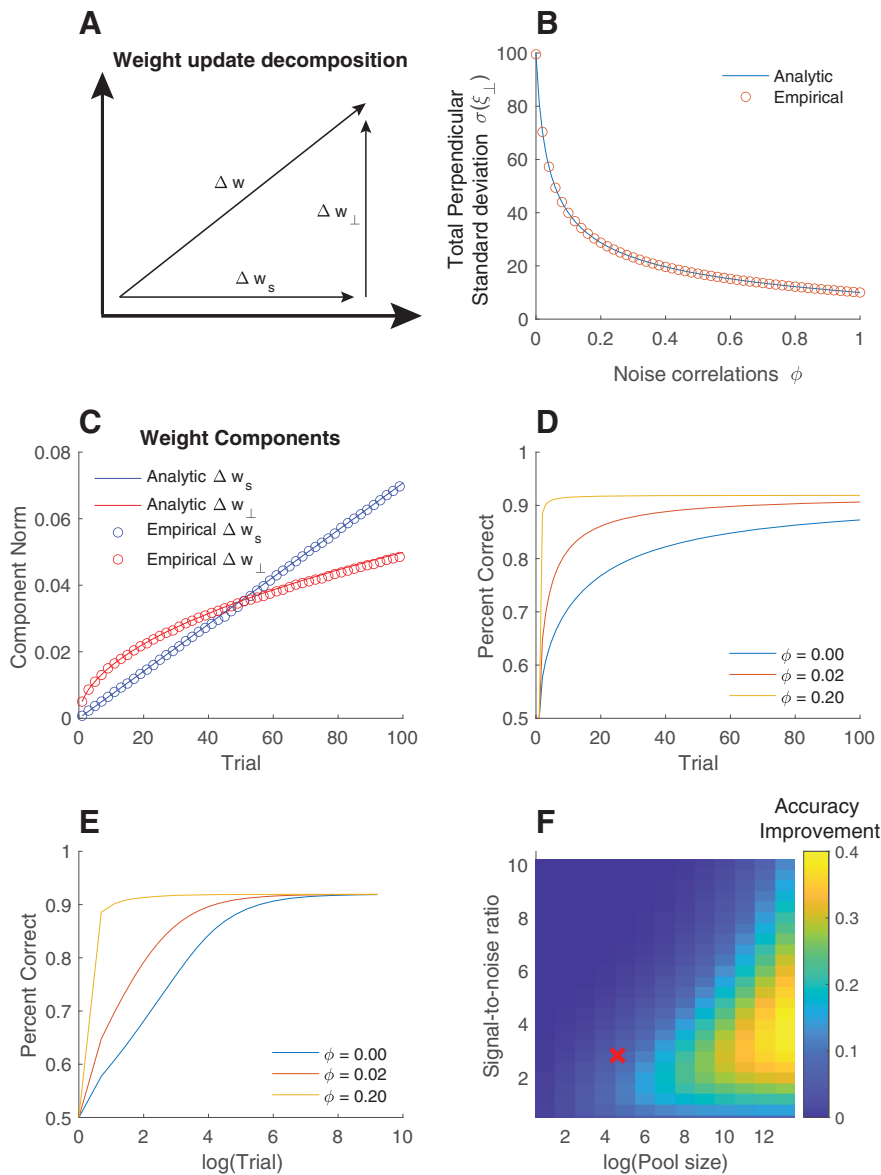
Given that noise correlations implemented in our previous simulations, like those observed in the brain, depended on the tuning of individual units, we tested whether such noise correlations might be produced via Hebbian plasticity. Specifically, we considered an extension of our neural network in which an additional intermediate layer is included between input and output neurons (Fig. 6A). Input units were again divided into two pools that differed in their encoding, but variability was uncorrelated across neurons within a given pool. Connections between the input layer and intermediate layer were initialized so that each input unit strongly activated one intermediate layer unit and were shaped over time using a Hebbian learning rule that strengthened connections between coactivated neuron pairs. Despite the lack of noise correlations in the input layer of this network [Fig. 6B; mean(std) in-pool residual correlation = 0.0015(0.10)], neurons in the intermediate layer developed tuning-specific noise correlations of the form that were beneficial for learning in the previous simulations [Fig. 6C; mean(std) in-pool residual correlation = 0.55(0.07);  $t$  test on difference from input layer correlations,  $t = 443$ ,  $df = 19,800$ ,  $p < 10e-50$ ]. Hebbian learning produced analogous noise correlation structure when initialized with random weights. The ability of Hebbian learning to reduce the dimensionality of the input units is consistent with previous theoretical work showing that it extracts the first principal component of the input vector, which in this case is the signal (Oja, 1982).



**Figure 3.** Correlated noise within similarly tuned populations leads to faster and more robust learning of a perceptual discrimination. **A**, A two-layer feedforward neural network was designed to solve a two-alternative forced-choice motion discrimination task at or near perceptual threshold. Input layer contains two homogenous pools of identically tuned neurons that provide evidence for alternate percepts (e.g., leftward motion vs rightward motion) and output neurons encode alternate courses of actions (e.g., saccade left vs saccade right). Layers are fully connected with weights randomized to small values and adjusted after each trial according to rewards (see above, Materials and Methods). **B**, Average learning curves for neural network models in which population signal-to-noise ratio in pools 1 and 2 were fixed, but noise correlations (gray-scale) were allowed to vary from small (dark) to large (light) values. **C**, **D**, Weight differences (left output–right output) for each input unit (color coded according to pool) after 100 timesteps of learning for low (**C**) and high (**D**) noise correlations. **E**, Accuracy in the last 20 training trials is plotted as a function of noise correlations for learned readouts (orange) and optimal readout (red). Lines/shading reflect mean/SEM. **F**, The shortest distance, in terms of neural activation, required to take the mean input for a given category (e.g., left or right) to the boundary that would result in misclassification is plotted for the final learned (orange) and optimal (red) weights for each noise correlation condition (abscissa). Lines/shading reflect mean/SEM.

### Dynamic, task-dependent noise correlations speed learning by constraining it to relevant feature dimensions

To understand how noise correlations might affect learning in mixed encoding populations, we extended our perceptual discrimination task to include two directions of motion discrimination (e.g., up/down and left/right). On each trial, a cue indicated



**Figure 4.** Analytic learning trajectories demonstrate advantage for noise correlations when pools are large and signal-to-noise ratio is low. **A**, Our analytical approach decomposed weight updates  $\Delta w$  into two components: updates in the signal dimension ( $\Delta w_s$ ) and updates perpendicular to the signal dimension ( $\Delta w_\perp$ ). **B**, SD of the variability in the dimension perpendicular to the signal (ordinate) decreased as a function of noise correlation (abscissa) as derived with our analytic approach (blue line; see above, Materials and Methods), and for the empirical simulations. **C**, For a given noise correlation (0.02 in this example) learning yielded weight changes in the signal dimension (blue circles) as well in the perpendicular dimension (red circles) that could be described analytically (blue and red lines). Circles represent average values from 20 empirical simulations. **D**, **E**, Theoretical accuracy derived from the analytical weights reproduces learning advantages observed in our simulations for higher levels of noise correlations (compare yellow to blue curves) and demonstrates convergence with sufficient observations (**E**; note abscissa in log units). **F**, Improvement in average accuracy over first 100 trials, derived analytically by taking the mean difference between yellow and blue curves in **D**, is indicated in color across a range of signal-to-noise ratios (ordinate) and neural population sizes (abscissa). The largest learning advantages for noise correlations were observed in large neural populations that contained limited stimulus information (moderately low SNR). Red X depicts parameters used for our simulations.

which of two possible motion discriminations should be performed (Fig. 7A, left; Cohen and Newsome, 2008). We extended our neural network to include four populations of 100 input units, each population encoding a conjunction of motion directions (up-right, up-left, down-right, down-left; Fig. 7A, input layer). Two additional inputs provided a perfectly reliable cue regarding the relevant feature for the trial (Fig. 7A, task units). Four output neurons encoded the four possible responses (up,

left, down, right) and were fully connected to the input layer (Fig. 7A, output layer). Task units were hard wired to eliminate irrelevant task responses, but weights of input units were learned over time as in our previous simulations.

Learning performance in the two-feature discrimination task depended not only on the level of noise correlations but also on the type. As in the previous simulation, adding noise correlations to each individual population of identically tuned units led to faster learning of the appropriate readout [Fig. 7B,C, compare blue and yellow; Fig. 7D,E, vertical axis; mean (std) accuracy across training: 0.54(0.05) and 0.70(0.05) for minimum (0) and maximum (0.2) in-pool correlations; *t* test for difference in accuracy, *t* = 226, *df* = 19,998, *p* < 10e-50].

However, the more complex task design also allowed us to test whether dynamic trial-to-trial correlations might further facilitate learning. Specifically, correlations that increase shared variability among units that contribute evidence to the same response have been observed previously (Cohen and Newsome, 2008) and could in principle focus learning on relevant dimensions (Fig. 2C,D), even when those dimensions change from trial to trial. Indeed, adding correlations among separate pools that share the same encoding of the relevant feature (e.g., up on a vertical trial) led to faster learning [Fig. 7B; mean (std) training accuracy for model with relevant pool correlations: 0.73(0.05); *t* test for difference from in-pool correlation only model, *t* = 34, *df* = 19,998, *p* < 10e-50] and weights that more closely approached the optimal readout (Fig. 7D, horizontal axis). In contrast, when positive noise correlations were introduced across separate encoding pools that shared the same tuning for the irrelevant dimension on each trial (e.g., up on a horizontal trial) learning was impaired dramatically [Fig. 7C; mean(std) training accuracy for model with irrelevant pool correlations, 0.51(0.05); *t* test for difference from in-pool correlation only model, *t* = -278, *df* = 19,998, *p* < 10e-50] and learned weights diverged from the optimal readout (Fig. 7E, horizontal axis). Model performance differences were completely attributable to learning the readout as all models performed similarly when using the optimal readout.

To test the idea that noise correlations might focus learning onto relevant dimensions, we extracted weight updates from each trial and projected these updates into a two-dimensional space where the first dimension captured the relative sensitivity to leftward versus rightward motion, and the second dimension captured relative sensitivity to upward versus downward motion.

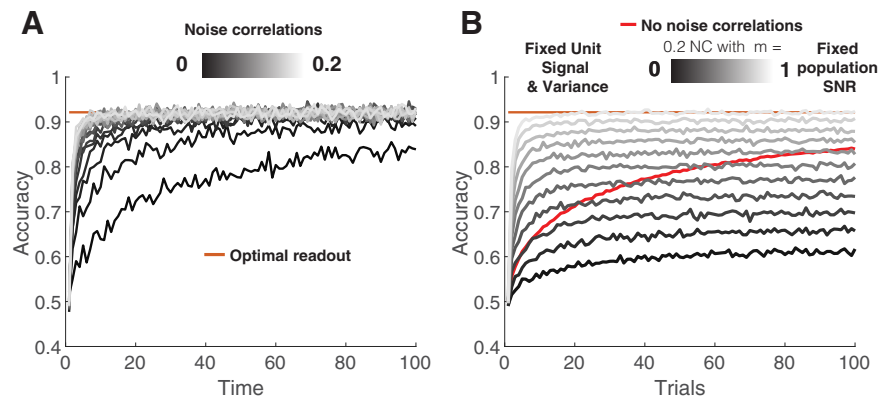


In the model where input units were only correlated within their identically tuned pool, weight updates projected in all directions more or less uniformly (Fig. 7G) and did not differ systematically across trial types (vertical vs horizontal). However, dynamic noise correlations that shared variability across the relevant dimension tended to push weight updates onto the appropriate dimension for a given trial [Fig. 7F; *t* test for difference in the magnitude of updating in up/down and left/right dimensions across conditions (up/down–left/right); *t* = 3.4, *df* = 98, *p* = 0.001]. In contrast, dynamic noise correlations that shared variability across the irrelevant dimension tended to push weight updates onto the wrong dimension [Fig. 7H; *t* test for difference in the magnitude of updating in up/down and left/right dimensions across conditions (up/down–left/right); *t* = –9.5, *df* = 98, *p* = 10e–14]. Both of these trends were consistent across simulations, providing an explanation for the performance improvements achieved by relevant noise correlations (projection of learning onto an appropriate dimension) and performance impairments produced by irrelevant noise correlations (projection of learning onto an inappropriate dimension).

## Discussion

Collectively, our results suggest that in settings where the population signal-to-noise ratio is externally limited and relevant task representations are low-dimensional, noise correlations can facilitate faster and more robust learning. We demonstrate this principle in a perceptual learning task (Fig. 3), where beneficial noise correlations emerged through simple Hebbian learning (Fig. 6). We extended our framework to a contextual learning task to demonstrate that dynamic noise correlations that bind task-relevant feature representations speed learning (Fig. 7B,D) by pushing learning onto task-relevant dimensions (Fig. 7F). Noise correlations among similarly tuned sensory neurons are pervasive (Zohary et al., 1994; Maynard et al., 1999; Bair et al., 2001; Averbeck and Lee, 2003; Cohen and Maunsell, 2009; Huang and Lisberger, 2009; Ecker et al., 2010; Gu et al., 2011; Adibi et al., 2013), and noise correlation dynamics that we show are beneficial for learning are observed *in vivo* (Cohen and Newsome, 2008). Therefore, we interpret our results as suggesting that noise correlations between similarly tuned neurons are a feature of neural coding architectures that ensures efficient readout learning rather than a bug that limits encoding potential.

This interpretation rests on several assumptions in our model. Of particular importance, is the assumption that the signal-to-noise ratio of our populations is fixed, meaning that our manipulation of noise correlations can focus variance on specific dimensions without gaining or losing information. This reflects conditions in which information is limited at the level of the inputs, for instance because of noisy peripheral sensors (Beck et al., 2012; Kanitscheider et al., 2015). In such conditions, even with optimal encoding, population information saturates at an upper bound determined by the information available in the inputs. Therefore, fixing the signal-to-noise ratio enabled us to examine noise correlation effects on readout learning in the absence of any

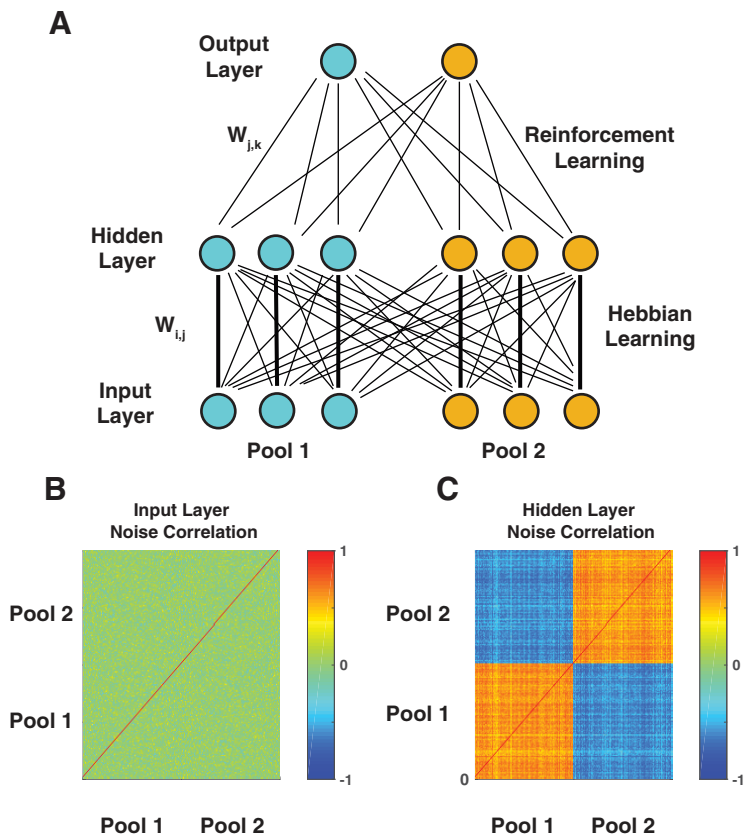


**Figure 5.** Impact of noise correlations is robust to single-unit variance but depends on assumptions about signal-to-noise ratio. *A*, Accuracy (ordinate) across trials (abscissa) for models in which signal-to-noise preserving noise correlations (gray-scale) were produced so that each unit maintains the same variance regardless of noise correlation magnitude (Eq. 25). Higher noise correlations (lighter colors) produced using this method also yielded faster learning. Orange line indicates accuracy of optimal readout. *B*, Accuracy (ordinate) as a function of trials (abscissa) for a model without noise correlations (red; equivalent to darkest line in *A*) and for several models that generate noise correlations (0.2) under different assumptions. The lightest color reflects a case where signal-to-noise ratio of the population is completely preserved, analogous to *A*. The darkest color reflects a case where the variance and signal of individual neurons is fixed, leading to a population signal-to-noise ratio that varies as a function of noise correlations. Intermediate colors indicate parametric mixtures of these assumptions created using Equation 26. Note that learning advantages depend critically on assumptions about signal-to-noise ratio and that noise correlations implemented using intermediate assumptions introduce a trade-off between faster learning (gray lines above red line for early trials) and lower asymptotic performance (gray lines below red line for later trials).

influence of noise correlations on the quantity of information contained in the population.

Previous theoretical work exploring the role of noise correlations in encoding has typically assumed that single neurons have a fixed variance and signal so that tilting the covariance of neural populations toward or away from the signal dimension would drastically affect the amount of information that can be encoded by a population (Fig. 1A; Averbeck et al., 2006; Moreno-Bote et al., 2014). Such assumptions lead to the idea that positive noise correlations among similarly tuned neurons limit encoding potential, raising the question of why they are so common in the brain (Cohen and Kohn, 2011). In considering the implications of this framework, one important question is the following: If information encoded by the population can be increased by changing the correlation structure among neurons, where does this additional information come from? In some cases, the neural population in question may indeed receive sufficient task-relevant information from upstream brain regions to reorganize encoding in this way, but in other cases information is likely limited by the inputs (Kanitscheider et al., 2015; Kohn et al., 2016). In cases where incoming information is limited, further increasing representational capacity is impossible, and formatting information for efficient readout is the best that the population code could do. Here we show that information-limiting noise correlations are exactly the type that format information most efficiently for readout under alternate assumptions. Between these two bookends of a fixed signal-to-noise ratio and fixed single-unit signal and variance, we also simulated intermediate regimes that do not perfectly preserve the signal-to-noise ratio. In these intermediate regimes, a trade-off emerges; noise correlations between similarly tuned neurons produce faster learning in the short term at the cost of lower asymptotic performance in the long run (Fig. 5B).

Jointly considering these perspectives on noise correlations provides a more nuanced view of how neural representations are likely optimized for learning. To optimize an objective function,



**Figure 6.** Hebbian learning produces correlations within similarly tuned populations in a perceptual discrimination task. **A**, Three-layer neural network architecture. Input layer feeds forward to hidden layer, which is fully connected to an output layer. Input layer provides uncorrelated inputs to hidden layer through projection weights that are adjusted according to a Hebbian learning rule. **B**, **C**, Noise correlations observed in hidden layer units at the beginning (**B**) and end (**C**) of training.

a neural population can reduce correlated noise in task-relevant dimensions to increase representational capacity up to some level constrained by its inputs (Fig. 8, left). But once all available task-relevant information is represented, populations can additionally optimize representations by pushing as much variance onto task-relevant dimensions as possible, thereby offering efficient downstream readout learning (Fig. 8, right). In short, the optimization of a neural population code depends critically on both upstream (e.g., input constraints) and downstream (e.g., readout) neural populations (Fig. 8). In this view, if a neural population is not fully representing the decision-relevant information made available to it, then learning could improve the efficiency of representations by reducing rate-limiting noise correlations as has been observed in some paradigms (Gu et al., 2011; Ni et al., 2018). In contrast, once available information is fully represented, readout learning could be further optimized by reformatting population codes so that variability is shared across neurons with similar tuning for the relevant task feature, producing the sorts of dynamic noise correlations observed in well-trained animals (Cohen and Newsome, 2008).

Beyond the form of noise correlations, our modeling included additional simplifying assumptions that are unlikely to hold up in real neural populations. For example, we consider pools of neurons identically tuned to discrete stimuli, rather than more realistic heterogeneous populations responding to continuous stimulus spaces. Previous work has shown that noise correlations do not necessarily limit encoding potential in heterogeneous populations with diverse tuning (Shamir and Sompolinsky, 2004, 2006; Chelaru and Dragoi,

2008; Ecker et al., 2011), and thus the degree to which the principles we reveal here will generalize to more realistic neural populations remains open. We hope that our results pave the way for future work employing mixed heterogeneous populations or more realistic architectures that go beyond the simple feedforward flow of information considered here.

### Model predictions

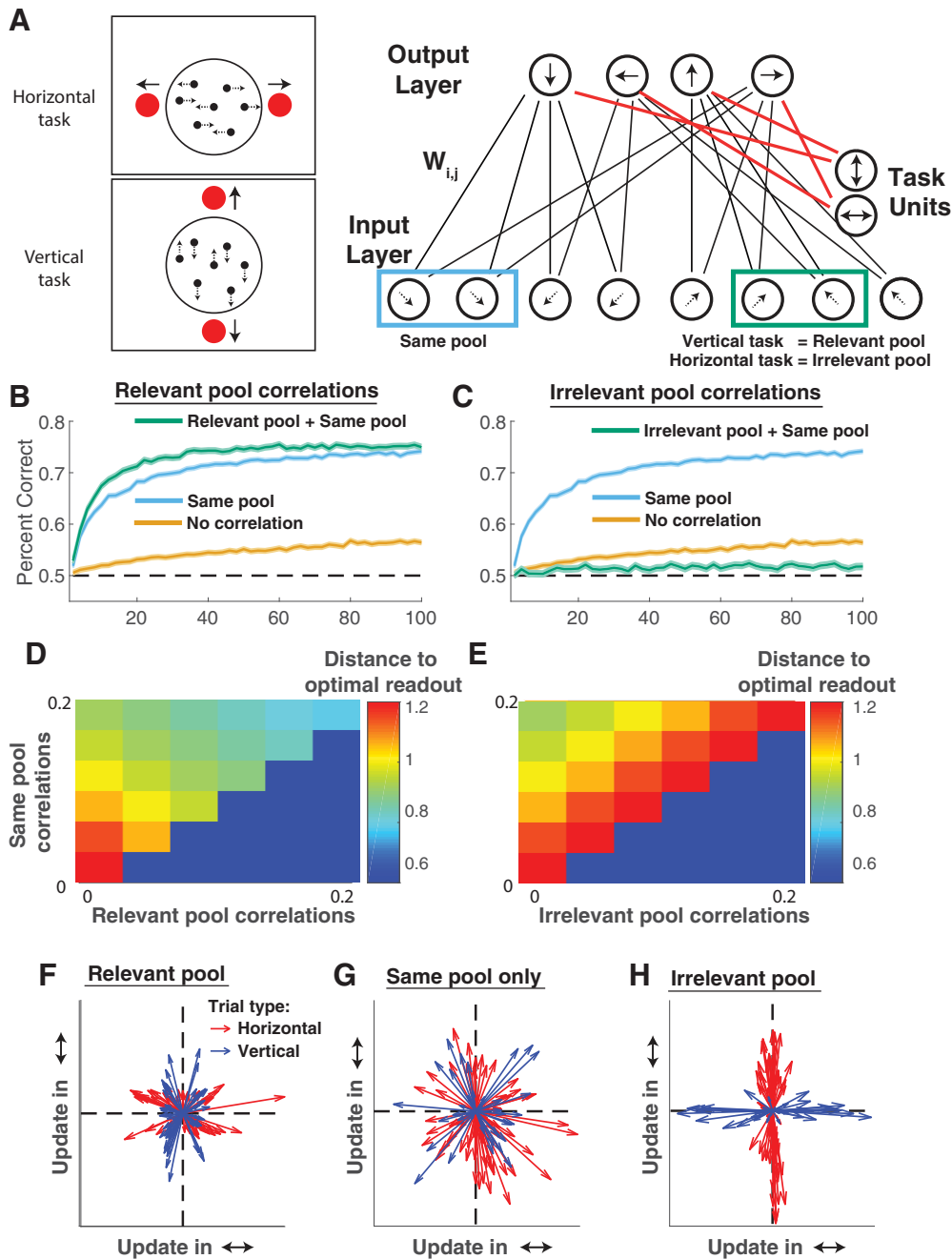
Our work shows that noise correlations can focus the gradient of learning onto the most appropriate dimensions. Thus, our model predicts that the degree to which similarly tuned neurons are correlated during a perceptual discrimination should be positively related to performance improvements experienced on subsequent discriminations. In contrast, our model predicts that the degree of correlation between neurons that are similarly tuned to a task-irrelevant feature should control the degree of learning on irrelevant dimensions and thus negatively relate to performance improvements on subsequent discriminations. These predictions are strongest for the earliest stages of learning where weight adjustments are critical for subsequent performance but may also hold for later stages of learning, when correlations on irrelevant dimensions, including independent noise channels, could potentially lead to systematic deviations from optimal readout (Figs. 2*F*, 4*D,E*). These predictions could be tested by recording neural responses to a stimulus set that differs across multiple features to characterize both signal-to-noise and correlated variability for each feature discrimination.

A strong prediction of our model is that correlated variability within neurons tuned to a given feature should be a predictor of subsequent learning of responses to that feature, above and beyond feature value discriminability.

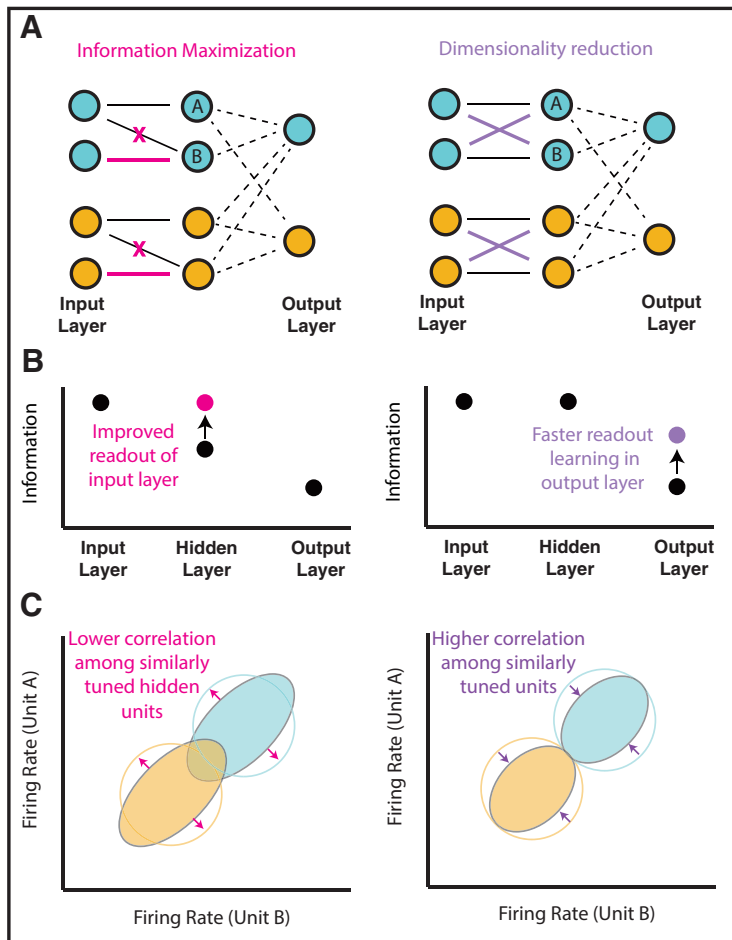
One interesting special case involves tasks where the relevant dimension changes in an un signaled manner (Birrell and Brown, 2000). In such tasks, noise correlations on the previously relevant dimension would, after such an extradimensional shift, force gradients into a task-irrelevant dimension and thus impair learning performance. Interestingly, learning after extradimensional shifts can be selectively improved by enhancing noradrenergic signaling (Devauges and Sara, 1990; Lapis and Morilak, 2006), which leads to increased arousal (Joshi et al., 2016; Reimer et al., 2016) and decreased pairwise noise correlations in the sensory and association cortex (Vinck et al., 2015; Joshi and Gold, 2020). Although these observations have been made in different paradigms, our model suggests that the reduction of noise correlations resulting from increased sustained levels of norepinephrine after an extradimensional shift (Bouret and Sara, 2005) could mediate faster learning by expanding the dimensionality of the learning gradients (compare Fig. 7*G* with 7*F*) to consider features that have not been task-relevant in the past.

### Origins of useful noise correlations

One important question stemming from our work is how noise correlations emerge in the brain. This question has been one of



**Figure 7.** Task-dependent noise correlations affect learning speed by projecting learning onto specific feature dimensions. **A**, A neural network was trained to perform two interleaved motion discrimination tasks (left; Cohen and Newsome, 2008). Network schematic (right) depicts two-layer feedforward network in which each homogenous pool of input units represents two dimensions of motion (up vs down, and left vs right), and output units produce responses in favor of alternative actions (up, down, left, right). Each homogenous pool of input units is identically tuned to one of four conjunctions of movement directions: up-left, down-left, up-right, down-right. Two additional input units provide cue information that biases output units to produce an output corresponding to the discrimination appropriate on this trial (e.g., horizontal or vertical). Noise correlations were manipulated among (1) identically tuned neurons (blue rectangle, same pool), (2) neurons that have similar encoding of the task relevant feature (green rectangle pair in vertical trials, relevant pool), and (3) neurons that have similar encoding of the task-irrelevant feature (green rectangle pair in horizontal trials, irrelevant pool). **B**, **C**, Learning curves showing accuracy (ordinate) over trials (abscissa) for models (1) lacking noise correlations (orange), (2) containing noise correlations that are limited to neurons that have same tuning for both features (same pool, blue), (3) containing same pool noise correlations along with correlations between neurons in different pools that have the same tuning for the task-relevant feature (in pool+rel pool, green in **B**), and (4) containing in-pool noise correlations along with correlations between neurons in different pools that have the same tuning for the task irrelevant feature (in pool+irrel pool, green in **C**). **D**, **E**, Distance between learned weights and the optimal readout (color) for models that differ in their level of in-pool correlations (ordinate, both plots), relevant-pool correlations (abscissa, **D**), and irrelevant-pool correlations (abscissa, **E**). **F**, **G**, **H**, Weight updates for example learning sessions were projected into a two-dimensional space in which net updates to the relative contribution of vertical motion information (e.g., up vs down) is represented on the abscissa, and updates to the relative contribution of horizontal motion information (e.g., left vs right) is represented on the ordinate. Arrows reflect single-trial weight updates and are colored according to the trial type (red = horizontal discrimination, blue = vertical discrimination). Weight updates for a model with only in-pool correlations look similar across trial types (**G**), but weight updates for a model with relevant-pool correlations indicate more weight updating on the relevant feature (**F**), whereas the opposite was observed in the case of irrelevant-pool correlations (**H**).



**Figure 8.** Information maximization and dimensionality reduction can be useful for learning under different situations and have opposite effects on noise correlations among similarly tuned units. **A**, A schematic representation of a three-layer neural network in which units provide evidence for one of two categorizations (blue/orange). In the left network, the hidden layer initially has access to information from only one of two independent units in each pool, but weights are subsequently adjusted to increase task-relevant information represented in the hidden layer (pink). In the right network, the hidden layer initially has access to all task-relevant information, but weights are subsequently adjusted to share signal and noise across similarly tuned units to provide dimensionality reduction (purple). Note that the information maximizing weight adjustments (left, pink) increase signal-to-noise ratio in the hidden layer but preserve the variance in the firing rate of individual neurons, whereas the dimensionality reducing weight adjustments (right, purple) maintain a fixed signal-to-noise ratio in hidden units but decrease the variance of individual units by averaging across multiple similarly tuned inputs. Dashed lines to output units reflect weights that need to be learned based on feedback. **B**, Task-relevant information (mutual information between unit activations and stimulus category, abscissa) is depicted for each layer (ordinate). Weight adjustments providing information maximization (left) increase task-relevant information in the hidden layer (pink), whereas weight adjustments that provide dimensionality reduction (right) do not affect task-relevant information in the hidden layer itself but instead increase the rate of learning in the output layer, thereby leading to more task-relevant information in the output layer (purple). **C**, Weight adjustments for information maximization (pink in **A**) decrease correlations among hidden units A and B by removing shared input from a single input unit and instead providing independent sources of input to each unit (pink arrows). In contrast, weight adjustments for dimensionality reduction increase noise correlations among hidden units A and B by providing them with the same mixture of information from the two identically tuned input units. We propose that both of these processes play a critical role in learning and that changes in noise correlations across learning will depend critically on which process dominates. As shown in **B**, this will depend critically on whether the neural population in question has already fully represented information available from its inputs. In principle, these processes could occur serially, with early learning maximizing information available in intermediate layers (left) and later learning compressing that information into a format allowing rapid readout learning (right).

long-standing debate, largely because there are so many potential mechanisms through which correlations could emerge (Kanitscheider et al., 2015; Kohn et al., 2016). Noise correlations could emerge from convergent and divergent feedforward wiring (Shadlen and Newsome, 1998), local connectivity patterns within a neural population (Hansen et

al., 2012; Smith et al., 2013), or top down inputs provided separately to different neural populations (Haefner et al., 2016). Here we show that static noise correlations that are useful for perceptual learning emerge naturally from Hebbian learning in a feedforward network. While this certainly suggests that useful noise correlations could emerge through feed forward wiring, it is also possible to consider our Hebbian learning as occurring in a one-step recurrence of the input units, and thus the same data support the possibility of noise correlations through local recurrence. The context dependent noise correlations that speed learning (Fig. 7), however, would not arise through simple Hebbian learning. Such correlations could potentially be produced through selective top-down signals from the choice neurons, as has been previously proposed (Wimmer et al., 2015; Haefner et al., 2016; Bondy et al., 2018; Lange et al., 2018). Moreover, top-down input may selectively target neuronal ensembles produced through Hebbian learning (Collins and Frank, 2013). Although previous work has suggested that such a mechanism could be adaptive for accumulating information over the course of a decision (Haefner et al., 2016), our work demonstrates that the same mechanism could effectively be used to tag relevant neurons for weight updating between trials, making efficient use of top-down circuitry. Haimerl et al. (2019) made a similar point, showing that stochastic modulatory signals shared across task-informative neurons can serve to tag them for a decoder.

### Noise correlations as inductive biases

Artificial intelligence (AI) has undergone a revolution over the past decade leading to human-level performance in a wide range of tasks (Mnih et al., 2015). However, a major issue for modern AI systems, which build heavily on neural network architectures, is that they require far more training examples than a biological system would (Hassabis et al., 2017). This biological advantage occurs despite the fact that the total number of synapses in the human brain, which could be thought of as the free parameters in our learning architecture, is much greater than the number of weights in even the most parameter-heavy deep learning architectures. Our work provides some insight into why this occurs; correlated variability across neurons in the brain constrain learning to specific dimensions, thereby limiting the effective complexity of the learning problem (Figs. 4A, 7F,G). We show that for simple tasks, this can be achieved using Hebbian learning rules (Fig. 6), but that contextual noise correlations, of the form that might be produced through top-down signals (Haefner et al., 2016), are critical for

appropriately focusing learning in more complex circumstances. In principle, algorithms that effectively learn and implement noise correlations might reduce the amount of data needed to train AI systems by limiting degrees of freedom to the most relevant dimensions. Furthermore, our work suggests that large-scale neural recordings in early stages of learning complex tasks might serve as indicators of the inductive biases that constrain learning in biological systems.

In summary, we show that under external constraints of task-relevant information, noise correlations that have previously been called rate limiting can serve an important role in constraining learning to task-relevant dimensions. In the context of the previous theory focusing on representation, our work suggests that neural populations are subject to competing forces when optimizing covariance structures; on the one hand reducing correlations between pairs of similarly tuned neurons can be helpful to fully represent available information, but increasing correlations among similarly tuned neurons can be helpful for assigning credit to task-relevant features. We believe that this view of the learning process not only provides insight to understanding the role of noise correlations in the brain but opens up the door to better understand the inductive biases that guide learning in biological systems.

## References

- Adibi M, McDonald JS, Clifford CWG, Arabzadeh E (2013) Adaptation improves neural coding efficiency despite increasing correlations in variability. *J Neurosci* 33:2108–2120.
- Averbeck BB, Lee D (2003) Neural noise and movement-related codes in the macaque supplementary motor area. *J Neurosci* 23:7630–7641.
- Averbeck BB, Latham PE, Pouget A (2006) Neural correlations, population coding and computation. *Nat Rev Neurosci* 7:358–366.
- Bair W, Zohary E, Newsome WT (2001) Correlated firing in macaque visual area MT: time scales and relationship to behavior. *J Neurosci* 21:1676–1697.
- Beck JM, Ma WJ, Pitkow X, Latham PE, Pouget A (2012) Non noisy, just wrong: the role of suboptimal inference in behavioral variability. *Neuron* 74:30–39.
- Birrell JM, Brown VJ (2000) Medial frontal cortex mediates perceptual attentional set shifting in the rat. *J Neurosci* 20:4320–4324.
- Bondy AG, Haefner RM, Cumming BG (2018) Feedback determines the structure of correlated variability in primary visual cortex. *Nat Neurosci* 21: 598–606.
- Bouret S, Sara SJ (2005) Network reset: a simplified overarching theory of locus coeruleus noradrenergic function. *Trends Neurosci* 28:574–582.
- Chelaru MI, Dragoi V (2008) Efficient coding in heterogeneous neuronal populations. *Proc Natl Acad Sci U S A* 105:16344–16349.
- Cohen MR, Newsome WT (2008) Context-dependent changes in functional circuitry in visual area MT. *Neuron* 60:162–173.
- Cohen MR, Maunsell JHR (2009) Attention improves performance primarily by reducing interneuronal correlations. *Nat Neurosci* 12:1594–1600.
- Cohen MR, Kohn A (2011) Measuring and interpreting neuronal correlations. *Nat Neurosci* 14:811–819.
- Collins AGE, Frank MJ (2013) Cognitive control over learning: creating, clustering, and generalizing task-set structure. *Psychol Rev* 120:190–229.
- Devauges V, Sara SJ (1990) Activation of the noradrenergic system facilitates an attentional shift in the rat. *Behav Brain Res* 39:19–28.
- Ecker AS, Berens P, Keliris GA, Bethge M, Logothetis NK, Tolias AS (2010) Decorrelated neuronal firing in cortical microcircuits. *Science* 327:584–587.
- Ecker AS, Berens P, Tolias AS, Bethge M (2011) The effect of noise correlations in populations of diversely tuned neurons. *J Neurosci* 31:14272–14283.
- Gu Y, Liu S, Fetsch CR, Yang Y, Fok S, Sunkara A, DeAngelis GC, Angelaki DE (2011) Perceptual learning reduces interneuronal correlations in macaque visual cortex. *Neuron* 71:750–761.
- Haefner RM, Berkes P, Fiser J (2016) Perceptual decision-making as probabilistic inference by neural sampling. *Neuron* 90:649–660.
- Haimerl C, Savin C, Simoncelli EP (2019) Flexible and accurate decoding of neural populations through stochastic comodulation. *Biorxiv* 624387.
- Hansen BJ, Chelaru MI, Dragoi V (2012) Correlated variability in laminar cortical circuits. *Neuron* 76:590–602.
- Hassabis D, Kumaran D, Summerfield C, Botvinick M (2017) Neuroscience-inspired artificial intelligence. *Neuron* 95:245–258.
- Hawkey DJC, Amitay S, Moore DR (2004) Early and rapid perceptual learning. *Nat Neurosci* 7:1055–1056.
- Hu Y, Zylberberg J, Shea-Brown E (2014) The sign rule and beyond: boundary effects, flexibility, and noise correlations in neural population codes. *PLoS Comput Biol* 10:e1003469.
- Huang X, Lisberger SG (2009) Noise correlations in cortical area MT and their potential impact on trial-by-trial variation in the direction and speed of smooth-pursuit eye movements. *J Neurophysiol* 101:3012–3030.
- Joshi S, Gold JJ (2020) Context-Dependent Relationships between Locus Coeruleus Firing Patterns and Coordinated Neural Activity in the Anterior Cingulate Cortex. *bioRxiv* 2020.09.26.314831.
- Joshi S, Li Y, Kalwani RM, Gold JJ (2016) Relationships between pupil diameter and neuronal activity in the locus coeruleus, colliculi, and cingulate cortex. *Neuron* 89:221–234.
- Kanitscheider I, Coen-Cagli R, Pouget A (2015) Origin of information-limiting noise correlations. *Proc Natl Acad Sci U S A* 112:E6973–E6982.
- Kohn A, Coen-Cagli R, Kanitscheider I, Pouget A (2016) Correlations and neuronal population information. *Annu Rev Neurosci* 39:237–256.
- Krotov D, Hopfield JJ (2019) Unsupervised learning by competing hidden units. *Proc Natl Acad Sci U S A* 116:7723–7731.
- Lange RD, Chatteraj A, Beck JM, Yates JL, Haefner RM (2018) A confirmation bias in perceptual decision-making due to hierarchical approximate inference. *Biorxiv* 440321.
- Lapiz MDS, Morilak DA (2006) Noradrenergic modulation of cognitive function in rat medial prefrontal cortex as measured by attentional set shifting capability. *Neuroscience* 137:1039–1049.
- Law C-T, Gold JJ (2009) Reinforcement learning can account for associative and perceptual learning on a visual-decision task. *Nat Neurosci* 12:655–663.
- Maynard EM, Hatsopoulos NG, Ojakangas CL, Acuna BD, Sanes JN, Normann RA, Donoghue JP (1999) Neuronal interactions improve cortical population coding of movement direction. *J Neurosci* 19:8083–8093.
- Mnih V, Kavukcuoglu K, Silver D, Rusu AA, Veness J, Bellemare MG, Graves A, Riedmiller M, Fidjeland AK, Ostrovski G, Petersen S, Beattie C, Sadik A, Antonoglou I, King H, Kumaran D, Wierstra D, Legg S, Hassabis D (2015) Human-level control through deep reinforcement learning. *Nature* 518:529–533.
- Moreno-Bote R, Beck J, Kanitscheider I, Pitkow X, Latham P, Pouget A (2014) Information-limiting correlations. *Nat Neurosci* 17:1410–1417.
- Ni AM, Ruff DA, Alberts JJ, Symmonds J, Cohen MR (2018) Learning and attention reveal a general relationship between population activity and behavior. *Science* 359:463–465.
- Oja E (1982) Simplified neuron model as a principal component analyzer. *J Math Biol* 15:267–273.
- Pouget A, Dayan P, Zemel R (2000) Information processing with population codes. *Nat Rev Neurosci* 1:125–132.
- Reimer J, McGinley MJ, Liu Y, Rodenkirch C, Wang Q, McCormick DA, Tolias AS (2016) Pupil fluctuations track rapid changes in adrenergic and cholinergic activity in cortex. *Nat Commun* 7:13289.
- Shadlen MN, Newsome WT (1998) The variable discharge of cortical neurons: implications for connectivity, computation, and information coding. *J Neurosci* 18:3870–3896.
- Shamir M, Sompolinsky H (2004) Nonlinear population codes. *Neural Comput* 16:1105–1136.
- Shamir M, Sompolinsky H (2006) Implications of neuronal diversity on population coding. *Neural Comput* 18:1951–1986.
- Smith MA, Jia X, Zandvakili A, Kohn A (2013) Laminar dependence of neuronal correlations in visual cortex. *J of Neurophysiology* 109:940–947.
- Stringer C, Michaelos M, Pachitariu M (2019) High precision coding in mouse visual cortex. *Biorxiv* 679324.
- Tsivids P, Pouncy T, Xu JL, Tenenbaum JB, Gershman SJ (2017) Human learning in Atari. Menlo Park, CA: AAAI SS-17-07.
- Vinck M, Batista-Brito R, Knoblich U, Cardin JA (2015) Arousal and locomotion make distinct contributions to cortical activity patterns and visual encoding. *Neuron* 86:740–754.
- Wimmer RD, Schmitt LI, Davidson TJ, Nakajima M, Deisseroth K, Halassa MM (2015) Thalamic control of sensory selection in divided attention. *Nature* 526:705–709.
- Zohary E, Shadlen MN, Newsome WT (1994) Correlated neuronal discharge rate and its implications for psychophysical performance. *Nature* 370:140–143.