1
2
3
4
5
6
7
8
9

10 **Noise correlations for faster and more robust learning**

11
12
13
14

15 Matthew R. Nassar[1,2], Daniel Scott[1,3], Apoorva Bhandari[1,3]

16
17
18
19
20
21

22 1. Robert J. & Nancy D. Carney Institute for Brain Science, Brown University, Providence RI
23    02912-1821, USA
24 2. Department of Neuroscience, Brown University, Providence RI 02912-1821, USA
25 3. Department of Cognitive, Linguistic, and Psychological Sciences, Providence RI, 02912-1821

26
27
28
29
30
31

39

40 *Competing interests:*

41

42 The authors have no financial or non-financial conflicts of interest related to this
43 work.

44
45

**Abstract:**

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 overall population signal-to-noise ratio in order 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 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.

**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

2

90  the population (Averbeck et al., 2006). In particular, positive noise correlations
91  between neurons that share the same stimulus tuning can reduce the amount of
92  decodable information in the neural population (Averbeck et al, 2006; Moreno-
93  Bote et al., 2014; Hu et al., 2014). Despite their detrimental effect on encoding,
94  noise correlations of this type are reliably observed, even after years of training
95  on perceptual tasks (Cohen and Kohn, 2011). Furthermore, noise correlations
96  between neurons are dynamically enhanced under conditions where two neurons
97  provide evidence for the same response in a perceptual categorization task
98  (Cohen and Newsome, 2008), raising questions about whether they might serve
99  a function rather than simply reflecting a suboptimal encoding strategy.
100
101  At the same time, learning to effectively read out a distributed code also poses a
102  significant challenge. Learning the appropriate weights for potentially tens of
103  thousands of neurons in a low signal-to-noise regime is a difficult, high-
104  dimensional problem, requiring a very large number of learning trials and
105  entailing considerable risk of "over fitting" to specific patterns of noise
106  encountered during learning trials. Nonetheless, people and animals can rapidly
107  learn to perform perceptual discrimination tasks, albeit with performance that
108  does not approach theoretically achievable levels (Hawkey et al., 2004; Stringer
109  et al., 2019). In comparison, deep neural networks capable of achieving human
110  level performance typically require a far greater number of learning trials than
111  would be required by humans and other animals (Tsividis et al., 2017). This
112  raises the question of how brains might implement inductive biases to enable
113  efficient learning in high dimensional spaces.
114
115  Here we address open questions about noise correlations and learning by
116  considering the possibility that noise correlations facilitate faster learning.
117  Specifically, we propose that noise correlations aligned to task relevant
118  dimensions could reduce the effective dimensionality of learning problems,
119  thereby making them easier to solve. For example, perceptual stimuli often
120  contain a large number of features that may be irrelevant to a given
121  categorization. At the level of a neural population, individual neurons may differ in
122  the degree to which they encode task irrelevant information, thus making the
123  learning problem more difficult. In principle, noise correlations in the relevant
124  dimension could reduce the effects of this variability on learned readout. Such an
125  explanation would be consistent with computational analyses of Hebbian learning
126  rules (Oja, 1982), which can both facilitate faster and more robust learning
127  (Krotov and Hopfield, 2019), and in turn may induce noise correlations. We
128  propose that faster learning of an approximate readout is made possible through
129  low dimensional representations that share both signal and noise across a large
130  neural population. In particular, we hypothesize that representations
131  characterized by enhanced noise correlations among similarly tuned neurons can
132  improve learning by focusing adjustments of the readout onto task relevant
133  dimensions.

134
135    We explore this possibility using neural network models of a two-alternative
136    forced choice perceptual discrimination task in which the correlation among
137    similarly tuned neurons can be manipulated independently of the overall
138    population signal-to-noise ratio. Within this framework, noise correlations, which
139    can be learned through Hebbian mechanisms, speed learning by forcing learned
140    weights to be similar across pools of similarly tuned neurons, thereby ensuring
141    learning occurs over the most task relevant dimension. We extend our framework
142    to a cued multidimensional discrimination task and show that dynamic noise
143    correlations similar to those observed in vivo (Cohen and Newsome, 2008),
144    speed learning by constraining weight updates to the relevant feature space. Our
145    results demonstrate that when information is extrinsically limited, noise
146    correlations can make learning faster and more robust by controlling the
147    dimensions over which learning occurs.
148
149

150    **Materials and Methods:**

151    Our goal was to understand the computational principles through which
152    correlations in the activity of similarly tuned neurons affect the speed with which
153    downstream neurons could learn an effective readout. Previous work has
154    demonstrated that manipulating noise correlations while maintaining a fixed
155    variance in the firing rates of individual neurons leads to changes in the
156    theoretical encoding capacity of a neural population (Averbeck et al., 2006;
157    Moreno-Bote et al., 2014). To minimize the potential impact of such encoding
158    differences, we took a different approach; rather than setting the variance of
159    individual neurons in our population to a fixed value, we set the signal-to-noise
160    ratio of our population to a fixed value. Thus, our approach does not ask how
161    maximum information can be packed into a given neural population's activity, but
162    rather how the strategy for packing a *fixed* amount of information in a population
163    affects the speed with which an appropriate readout of that information can be
164    learned. We implement this approach in a set of neural networks described in
165    more detail below.

166    *Learning readout in perceptual learning task*

167    Simulations and analyses for a simple perceptual discrimination task were
168    performed with a simplified and statistically tractable two-layer feed-forward
169    neural network (figure 3A). The input layer consisted of two homogenous pools of
170    100 units that were each identically "tuned" to one of two motion directions (left,
171    right). On each trial normalized firing rates for the neural population were drawn
172    from a multivariate normal distribution that was specified by a vector of stimulus-
173    dependent mean firing rates (signal: +1 for preferred stimulus, -1 for non-
174    preferred stimulus) and a covariance matrix. All elements of the covariance
175    matrix corresponding to covariance between units that were "tuned" to different

176 stimuli were set to zero. The key manipulation was to systematically vary the
177 magnitude of diagonal covariance components (eg. noise in the firing of
178 individual units) and the "same pool" covariance elements (eg. shared noise
179 across identically tuned neurons) while maintaining a fixed level of variance in
180 the summed population response for each pool:

181
$$\sigma_{pool}^2 = n\sigma_{unit}^2 + n(n-1)Cov(within\ pool) \quad Eq.1$$

182 Where $\sigma_{pool}^2$ is the variance on the sum of normalized firing rates from neurons
183 within a given pool, n is the number of units in the pool and the within pool
184 covariance ($Cov(within\ pool)$) specifies the covariance of pairs of units
185 belonging to the same pool. Signal-to-noise ratio (SNR) was defined as the
186 population signal (preferred-antipreferred) divided by the standard deviation of
187 the population response in the signal dimension. SNR was set to be 2 for each
188 individual pool of neurons, leading to a signal-to-noise ratio for the entire
189 population (both pools) equal to $2\sqrt{2}$. Given this constraint, the fraction of noise
190 that was shared across neurons within the same pool was manipulated as
191 follows:
192

193
$$\sigma_{unit}^2 = \frac{\sigma_{pool}^2}{n + n(n-1)\phi} \quad Eq.2$$

194
$$Cov(within\ pool) = \phi\sigma_{unit}^2 \quad Eq.3$$

195 Where $\phi$ reflects the fraction of noise that is correlated across units, which we
196 refer to in the text as noise correlations. Noise correlations ($\phi$) were manipulated
197 across values ranging from 0 to 0.2 for simulations. Note that, since $\phi$ appears in
198 the denominator of equation 2, adding noise correlations while sustaining a fixed
199 population signal-to-noise ratio leads to lower variance in the firing rates of single
200 neurons, differing from previous theoretical assumptions (compare figure 2a&b).
201

202 The input layer of the neural network was fully connected to an output layer
203 composed of two output units representing left and right responses. Output units
204 were activated on a given trial according to a weighted function of their inputs:
205

206
$$\boldsymbol{F}_{output} = \boldsymbol{w}\boldsymbol{F}_{input} \quad Eq.4$$

207

208 Where $F_{output}$ is a vector of firing rates of output units, $F_{input}$ is a vector of firing
209 rates of the input units, and w is the weight matrix. Firing of an individual output
210 unit can also be written as a weighted sum over input unit activity:
211

212
$$F_j = \sum_{i=1}^{200} w_{i,j} F_i \quad Eq.5$$

213

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

218

219
$$p(A_j) = \frac{e^{\beta F_j}}{\sum_k e^{\beta F_k}} \quad Eq.6$$

220　where $\beta$ is an inverse temperature, which was set to a relatively deterministic
221　value (10000). Learning was implemented through reinforcement of weights to
222　the selected output neuron (subscripted j below):

223

224
$$\Delta w_{i,j} = \alpha \delta F_i \quad Eq.7$$

225　Where $F_i$ is the normalized firing rate of the $i^{th}$ input neuron, $\delta$ is the reward
226　prediction error experienced on a given trial [+0.5 for correct trials and -0.5 for
227　error trials], and $\alpha$ is a learning rate (set to 0.0001 for simulations in figure 2). The
228　network was trained to correctly identify two stimuli (each of which was preferred
229　by a single pool of input neurons) over 100 trials (the last 20 trials of which were
230　considered testing). Simulations were repeated 1000 times for each level of $\phi$
231　and performance measures were averaged across all repetitions. Mean accuracy
232　per trial across all simulations was convolved with a Gaussian kernel (standard
233　deviation = 0.5 trials) for plotting in figure 2b. Mean accuracy across the final 20
234　trials was used as a measure of final accuracy (figure 2e). Statistics on model
235　performance were computed as Pearson correlations between noise correlations
236　$\phi$ and performance measures across all simulations and repetitions.

237　*Analytical learning trajectories*

238　One advantage of our simple network architecture is its mathematical tractability.
239　To complement the simulations described above, we also explored learning in
240　the network analytically. Specifically, we decomposed weight updates into two
241　categories: weight updates in the signal dimension, and weight updates
242　perpendicular to the signal dimension. Weight updates in the signal dimension
243　improved performance through alignment with the signal itself, whereas weight
244　updates in the perpendicular dimension limited performance through chance

6

245 alignment with trial-to-trial noise. An intuition for our approach and derivation are
246 provided below.
247
248 The two-alternative discrimination task is a one dimensional signal detection
249 problem, because it depends only on the difference between two scalars. In
250 particular, if $y = [y_1, y_2]$ denotes the readout activity the pair of pools and r
251 denotes the response (e.g. r=-1 is "respond left" and r=1 is "respond right"), then
252 $r = r(y_1 - y_2) = r(\Delta y)$. In addition, $\Delta y = w_1 x - w_2 x \equiv \Delta w x$, where x reflects
253 the firing rates of the input units and $w_1$ reflects the vector of weights mapping
254 input activation onto output unit 1 ($y_1$). To determine how accuracy is impacted
255 by noise correlations, we ask how Mahalanobis distance (d'), mean separation
256 (d), and signal variance ($\sigma_{s*}^2$) diverge over training time for the different noise
257 correlation conditions. The effective variance, $\sigma_{s*}^2$, differs from the true noise
258 variance in the signal dimension due to the fact that out-of-signal-dimension
259 noise is transferred into the signal dimension by imperfect readout weights.
260 Intuitively, learning speed may be improved by noise correlations because less
261 out-of-dimension noise is "learned into" the weights, thereby reducing the transfer
262 out-of-dimension noise into the signal space on any given trial.
263
264 The logic of training is as follows: On a correct trial, the weights to the chosen
265 unit are incremented by a multiple of the input vector x. That is:
266

$$w_i \rightarrow w_i + \alpha \delta x \quad Eq. 8$$

268
269 Here $\alpha$ reflects a positive learning rate, x reflects the activity of the input units and
270 $\delta$ is the reward prediction error, which we use as the absolute reward prediction
271 error instead of the signed one in this section for convenience.
272
273 Now the input is a sum of signal and zero mean noise:
274

$$x = \mu + \xi \quad Eq. 9$$

276
277 The expectation of noise is zero ($E(\xi) = 0$) and the signal $\mu$ can take only two
278 values $\mu \in \{\pm\mu_0\}$. Therefore if the weights start from some value $\Delta w(0)$, we will
279 find that:
280

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

282
283 Where t reflects the current timestep of learning. In words, we expect the amount
284 of signal in the weights to increase linearly over time. This means that we expect
285 the response to a noise-free signal ($\mu_0$) after t timesteps to be:
286

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

288

289  This is the measure d between the two Gaussian peaks in the one dimensional
290  signal detection problem described above. Below, we ignore the initial weight
291  term, since it does not change over time. To compute accuracy and d' over
292  training time we also need to compute the effective variance along the signal
293  dimension. First we note that the noise can be decomposed as:
294
295
$$\xi = \xi_s + \xi_\perp \quad Eq.\,12$$
296
297  where $\xi_s$ and $\xi_\perp$ are orthogonal components of the noise in the signal dimension
298  ($\xi_s$) and perpendicular to the signal dimension ($\xi_\perp$). Here we consider cases
299  where the noise along the signal dimension ($\xi_s$) has constant variance, following
300  on the assumptions that SNR is set to a constant value and that the mean signal
301  is the same for all noise correlation conditions.
302
303  The difference $\Delta y$ on any given trial decomposes into a sum of terms, one
304  reflecting weight-based transfer of signal and one reflecting the transfer of
305  orthogonal noise. This latter term arises because the weights are not, at any
306  finite time, a perfect matched filter for the signal. Letting subscripts s and $\perp$
307  continue to denote "signal" and "perpendicular" dimensions, we have:
308
309
$$\Delta y = \Delta w x \quad Eq.\,13$$
310
311
$$\Delta y = (\Delta w_s + \Delta w_\perp)(\mu + \xi_s + \xi_\perp) \quad Eq.\,14$$
312
313
$$\Delta y = \Delta w_s (\mu + \xi_s) + \Delta w_\perp \xi_\perp \quad Eq.\,15$$
314  where the final equation reflects the absence of terms that have zero products by
315  definition of the perpendicular subspaces. The variance of $\Delta y$ can be computed
316  using independence and orthogonality properties:
317
318
$$Var(\Delta y) = Var(\Delta w_s (\mu + \xi_s) + \Delta w_\perp \xi_\perp) \quad Eq.\,16$$
319
320
$$Var(\Delta y) = \Delta w_s^2 E[\xi_s^2] + \Delta w_\perp^2 E[\xi_\perp^2] \quad Eq.\,17$$
321
322
323  For any given network, the term $\Delta w(t)_\perp^2$ is a mean-zero diffusion process arising
324  from the fact that noise is added to the weights at every time step. For the
325  Gaussian white noise case, $\Delta w(t)_\perp^2$ is equivalent to Brownian motion in the (n-1)
326  dimensions perpendicular to the signal. Because (n-1) is not small, the summed
327  empirical variance of these processes, operative on each component, is likely to
328  be close to the theoretical total variance. If we split the term $\Delta w(t)_\perp^2$ into the (n-1)
329  components and index them with i, this gives:
330

8

331
$$\Delta w_\perp^i = \alpha\delta \sqrt{\frac{t}{n-1}}\, \sigma_\perp \quad Eq.\,18$$

332

333 The denominator of (n-1) appears here because Brownian motion determines
334 growth in the variance of each of the (n-1) perpendicular noise directions among
335 which the total variance $\sigma_\perp$ is distributed. Technically, our manipulation of the
336 noise covariance fixes the variance in a second direction of the space as well, so
337 that noise variance is actually evenly distributed over only (n-2) of the (n-1)
338 perpendicular dimensions, but this inhomogeneity is inconsequential if n is not
339 small; In effect, we are ignoring an order 1 term relative to an order n term for
340 simplicity. In order to understand how perpendicular weights grow with time, we
341 need only to determine $\sigma_\perp(\phi)$, where $\phi$ is the parameter controlling the noise
342 covariance matrix in our simulations. Specifically, the first row of the covariance
343 matrix takes the form:

344

345
$$\Sigma(\xi)_1 = [b, \phi b, \phi b, \dots \phi b, 0, \dots, 0] \quad Eq.\,19$$

346

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

349

350
$$b + \left(\frac{n}{2} - 1\right)\phi b = \sigma_s^2 \quad Eq.\,20$$

351

352
$$b = \frac{2\sigma_s^2}{2 + (n-2)\phi} \quad Eq.\,21$$

353

354 Since the eigenvalues of $\Sigma(\xi)$ are the variances in different dimensions of the
355 space, we can find the total variance perpendicular to the signal by subtracting
356 the known signal variance from the trace of $\Sigma(\xi)$:

357

358
$$\sigma_\perp^2 = Var(\xi_\perp) = Tr\big(\Sigma(\xi)\big) - Var(\xi_s) \quad Eq.\,22$$

359

360
$$\sigma_\perp^2 = Var(\xi_\perp) = nb - \sigma_s^2 \quad Eq.\,23$$

361

362 Putting this together with previous results, we have:

363

364
$$Var(\Delta y) = (t\alpha\delta\mu\sigma_s)^2 + \frac{t(\alpha\delta)^2\sigma_\perp^4}{n-1} \quad EQ\ 24$$

365

366 This provides analytic prediction for the variance of our readout decision variable
367 $\Delta y$ after learning for t trials, using a learning rate $\alpha$ to learn from from prediction
368 errors of magnitude $\delta$. Note that $\sigma_s$ was fixed in our simulations, but that $\sigma_\perp^4$
369 depends on $\phi$ through b, such that larger values of $\phi$ lead to smaller values of b,

9

370  and thus a smaller $\sigma_{\perp}^2$, reducing the second term in Eq 24. Furthermore, since the
371  first term in Eq 24 scales with $t^2$ its contributions dominate as more trials are
372  observed. This leads to identical asymptotic variance in the limit of large t, since
373  the first term does not depend on $\phi$.

374

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

379


380  *Noise correlations with fixed signal-to-noise ratio and single unit variance*

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

388

389
$$S_{neuron} = \sqrt{\frac{\sigma_{unit}^2\,(1 + (n-1)\phi)}{n}} \quad Eq.\,25$$

390

391  where $S_{neuron}$ reflects the amount of signal provided by each unit, $\sigma_{unit}^2$ reflects a
392  fixed variance assigned to each unit, n reflects the number of units in the pool,
393  and $\phi$ reflects the level of noise correlations. Thus, when we simulated correlated
394  noise using this equation, neurons maintained the same variance ($\sigma_{unit}^2$), but
395  increased their signal relative to the zero noise correlation condition ($\phi = 0$).

396

397  *Noise correlations that are bounded to a maximum signal-to-noise ratio*

398

399  In order to examine the importance of our assumption regarding fixed signal-to-
400  noise ratio, we also considered a parameterized model where signal ($S_{neuron}$)
401  was set according to a linear mixture:

402

403
$$S_{neuron} = m\sqrt{\frac{\sigma_{unit}^2\,(1 + (n-1)\phi)}{n}} + (1 - m)\sqrt{\frac{\sigma_{unit}^2}{n}} \quad Eq.\,26$$

404

405  where m is a mixing parameter that combines the signal producing a fixed signal-
406  to-noise ratio (first term) with a fixed signal that does not depend on the level of
407  noise correlations (second term). When m is set to 1, this parameterized model

408 obeys our assumptions regarding fixed signal to noise ratio, but when m is set to
409 0, the model conforms to more standard assumptions regarding fixed single unit
410 variance and signal.
411

412 *Hebbian learning of noise correlations in three-layer network*

413 We extended the two-layer feed-forward architecture described above to include
414 a third hidden layer in order to test whether Hebbian learning could facilitate
415 production of noise correlations among similarly tuned neurons (figure 5A). The
416 input layer was fully connected to the hidden layer, and each layer contained 200
417 neurons. In the input layer, neurons were homogenously tuned (100 leftward,
418 100 rightward) as described above, with $\phi$ set to zero (eg. no noise correlations).
419 Weights to the hidden layer were initialized to favor one-to-one connections
420 between input layer units and hidden layer units by adding a small normal
421 random weight perturbation (mean=0, standard deviation = 0.01) to an identity
422 matrix (though an alternate initialization was used to produce figure 6-1). During
423 learning, weights between the input and hidden layer were adjusted according to
424 a normalized Hebbian learning rule:
425

$$\Delta W = \alpha_{hebb} F'_1 F_2 \quad Eq. 27$$

427

428 Where $F'_1$ is a normalized vector of firing rates corresponding to the input layer
429 and $F_2$ is a normalized vector of firing rates corresponding to the hidden layer
430 units. The learning rate for Hebbian plasticity ($\alpha_{hebb}$) was set to 0.00005 for
431 simulations in figure 4 and 0.0005 for simulations in figure extended data figure
432 6-1. Weights were normalized after Hebbian learning to ensure that the
433 Euclidean norm of the incoming weights to each unit in layer two was equal to
434 one. The model was "trained" over 100 trials in the same perceptual
435 discrimination task described above and an additional 100 trials of the task were
436 completed to measure emergent noise correlations in the hidden layer. Noise
437 correlations were measured by regressing out variance attributable to the
438 stimulus on each trial, and then computing the Pearson correlation of residual
439 firing rate across each pair of neurons for the 100 testing trials (figure 4B&C).
440

441 *Learning readout in multiple discrimination task*

442 In order to test the impact of contextual noise correlations on learning (Cohen
443 and Newsome, 2008), the perceptual discrimination task was extended to include
444 two dimensions and two interleaved trial types: one in which an up/down
445 discrimination was performed (vertical), and one in which a right/left
446 discrimination was performed (horizontal). Each trial contained motion on the
447 vertical axis (up or down) and on the horizontal axis (left or right), but only one of
448 these motion axes was relevant on each trial as indicated by a cue.

449
450    In order to model this task, we extended our two-layer feed-forward network to
451    include 4 pools of input units, 4 output units, and 2 task units (figure 5A). Each
452    homogenous pool of 100 input units encoded a conjunction of the movement
453    directions (up-right, up-left, down-right, down-left). On each trial, the mean firing
454    rate of each input unit population was determined according to their tuning
455    preferences:
456
457
458
$$\mu = V + H \quad Eq. 28$$

459
460    where V was +1/-1 for trials with the preferred/anti-preferred vertical motion
461    direction H was +1/-1 for trials with the preferred/anti-preferred horizontal motion
462    direction. Firing rates for individual neurons were sampled from a multivariate
463    Gaussian distribution with mean $\mu$ and a covariance matrix that depended on trial
464    type (vertical versus horizontal) and the level of same pool, relevant pool, and
465    irrelevant pool correlations.
466
467    In order to create a covariance matrix, we stipulated a desired standard error of
468    the mean for summed population activity (SEM=20 for simulations in figure 5)
469    and determined the summed population variance that would correspond to that
470    value ($\sigma^2_{pool}$). We then determined the variance on individual neurons that would
471    yield this population response under a given noise correlation profile as follows:
472

473
$$\sigma^2_{unit} = \frac{\sigma^2_{pool}}{n + n(n-1)\phi_{same} + n^2\phi_{relevant} - n^2\phi_{irrelevant}} \quad Eq. 29$$

474
475    where $\phi_{same}$ is the level of same pool correlations (range: 0-0.2 in our
476    simulations), $\phi_{relevant}$ is the level of relevant pool correlations (range: 0-0.2 in our
477    simulations), $\phi_{irrelevant}$ is the level of irrelevant pool correlations (range: 0-0.2 in
478    our simulations. Note that increasing the same pool or in pool correlations
479    reduces the overall variance in order to preserve the same level of variance on
480    the task relevant dimension in the population response, but that increasing
481    irrelevant pool correlations has the opposite effect. Covariance elements of the
482    covariance matrix were determined as follows:
483

484
$$Cov(same\ pool) = \phi_{same}\sigma^2_{unit} \quad Eq. 30$$
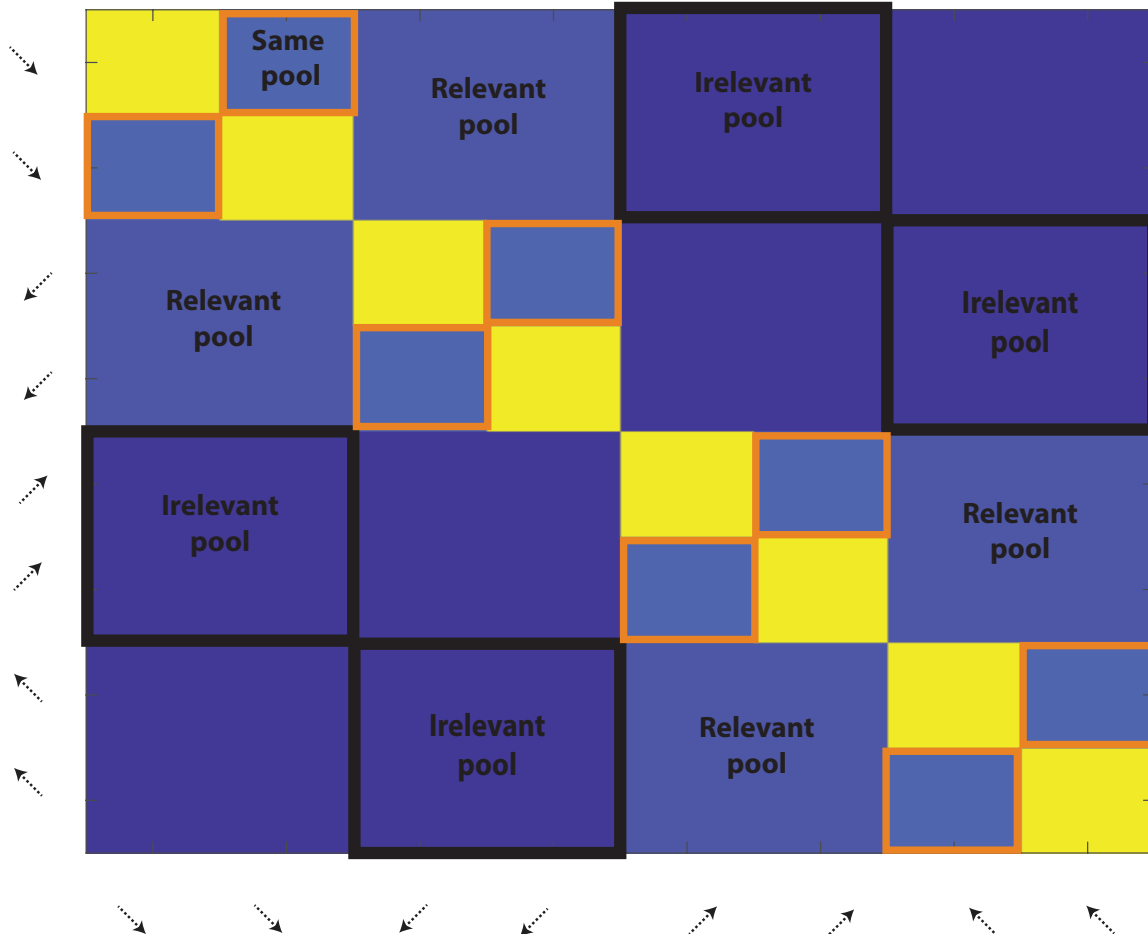
485
$$Cov(relevant\ pool) = \phi_{relevant}\sigma^2_{unit} \quad Eq. 31$$

486
$$Cov(irrelevant\ pool) = \phi_{irrelevant}\sigma^2_{unit} \quad Eq. 32$$

487 Variance and covariance values above were used to construct a covariance
488 matrix for each trial type (vertical/horizontal) as depicted in figure 1.
489

**Covariance matrix: vertical trials**



490
491 **Figure 1: Schematic of covariance matrix for two-dimensional motion discrimination task.**
492 The covariance between units with different motion tuning (reflected by the arrows labeling
493 columns and rows) is schematically represented for a simplified input layer, where only two
494 identically tuned neurons are in each pool (in actual simulations there were 100 units per pool).
495 Same pool correlations are controlled by covariance elements between neurons with identical
496 tuning (orange boxes). Relevant pool correlations are controlled by covariance elements between
497 neurons that are similarly tuned to the task-relevant feature. Task irrelevant correlations are
498 controlled by covariance elements between neurons that are similarly tuned to the task-irrelevant
499 feature. The covariance matrix shown here is for a vertical trial – on a horizontal trial the irrelevant
500 pool and relevant pool locations would be reversed. Covariance elements for pairs of neurons
501 that differed in tuning on both dimensions were set to zero. Each input population has been
502 depicted as two units here for presentation purposes. Background color reflects the case where
503 same pool correlations = 0.2 and relevant pool correlations = 0.1.
504
505
506 Output units corresponded to the four possible task responses (up, down, left,
507 right) and were activated according to a weighted sum of their inputs as

508    described previously. Task units were modeled as containing perfect information
509    about the task cue (vertical versus horizontal) and each task unit projected with
510    strong fixed weights (1000) to both responses that were appropriate for that task.
511    Decisions were made on each trial by selecting the output unit with the highest
512    activity level. Weights to chosen output unit were updated using the same
513    reinforcement learning procedure described in the two alternative perceptual
514    learning task.
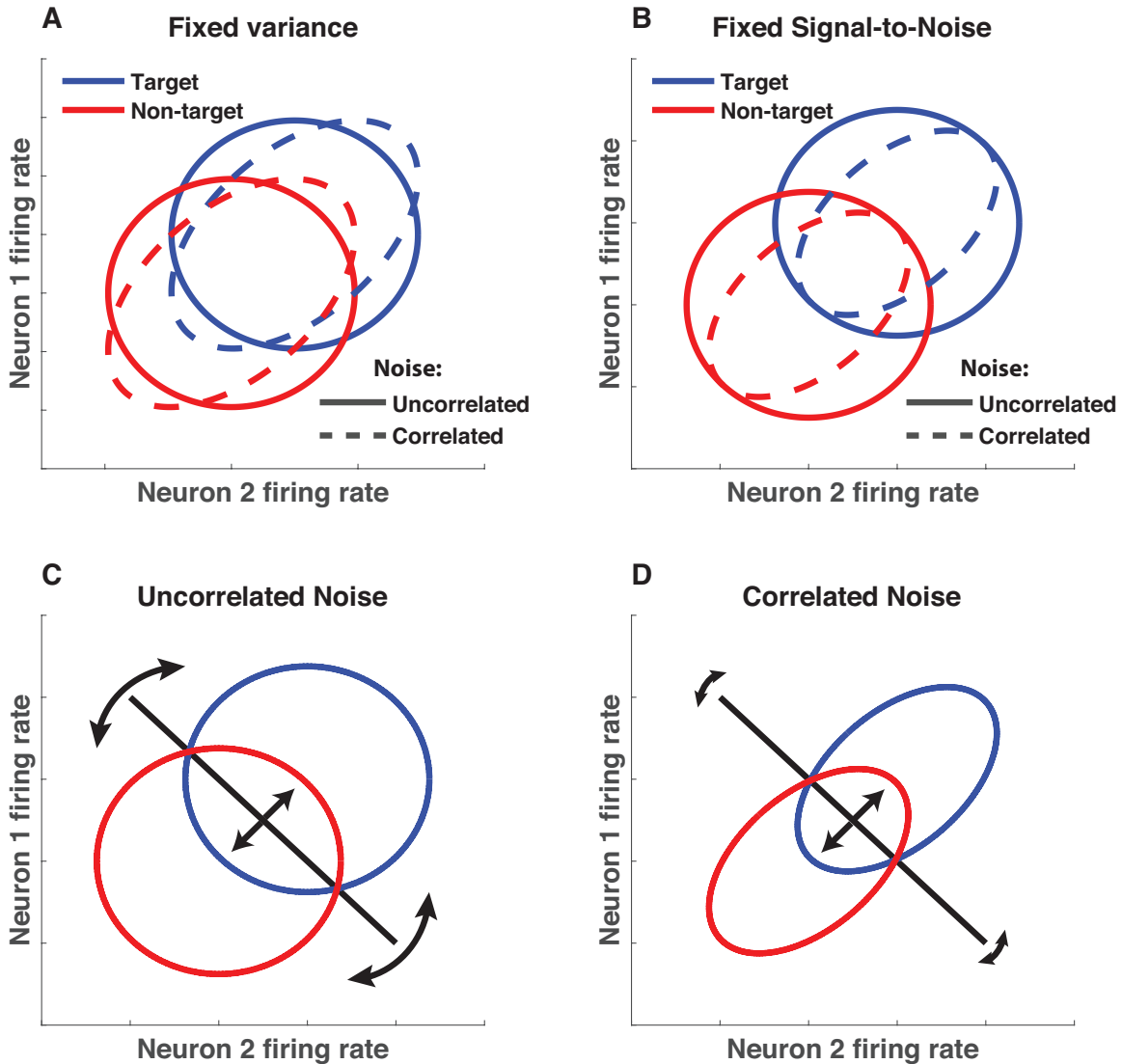515
516
517
518    **Results:**
519
520    We examine how noise correlations affect learning in a simplified neural network
521    where the appropriate readout of hundreds of weakly tuned units is learned over
522    time through reinforcement. In order to isolate the effects of noise correlations on
523    learning, rather than their effects on other factors such as representational
524    capacity, we consider population encoding schemes at the input layer that can be
525    constrained to a fixed signal-to-noise ratio. This assumption differs from previous
526    work on noise correlations where the *variance* of the neural population is
527    assumed to be fixed and covariance is changed to produce noise correlations,
528    thereby affecting the representational capacity of the population (figure 2A;
529    (Averbeck et al., 2006; Moreno-Bote et al., 2014)). Under our assumptions, a
530    fixed signal-to-noise ratio can be achieved for any level of noise correlations by
531    scaling the variance (figure 2B; equations 1-3), or, alternately scaling the
532    magnitude of the signal (equation 25). While we do not discount the degree to
533    which noise correlations affect the encoding potential of neural populations, we
534    believe that in many cases the relevant information is limited by extrinsic factors
535    (eg. the stimulus itself, or upstream neural populations providing input (Ecker et
536    al., 2011; Beck et al., 2012; Kanitscheider et al., 2015)). Under such conditions,
537    reducing noise correlations can increase information only until it saturates
538    because all of the available incoming information is encoded. Beyond that,
539    increasing encoding potential is not possible as it would be tantamount to the
540    population "creating new information" that was not communicated by inputs to the
541    population. Therefore, our framework can be thought of as testing how best to
542    format limited available information in a neural population in order to ensure that
543    an acceptable readout can be rapidly and robustly learned.
544
545    We propose that within this framework, noise correlations of the form that have
546    previously been shown to limit encoding are beneficial because they constrain
547    learning to occur over the most relevant dimensions. In general, a linear readout
548    can be thought of as hyperplane serving as a classification boundary in an N
549    dimensional space, where N reflects the number of neurons in a population.
550    Learning in such a framework involves adjustments of the hyperplane to
551    minimize classification errors. The most useful adjustments are in the dimension

14

552 that best discriminates signal from noise (central arrows in figure 2C&D), but
553 adjustments may also occur in dimensions orthogonal to the relevant one (such
554 as "twisting" of the hyperplane depicted by curved arrows in figure 2C&D) that
555 could potentially impair performance, or slow down learning. Our motivating
556 hypothesis is that by focusing population activity into the task relevant dimension,
557 noise correlations can increase the fraction of hyperplane adjustments that occur
558 in the task relevant dimension (figure 2D), thus reducing the effective
559 dimensionality of readout learning.

560



561
562
563
564 **Figure 2: Modeling noise correlations with extrinsic constraint on signal-to-noise ratio. A)**
565 Previous work has modeled noise correlations by assuming that population variance is fixed and
566 that covariance is manipulated to produce noise correlations. Under such assumptions, the firing
567 rate of two similarly tuned neurons is plotted in the absence (solid) or presence (dotted) of
568 information-limiting noise correlations. **B)** Here we assume that the signal-to-noise ratio of the
569 neural population is limited to a fixed value such that noise correlations between similarly tuned

15

570 neurons do not affect theoretical performance. Thus, the percent overlap of blue (target) and red
571 (non-target) activity profiles does not differ in the presence (dotted) or absence (solid) of noise
572 correlations. **C&D)** Under this assumption, noise correlations among similarly tuned neurons
573 could compress the population activity to a plane orthogonal to the optimal decision boundary,
574 thereby minimizing boundary adjustments in irrelevant dimensions (**C**) and maximizing boundary
575 adjustments on relevant ones (**D**).
576
577
578 *Noise correlations enable faster learning in a fixed signal-to-noise regime*
579
580 In order to test our overarching hypothesis, we constructed a fully connected two-
581 layer feed-forward neural network in which input layer units responded to one of
582 two stimulus categories (pool 1 & pool 2) and each output unit produced a
583 response consistent with a category perception (left/right units in figure 3A). On
584 each trial, the network was presented with one stimulus at random, and input
585 firing for each pool was drawn from a multivariate Gaussian with a covariance
586 that was manipulated while preserving the population signal-to-noise ratio.
587 Output units were activated according to a weighted average of inputs and a
588 response was selected according to output unit activations. On each trial,
589 weights to the selected action were adjusted according to a reinforcement
590 learning rule that strengthened connections that facilitated a rewarded action and
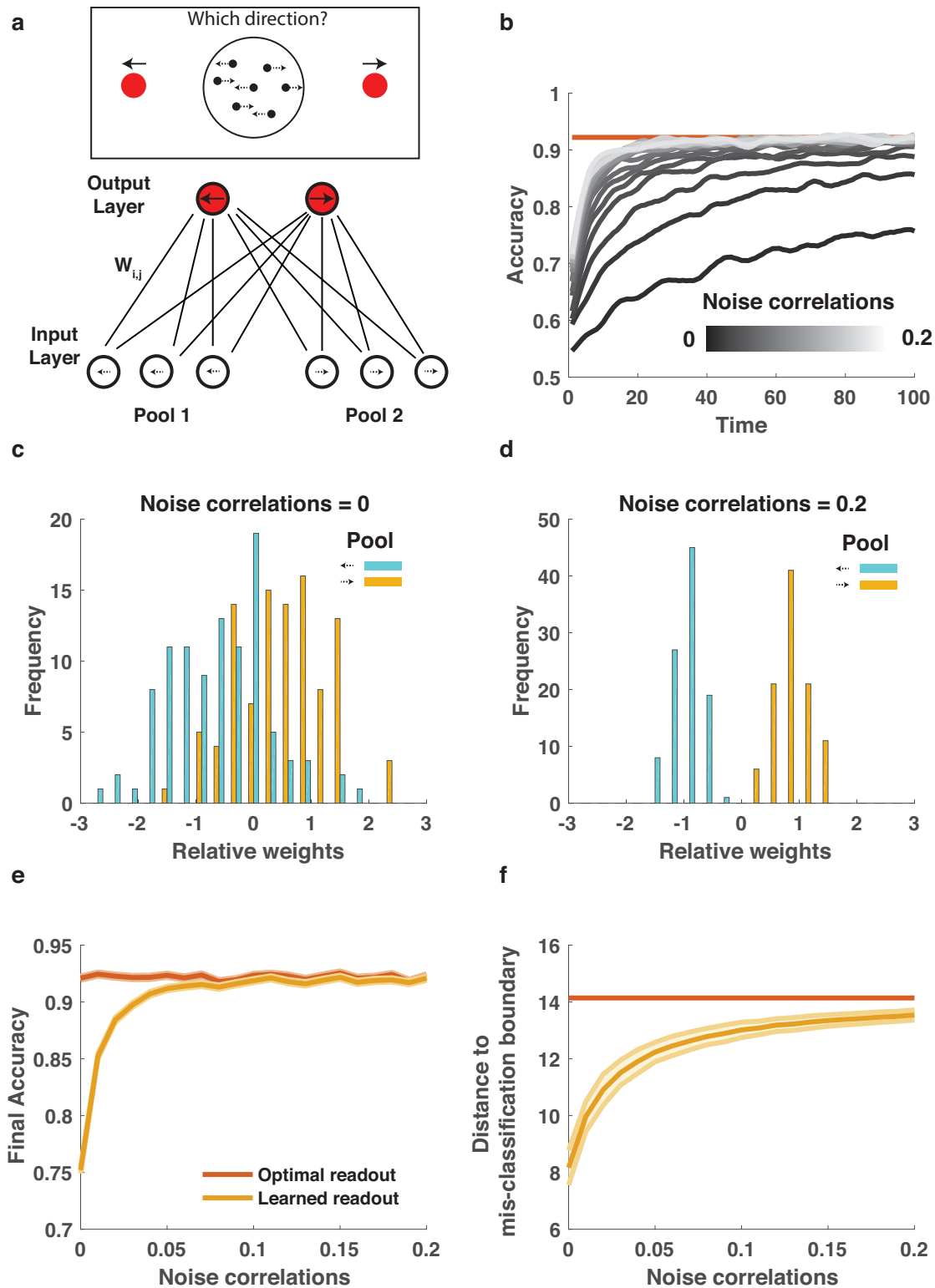591 weakened connections that facilitated an unrewarded action (Law and Gold,
592 2009).
593
594 Noise correlations led to faster and more robust learning of the appropriate
595 stimulus-response mapping. All neural networks learned to perform the requisite
596 discrimination, but neural networks that employed correlations among similarly
597 tuned neurons learned more rapidly (figure 3B). After learning, networks that
598 employed such noise correlations assigned more homogenous weights to input
599 units of a given pool than did networks that lacked noise correlations (compare
600 figure 3C&D). This led to better trained-task performance (figure 3E; Pearson
601 correlation between noise correlations and test performance: R = 0.29, p < 10e-
602 50) and greater robustness to adversarial noise profiles (figure 3F; R = 0.81, p <
603 10e-50) in the networks that employed noise correlations. Critically, these
604 learning advantages emerged despite the fact that optimal readout of all
605 networks achieved similar levels of performance and robustness (figure 3E&F,
606 compare optimal readout across conditions).
607
608
609

**Figure 3: Correlated noise within similarly tuned populations leads to faster and more robust learning of a perceptual discrimination. A)** A two-layer feed-forward neural network was designed to solve a two alternative forced choice motion discrimination task at or near

17

616   perceptual threshold. Input layer contains two homogenous pools of identically tuned neurons
617   that provide evidence for alternate percepts (eg. leftward motion versus rightward motion) and
618   output neurons encode alternate courses of actions (eg. saccade left versus saccade right).
619   Layers are fully connected with weights randomized to small values and adjusted after each trial
620   according to rewards (see methods). **B)** Average learning curves for neural network models in
621   which population signal-to-noise ratio in pools 1 and 2 were fixed, but noise correlations
622   (grayscale) were allowed to vary from small (dark) to large (light) values. **C&D)** Weight
623   differences (left output – right output) for each input unit (color coded according to pool) after 100
624   timesteps of learning for low (**C**) and high (**D**) noise correlations. **E**) Accuracy in the last 20
625   training trials is plotted as a function of noise correlations for learned readouts (orange) and
626   optimal readout (red). Lines/shading reflect Mean/SEM. F) The shortest distance, in terms of
627   neural activation, required to take the mean input for a given category (eg. left or right) to the
628   boundary that would result in misclassification is plotted for the final learned (orange) and optimal
629   (red) weights for each noise correlation condition (abscissa). Lines/shading reflect Mean/SEM.
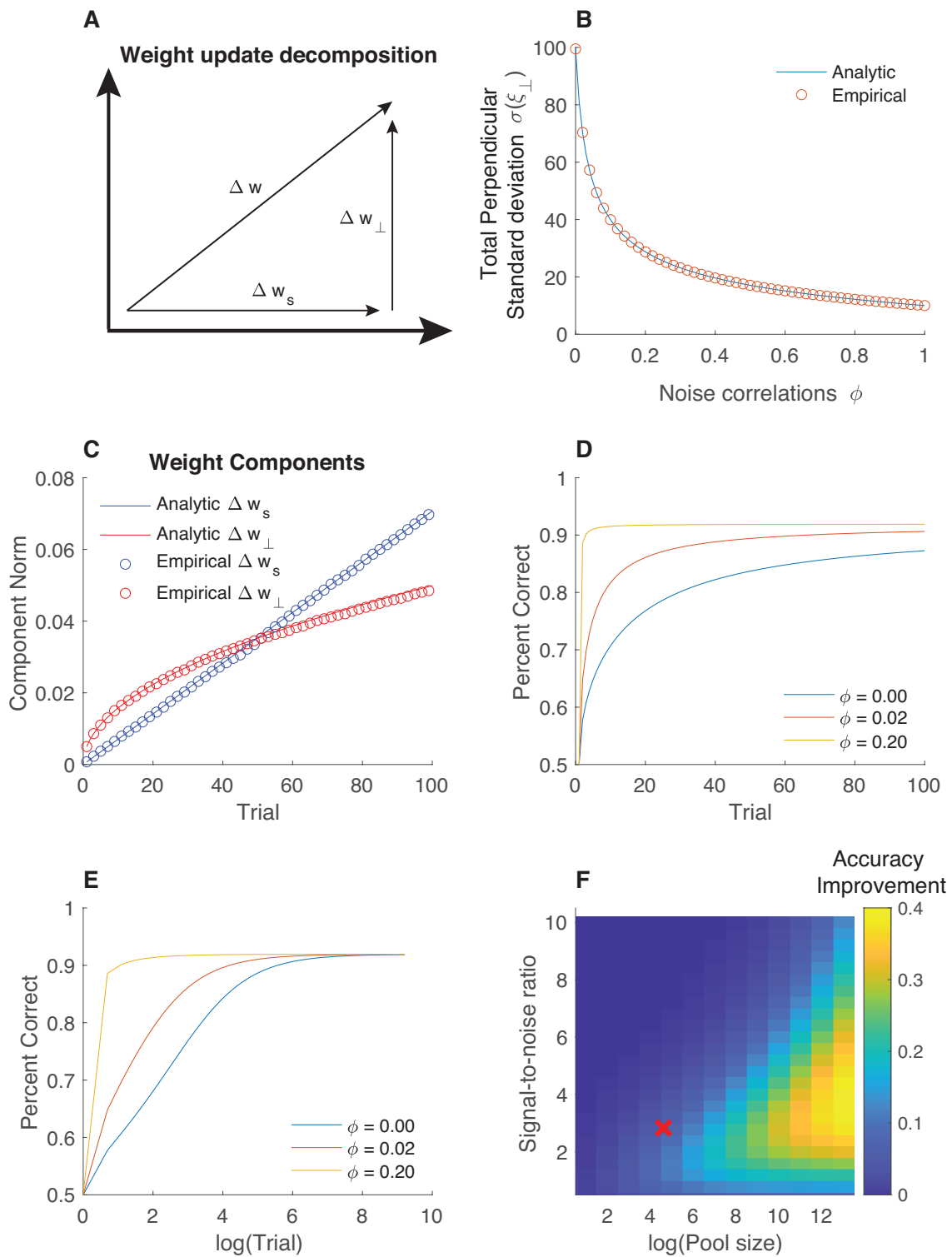630
631
632   *Learning benefits from noise correlations are greatest for large, low SNR populations.*
633
634   In order to better understand how noise correlations promoted faster learning we
635   developed an analytical method for describing learning trajectories (see
636   methods). Our method considered the impacts of two influences on weight
637   updates over time: 1) weight updates in the signal dimension that tend to align
638   with the signal and improve performance and 2) weight updates perpendicular to
639   the signal dimension, which through chance alignment with trial-to-trial firing rate
640   variability allow noise to impact decisions and therefore and hinder performance
641   (fig 4a). Noise correlations implemented using our methods decreased the latter
642   form of weight updates (fig 4b), leading updates in the signal dimension to more
643   quickly dominate performance (fig 4c), thereby speeding analytical predictions for
644   learning (fig 4d&e). The analytically derived learning advantage for fixed-SNR
645   noise correlations was greatest for situations in which SNR was relatively low
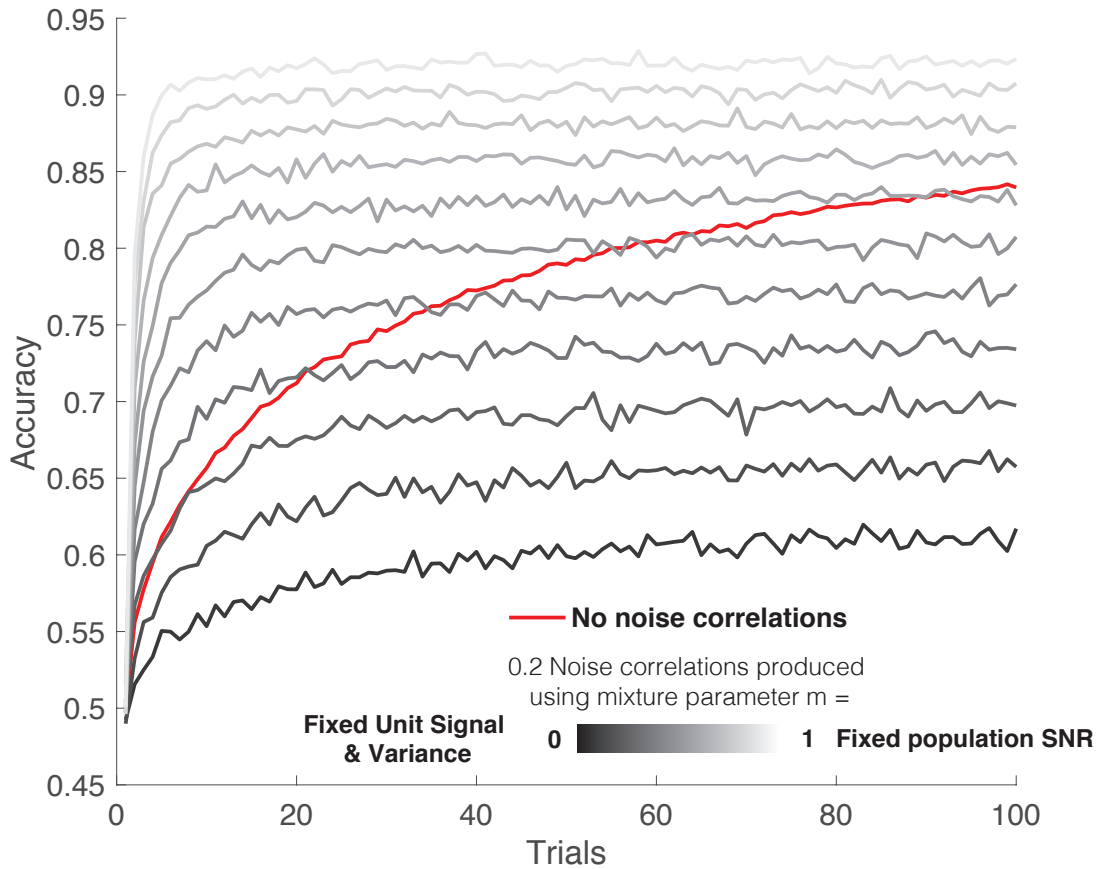646   and neural populations were large (fig 4f).
647

**A** Weight update decomposition

**B**

**C** Weight Components

**D**

**E**

**F** Accuracy Improvement

648
649

**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)** Standard deviation 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 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 twenty 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.

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, in order to maintain SNR for higher levels of noise correlation (equation 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, similarly to in our previous simulations (Extended data Fig 3-1).

*Noise correlations that do not maintain signal-to-noise ratio can introduce a tradeoff 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. In order 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 signal and single unit variance, analogous to assumptions of Averbeck and colleagues (Averbeck et al., 2006) (m=0). Consistent with our previous results, noise correlations improve learning in the m=1 case, and consistent with Averbeck 2006, asymptotic performance is reduced by noise correlations in the m=0 case (Fig 5). 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 tradeoff between short-term gains from rapid learning and long term gains from higher asymptotic performance.

698



**Figure 5: Impact of noise correlations on learning depends on the assumption that signal-to-noise ratio is fixed.** Accuracy (ordinate) as a function of trials (abscissa) for a model without noise correlations (red) 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 our previous simulations. The darkest color

706     reflects a case where the variance and signal of individual neurons is fixed, leading to a
707     population signal-to-noise ratio that varies as a function of noise correlations. Intermediate colors
708     indicate parametric mixtures of these assumptions created using equation 26. Note that learning
709     advantages depend critically on assumptions about signal-to-noise ratio, and that noise
710     correlations implemented using intermediate assumptions introduce a tradeoff between faster
711     learning (gray lines above red line for early trials) and lower asymptotic performance (gray lines
712     below red line for later trials).
713
714
715     _Hebbian learning can produce useful noise correlation structure._
716
717     Given that noise correlations implemented in our previous simulations, like those
718     observed in the brain, depended on the tuning of individual units, we tested
719     whether such noise correlations might be produced via Hebbian plasticity.
720     Specifically, we considered an extension of our neural network in which an
721     additional intermediate layer is included between input and output neurons
722     (figure 6a). Input units were again divided into two pools that differed in their
723     encoding, but variability was uncorrelated across neurons within a given pool.
724     Connections between the input layer and intermediate layer were initialized such
725     that each input unit strongly activated one intermediate layer unit, and shaped
726     over time using a Hebbian learning rule that strengthened connections between
727     co-activated neuron pairs. Despite the lack of noise correlations in the input layer
728     of this network (figure 6b; mean[std] in pool residual correlation = 0.0015[0.10]),
729     neurons in the intermediate layer developed tuning-specific noise correlations of
730     the form that were beneficial for learning in the previous simulations (figure 6c;
731     mean[std] in pool residual correlation = 0.55[0.07]; $t$-test on difference from input
732     layer correlations $t = 443$, $dof = 19800$, $p < 10e\text{-}50$). Hebbian learning produced
733     analogous noise correlation structure when initialized with random weights
734     (Extended data figure 6-1). The ability of Hebbian learning to reduce the
735     dimensionality of the input units is consistent with previous theoretical work
736     showing that it extracts the first principal component of the input vector, which in
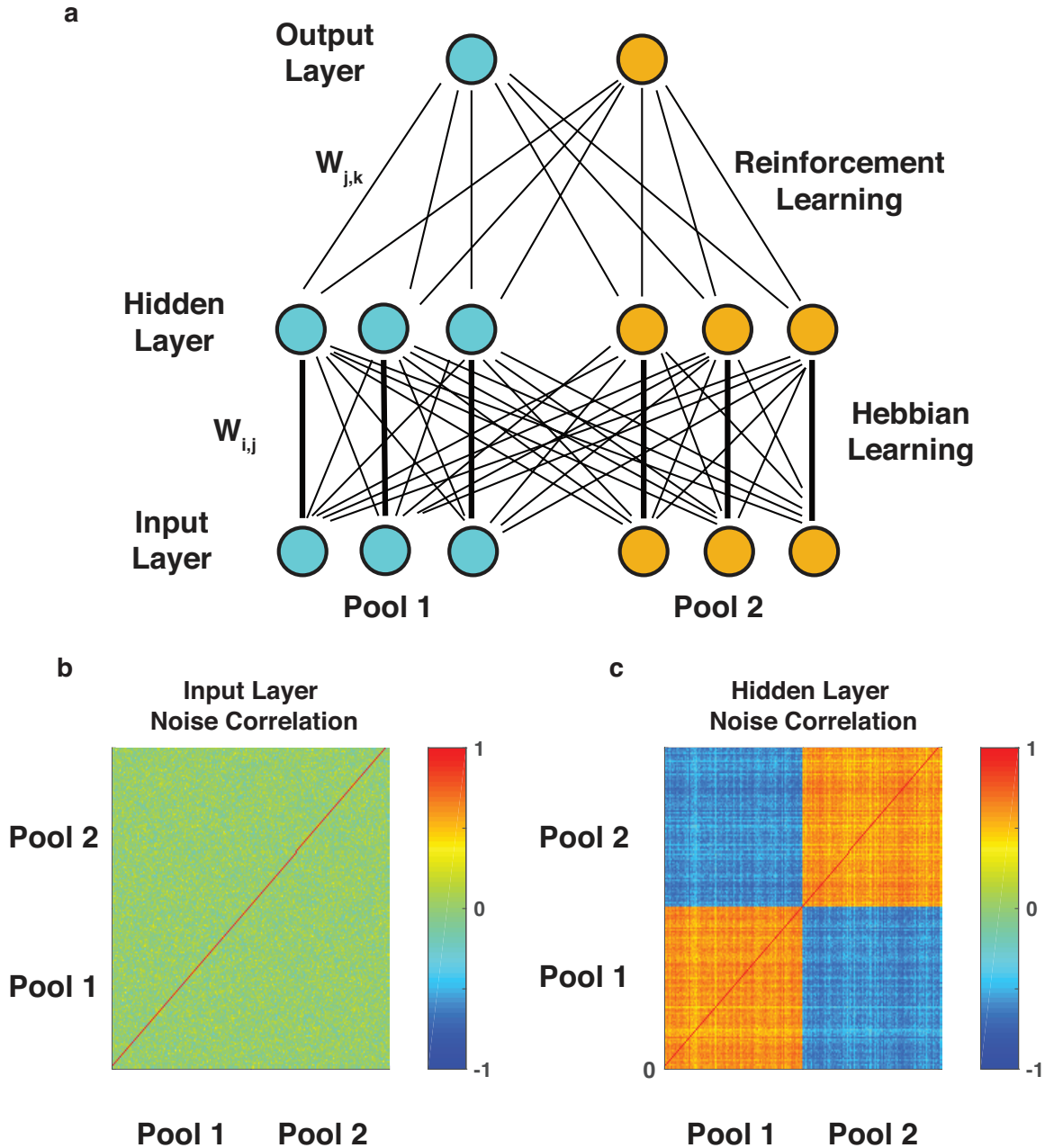737     this case, is the signal (Oja, 1982)
738
739
740
741
742
743

**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.

758

759 *Dynamic, task-dependent noise correlations speed learning by constraining it to relevant*
760 *feature dimensions.*

761

762 In order to understand how noise correlations might impact learning in mixed
763 encoding populations, we extended our perceptual discrimination task to include
764 two directions of motion discrimination (eg. up/down and left/right). On each trial,
765 a cue indicated which of two possible motion discriminations should be
766 performed (figure 7A, left; (Cohen and Newsome, 2008)). We extended our
767 neural network to include four populations of one hundred input units, each
768 population encoding a conjunction of motion directions (up-right, up-left, down-
769 right, down-left; figure 7A; input layer). Two additional inputs provided a perfectly
770 reliable "cue" regarding the relevant feature for the trial (figure 7A; task units).
771 Four output neurons encoded the four possible responses (up, left, down, right)
772 and were fully connected to the input layer (figure 7A; output layer). Task units
773 were hard wired to eliminate irrelevant task responses, but weights of input units
774 were learned over time as in our previous simulations.

775

776 Learning performance in the two-feature discrimination task depended not only
777 on the level of noise correlations, but also on the type. As in the previous
778 simulation, adding noise correlations to each individual population of identically
779 tuned units led to faster learning of the appropriate readout (Figure 7B&C,
780 compare blue and yellow; Figure 7D&E, vertical axis; mean[std] accuracy across
781 training: 0.54[0.05] and 0.70[0.05] for minimum (0) and maximum (0.2) in pool
782 correlations, t-test for difference in accuracy: $t = 226$, $dof = 19998$, $p < 10e\text{-}50$).

783

784 However, the more complex task design also allowed us to test whether dynamic
785 trial-to-trial correlations might further facilitate learning. Specifically, correlations
786 that increase shared variability among units that contribute evidence to the same
787 response have been observed previously (Cohen and Newsome, 2008), and
788 could in principle focus learning on relevant dimensions (figure 2C&D) even
789 when those dimensions change from trial to trial. Indeed, adding correlations
790 among separate pools that share the same encoding of the relevant feature (eg.
791 UP on a vertical trial) led to faster learning (figure 7B; mean[std] training
792 accuracy for model with relevant pool correlations: 0.73[0.05], *t*-test for difference
793 from in pool correlation only model: $t = 34$, $dof = 19998$, $p < 10e\text{-}50$) and weights
794 that more closely approached the optimal readout (figure 7E, horizontal axis). In
795 contrast, when positive noise correlations were introduced across separate
796 encoding pools that shared the same tuning for the irrelevant dimension on each
797 trial (eg. UP on a horizontal trial) learning was impaired dramatically (figure 7C;
798 mean[std] training accuracy for model with irrelevant pool correlations:
799 0.51[0.05], *t*-test for difference from in pool correlation only model: $t = -278$, $dof$
800 $= 19998$, $p < 10e\text{-}50$) and learned weights diverged from the optimal readout
801 (figure 7F, horizontal axis). Model performance differences were completely

24

802    attributable to learning the readout, as all models performed similarly when using
803    the optimal readout (extended data figure 7-1).
804
805    In order to test the idea that noise correlations might focus learning onto relevant
806    dimensions, we extracted weight updates from each trial and projected these
807    updates into a two-dimensional space where the first dimension captured the
808    relative sensitivity to leftward versus rightward motion and the second dimension
809    captured relative sensitivity to upward versus downward motion. In the model
810    where input units were only correlated within their identically tuned pool, weight
811    updates projected in all directions more or less uniformly (figure 7G), and did not
812    differ systematically across trial types (vertical versus horizontal). However,
813    dynamic noise correlations that shared variability across the relevant dimension
814    tended to push weight updates onto the appropriate dimension for a given trial
815    (figure 4F; *t*-test for difference in the magnitude of updating in up/down and
816    left/right dimensions across conditions [up/down – left/right]: $t = 3.4$, *dof*=98, *p =*
817    0.001). In contrast, dynamic noise correlations that shared variability across the
818    irrelevant dimension tended to push weight updates onto the wrong dimension
819    (figure 4H; t-test for difference in the magnitude of updating in up/down and
820    left/right dimensions across conditions [up/down – left/right]: $t = -9.5$, *dof*=98, *p =*
821    10e-14). Both of these trends were consistent across simulations, providing an
822    explanation for the performance improvements achieved by relevant noise
823    correlations (projection of learning onto an appropriate dimension) and
824    performance impairments produced by irrelevant noise correlations (projection of
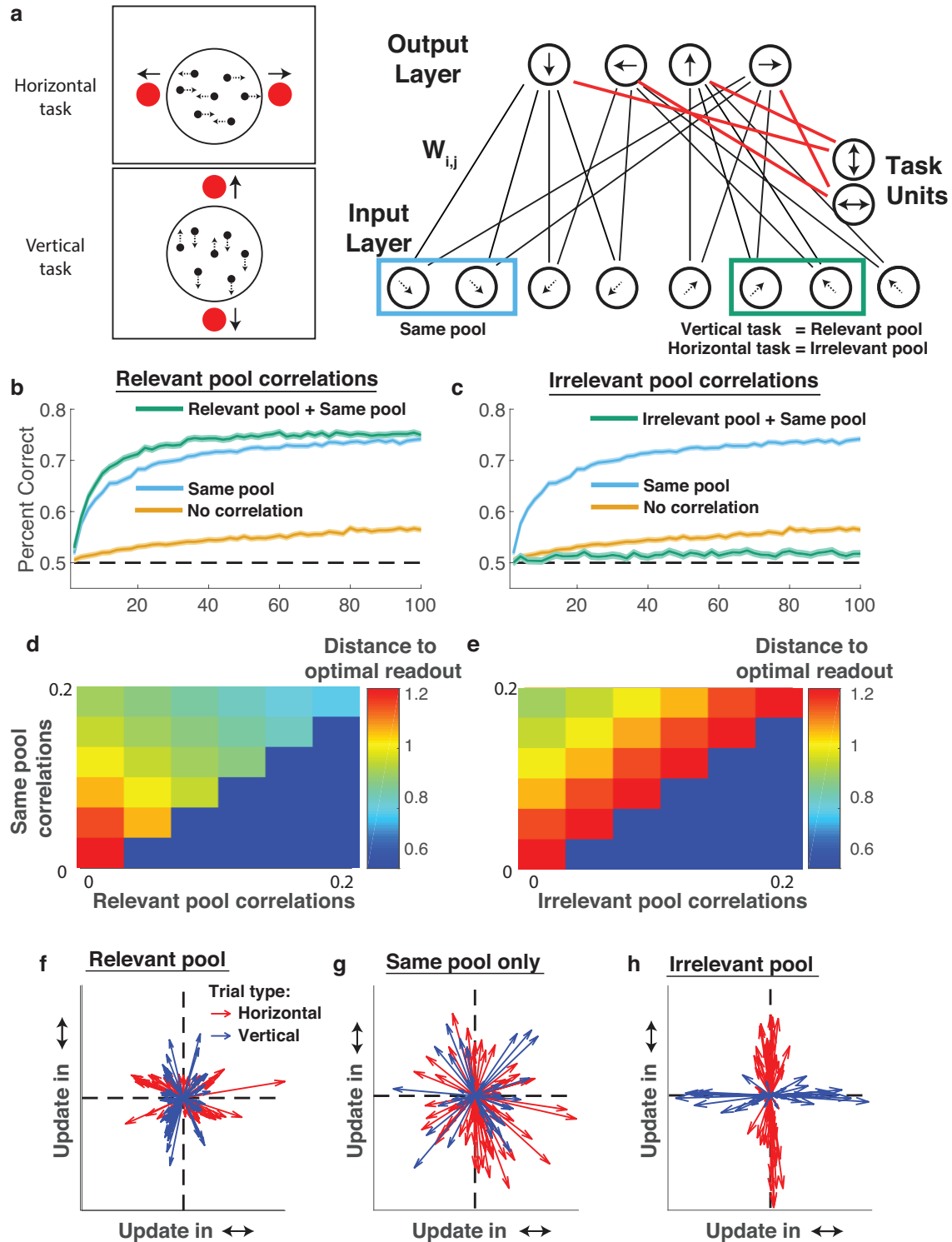825    learning onto an inappropriate dimension).
826
827
828
829
830
831
832
833

**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 feed-forward network in which each homogenous pool of input units represents two dimensions of motion (up versus down, and left versus right), and output units produce

842     responses in favor of alternative actions (up, down, left, right). Each homogenous pool of input
843     units is identically tuned to one of four conjunctions of movement directions: up-left, down-left, up-
844     right, down-right. Two additional input units provide cue information that biases output units to
845     produce an output corresponding to the discrimination appropriate on this trial (eg. horizontal or
846     vertical). Noise correlations were manipulated among 1) identically tuned neurons (blue
847     rectangle; same pool), 2) neurons that have similar encoding of the task relevant feature (green
848     rectangle pair in vertical trials; relevant pool), and 3) neurons that have similar encoding of the
849     task irrelevant feature (green rectangle pair in horizontal trials; irrelevant pool). **B&C**) Learning
850     curves showing accuracy (ordinate) over trials (abscissa) for models 1) lacking noise correlations
851     (orange), 2) containing noise correlations that are limited to neurons that have same tuning for
852     both features (same pool; blue), 3) containing same pool noise correlations along with
853     correlations between neurons in different pools that have the same tuning for the task-relevant
854     feature (in pool+rel pool; green in **B**), and 4) containing in-pool noise correlations along with
855     correlations between neurons in different pools that have the same tuning for the task irrelevant
856     feature (in pool+irrel pool; green in **C**). **D&E**) Distance between learned weights and the optimal
857     readout (color) for models that differ in their level of "in pool" correlations (ordinate, both plots),
858     "relevant pool" correlations (abscissa, **D**), and "irrelevant pool" correlations (abscissa, **E**). **F,G,H**)
859     Weight updates for example learning sessions were projected into a two dimensional space in
860     which net updates to the relative contribution of vertical motion information (eg. up versus down)
861     is represented on the abscissa and updates to the relative contribution of horizontal motion
862     information (eg. left versus right) is represented on the ordinate. Arrows reflect single trial weight
863     updates and are colored according to the trial type (red = horizontal discrimination, blue = vertical
864     discrimination). Weight updates for a model with only "in pool" correlations look similar across trial
865     types (**G**), but weight updates for a model with "relevant pool" correlations indicate more weight
866     updating on the relevant feature (**F**), whereas the opposite was observed in the case of "irrelevant
867     pool" correlations (**H**).
868
869

870 ## Discussion:

871

872     Taken together, our results suggest that in settings where the population signal-
873     to-noise ratio is limited by external factors (eg. inputs) and relevant task
874     representations are low dimensional, noise correlations can make learning faster
875     and more robust by focusing learning on the most relevant dimensions. We
876     demonstrate this basic principle in a simple perceptual learning task (figure 3),
877     where beneficial noise correlations between similarly tuned units could be
878     produced through a simple Hebbian learning rule (figure 6). We extended our
879     framework to a contextual learning task to demonstrate that dynamic noise
880     correlations that bind task relevant feature representations facilitate faster
881     learning (figure 7b&d) by pushing learning onto task-relevant dimensions (figure
882     7f). Given the pervasiveness of noise correlations among similarly tuned sensory
883     neurons (Zohary et al., 1994; Maynard et al., 1999; Bair et al., 2001; Averbeck
884     and Lee, 2003; Cohen and Maunsell, 2009; Huang and Lisberger, 2009; Ecker et
885     al., 2010; Gu et al., 2011; Adibi et al., 2013), and that the noise correlations
886     dynamics beneficial for learning in our simulations are similar to those that have
887     been observed *in vivo* (Cohen and Newsome, 2008), we interpret our results as
888     suggesting that noise correlations between similarly tuned neurons are a feature
889     of neural coding architectures that ensures efficient readout learning, rather than
890     a bug that limits encoding potential.

27

891

892

893 This interpretation rests on several assumptions in our model. Of particular
894 importance is the assumption that signal-to-noise ratio of our populations is fixed,
895 meaning that our manipulation of noise correlations can focus variance on
896 specific dimensions without gaining or losing information. This assumption
897 reflects conditions in which information is limited at the level of the inputs to the
898 population, for instance due to noisy peripheral sensors (Beck et al., 2012;
899 Kanitscheider et al., 2015). In such conditions, even with optimal encoding,
900 population information saturates at an upper bound determined by the
901 information available in the inputs to the population. Therefore, fixing the signal-
902 to-noise ratio enabled us to examine the effect of noise correlations on
903 downstream processes that learn to read-out the population code in the absence
904 of any influence of noise correlations on the quantity of information contained
905 within that population code.

906

907 Previous theoretical work exploring the role of noise correlations in encoding has
908 typically assumed that single neurons have a fixed variance, such that tilting the
909 covariance of neural populations towards or away from the dimension of signal
910 encoding would have a large impact on the amount of information that can be
911 encoded by a population (figure 1a; (Averbeck et al., 2006; Moreno-Bote et al.,
912 2014)). Such assumptions lead to the idea that positive noise correlations among
913 similarly tuned neurons limit encoding potential, raising the question of why they
914 are so common in the brain (Cohen and Kohn, 2011). In considering the
915 implications of this framework, one important question is: if information encoded
916 by the population can be increased by changing the correlation structure among
917 neurons, where does this additional information come from? In some cases, the
918 neural population in question may indeed receive sufficient task relevant
919 information from upstream brain regions to reorganize its encoding in this way,
920 but in other cases it is likely that information is limited by the inputs to a neural
921 population (Kanitscheider et al., 2015; Kohn et al., 2016). In cases where
922 incoming information is limited, further increasing representational capacity is not
923 possible, and formatting information for efficient readout is essentially the best
924 that the population code could do. Here we show that the noise correlations that
925 have previously been described as "information limiting" are exactly the type of
926 correlations that format information most efficiently for readout learning under
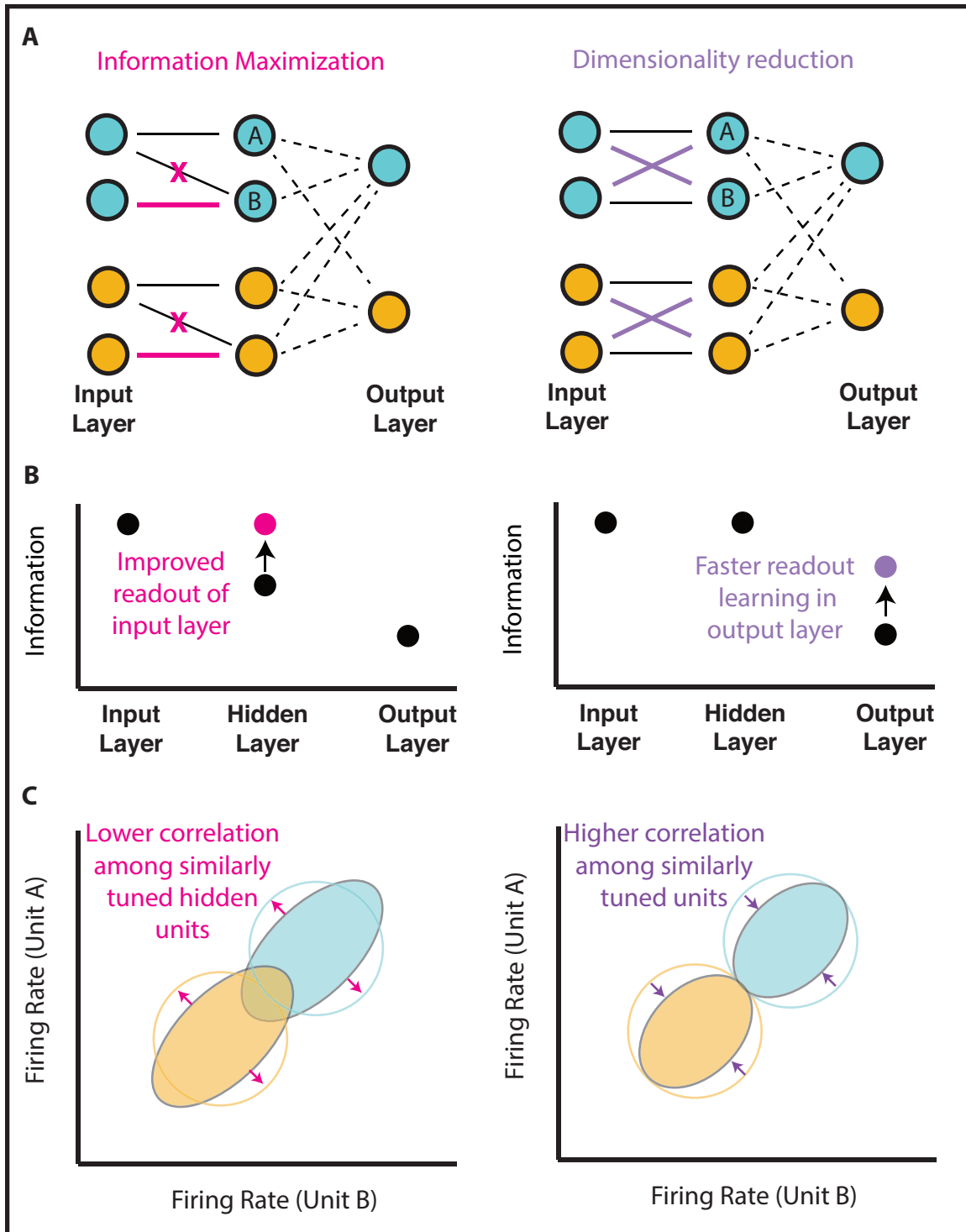927 such conditions.

928

929 Between these two bookends of a fixed signal-to-noise ratio and fixed single unit
930 variance, we also simulated intermediate regimes which do not perfectly
931 preserve the signal-to-noise ratio. In these intermediate regimes, a tradeoff
932 emerges: noise correlations between similarly tuned neurons produce faster
933 learning in the short term, at the cost of lower levels of asymptotic performance in
934 the long run (Fig 5).

935
936     Jointly considering these perspectives on noise correlations provides a more
937     nuanced view of how neural representations are likely optimized for learning. In
938     order to optimize an objective function, a neural population can reduce correlated
939     noise in task relevant dimensions to increase its representational capacity up to
940     some level constrained by its inputs (Figure 8, left). But once the population is
941     fully representing all task relevant information that has been provided to it, it can
942     additionally optimize representations by pushing as much variance onto task
943     relevant dimensions as possible, thereby affording efficient learning in
944     downstream neural populations (Figure 8, right). In short, optimization of a neural
945     population code does not occur in a vacuum, and instead depends critically on
946     both upstream (eg. input constraints) and downstream (eg. readout) neural
947     populations (Figure 8). In this view, if a neural population is *not* fully representing
948     the decision relevant information made available to it, then learning could
949     improve the efficiency of representations by reducing rate limiting noise
950     correlations as has been observed in some paradigms (Gu et al., 2011; Ni et al.,
951     2018). In contrast, once available information is fully represented, readout
952     learning could be further optimized by reformatting population codes such that
953     variability is shared across neurons with similar tuning for the relevant task
954     feature, producing the sorts of dynamic noise correlations that have been
955     observed in well trained animals (Cohen and Newsome, 2008).
956
957

958
959 **Figure 8: Information maximization and dimensionality reduction can be useful for**
960 **learning under different situations and have opposite effects on noise correlations among**
961 **similarly tuned units. A)** A schematic representation of a three layer neural network in which
962 units provide evidence for one of two categorizations (blue/orange). In the left network, the hidden
963 layer initially has access to information from only one of two independent units in each pool, but
964 weights are subsequently adjusted to increase task-relevant information represented in the
965 hidden layer (pink). In the right network, the hidden layer initially has access to all task-relevant
966 information, but weights are subsequently adjusted to share signal and noise across similarly

tuned units to afford 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 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 affording information maximization (left) increase task relevant information in the hidden layer (pink), whereas weight adjustments that afford 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 panel A) *decrease* correlations among hidden units A&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&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 panel 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).

In addition to key assumptions about an external limitation on signal-to-noise, our modeling included a number of 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 a continuous space of stimuli and the heterogeneous populations observed in sensory cortical regions of the brain. Previous work has shown that noise correlations do not necessarily limit encoding potential in heterogenous populations with diverse tuning (Shamir and Sompolinsky, 2004; 2006; Chelaru and Dragoi, 2008; Ecker et al., 2011). A primary goal of our work was to identify the computational principles that control the speed at which readout can be learned, and our simplified populations are considerably more tractable and transparent than realistic neural populations. The principles that we identify here are certainly at play in real neural populations, albeit with implications that are far less transparent. In particular, in a population with diverse tuning profiles, the degree to which individual neurons are informative about a task-relevant discrimination will vary. To benefit learning through coordinated weighted changes, the correlation structure in such populations should reflect this variability. Empirical studies in macaques suggest that such variability is indeed present in real populations (Cohen and Maunsell, 2009; Rabinowitz et al., 2015; Ruff and Cohen, 2016; Bondy et al., 2018). We hope that our simplified results pave the way for future work to assess nuances that can emerge in mixed heterogeneous populations, or in more realistic architectures that go beyond the simple feed forward flow of information considered here.

1016  *Model predictions*

1017

1018  Our work shows that noise correlations can focus the gradient of learning onto
1019  the most appropriate dimensions. Thus, our model predicts that the degree to
1020  which similarly tuned neurons are correlated during a perceptual discrimination
1021  should be positively related to performance improvements experienced on
1022  subsequent discriminations. In contrast, our model predicts that the degree of
1023  correlation between neurons that are similarly tuned to a task irrelevant feature
1024  should control the degree of learning on irrelevant dimensions, and thus
1025  negatively relate to performance improvements on subsequent discriminations.
1026  These predictions are strongest for the earliest stages of learning where weight
1027  adjustments are critical for subsequent performance, but they may also hold for
1028  later stages of learning, when correlations on irrelevant dimensions, including
1029  independent noise channels, could potentially lead to systematic deviations from
1030  optimal readout (figure 2f, 4d&e). These predictions could be tested by recording
1031  neural responses to a stimulus set that differs across multiple features to
1032  characterize both signal-to-noise and correlated variability for each feature
1033  discrimination. A strong prediction of our model is that correlated variability within
1034  neurons tuned to a given feature should be a predictor of subsequent learning of
1035  responses to that feature – above and beyond feature value discriminability.

1036

1037  One interesting special case involves tasks where the relevant dimension
1038  changes in an unsignaled manner (Birrell and Brown, 2000). In such tasks, noise
1039  correlations on the previously relevant dimension would, after such an
1040  "extradimensional shift", force gradients into a task-irrelevant dimension and thus
1041  impair learning performance. Interestingly, learning after extra-dimensional shifts
1042  can be selectively improved by enhancing noradrenergic signaling (Devauges
1043  and Sara, 1990; Lapiz and Morilak, 2006), which leads to increased arousal
1044  (Joshi et al., 2016; Reimer et al., 2016) and decreased cortical pairwise noise
1045  correlations in sensory and higher order cortex (Vinck et al., 2015; Joshi and
1046  Gold, n.d.). While these observations have been made in different paradigms,
1047  our model suggests that the reduction of noise correlations resulting from
1048  increased sustained levels of norepinephrine after an extradimensional shift
1049  (Bouret and Sara, 2005) could mediate faster learning by expanding the
1050  dimensionality of the learning gradients (compare figure 7G to 7F) to consider
1051  features that have not been task-relevant in the past.

1052

1053  *Relation to attentional effects on noise correlations*

1054

1055  In broad strokes, our finding that manipulation of noise correlations can focus
1056  variance on specific dimensions is in line with specific models of attention. In
1057  particular, noise reduction in task irrelevant dimensions might be considered in
1058  the same light that is often cast on suppression of task irrelevant dimensions by
1059  attentional mechanisms (Zanto and Gazzaley, 2009), in particular for purposes of

1060 accurate credit assignment (Akaishi et al., 2016; Leong et al., 2017). One
1061 possibility is that compressed low-dimensional task representations in higher-
1062 order decision regions (Mack et al., 2019) may pass accumulated decision
1063 related information back to sensory regions in order to approximate Bayesian
1064 inference (Haefner et al., 2016; Bondy et al., 2018; Lange et al., 2018). As task
1065 relevant features are learned, such a process would promote noise correlations
1066 between neurons coding those relevant features. In other words, noise
1067 correlations may reflect a chosen hypothesis about which feature is relevant for
1068 predicting outcomes. Such a signal would be beneficial if it could persist (and
1069 thus preserve correlations between neurons tuned to the same task relevant
1070 feature value) until the time of feedback or reinforcement. Recent work showing
1071 strengthened noise correlations between similarly tuned neurons during working
1072 memory maintenance suggests that this might very well be the case (Merrikhi et
1073 al., 2018).

1074

1075 One observation that seems at odds with this interpretation is that manipulations
1076 of attention that cue a particular location or feature tend to decrease noise
1077 correlations among neurons that encode that location or feature (Cohen and
1078 Maunsell, 2009; Mitchell et al., 2009; Cohen and Maunsell, 2011; Herrero et al.,
1079 2013; Doiron et al., 2016). The effects of attentional cuing on noise correlations
1080 are dynamic in that cues change from one trial to the next, and contextual, in that
1081 noise correlations are reduced most dramatically among neurons that contribute
1082 evidence toward the same response in a manner consistent with increasing the
1083 amount of task relevant information in the population code (Ruff and Cohen,
1084 2014; Downer et al., 2015). These effects are generally observed in well-trained
1085 animals during task performance and may not result from the same processes as
1086 the longer timescale noise correlation structure. Indeed, there may be a tradeoff
1087 between learning and performance, particularly if the computations giving rise to
1088 noise correlations do so without perfectly preserving signal-to-noise ratio. Our
1089 model does not account for these attentional effects, as we intentionally
1090 constrained the signal-to-noise ratio of our neural populations, thereby
1091 eliminating any potential changes in information encoding potential. However, we
1092 hope that our work motivates future studies to jointly consider the impacts of
1093 noise correlations on both learning and immediate performance in order to better
1094 understand the potentially competing imperatives that the brain faces in
1095 dynamically controlling the correlation structure of its own representations (see
1096 (Haimerl et al., 2019) for one attempt to do so).

1097

1098

1099 *Origins of useful noise correlations*

1100

1101 One important question stemming from our work is how noise correlations
1102 emerge in the brain. This question has been one of longstanding debate, largely
1103 because there are so many potential mechanisms through which correlations

1104    could emerge (Kanitscheider et al., 2015; Kohn et al., 2016). Noise correlations
1105    could emerge from convergent and divergent feed forward wiring (Shadlen and
1106    Newsome, 1998), local connectivity patterns within a neural population (Hansen
1107    et al., 2012; Smith et al., 2013), or top down inputs provided separately to
1108    different neural populations (Haefner et al., 2016). Here we show that static noise
1109    correlations that are useful for perceptual learning emerge naturally from
1110    Hebbian learning in a feed-forward network. While this certainly suggests that
1111    useful noise correlations could emerge through feed forward wiring, it is also
1112    possible to consider our Hebbian learning as occurring in a one-step recurrence
1113    of the input units, and thus the same data support the possibility of noise
1114    correlations through local recurrence. The context dependent noise correlations
1115    that speed learning (figure 7), however, would not arise through simple Hebbian
1116    learning. Such correlations could potentially be produced through selective top-
1117    down signals from the choice neurons, as has been previously proposed
1118    (Wimmer et al., 2015; Haefner et al., 2016; Bondy et al., 2018; Lange et al.,
1119    2018). Moreover, top-down input may selectively target neuronal ensembles
1120    produced through Hebbian learning (Collins and Frank, 2013). While previous
1121    work has suggested that such a mechanism could be adaptive for accumulating
1122    information over the course of a decision (Haefner et al., 2016), our work
1123    demonstrates that the same mechanism could effectively be used to tag relevant
1124    neurons for weight updating between trials, making efficient use of top-down
1125    circuitry. Haimerl et al. recently made a similar point, showing that stochastic
1126    modulatory signals shared across task-informative neurons can serve to tag
1127    them for a decoder (Haimerl et al., 2019).
1128
1129    *Noise correlations as inductive biases*
1130
1131    Artificial intelligence has undergone a revolution over the past decade leading to
1132    human level performance in a wide range of tasks (Mnih et al., 2015). However,
1133    a major issue for modern artificial intelligence systems, which build heavily on
1134    neural network architectures, is that they require far more training examples than
1135    a biological system would (Hassabis et al., 2017). This biological advantage
1136    occurs despite the fact that the total number of synapses in the human brain,
1137    which could be thought of as the free parameters in our learning architecture, is
1138    much greater than the number of weights in even the most parameter-heavy
1139    deep learning architectures. Our work provides some insight into why this occurs;
1140    correlated variability across neurons in the brain constrain learning to specific
1141    dimensions, thereby limiting the effective complexity of the learning problem
1142    (figures 4A, 7F-G). We show that, for simple tasks, this can be achieved using
1143    Hebbian learning rules (figure 6), but that contextual noise correlations, of the
1144    form that might be produced through top-down signals (Haefner et al., 2016), are
1145    critical for appropriately focusing learning in more complex circumstances. In
1146    principle, algorithms that effectively learn and implement noise correlations might
1147    reduce the amount of data needed to train AI systems by limiting degrees of

1148  freedom to those dimensions that are most relevant. Furthermore, our work
1149  suggests that large scale neural recordings in early stages of learning complex
1150  tasks might serve as indicators of the inductive biases that constrain learning in
1151  biological systems.

1153  In summary, we show that under external constraints of task-relevant
1154  information, noise correlations that have previously been called "rate limiting" can
1155  serve an important role in constraining learning to task-relevant dimensions. In
1156  the context of previous theory focusing on representation, our work suggests that
1157  neural populations are subject to competing forces when optimizing covariance
1158  structures; on one hand reducing correlations between pairs of similarly tuned
1159  neurons can be helpful to fully represent available information, but increasing
1160  correlations among similarly tuned neurons can be helpful for assigning credit to
1161  task relevant features. We believe that this view of the learning process not only
1162  provides insight to understanding the role of noise correlations in the brain, but
1163  opens up the door to better understand the inductive biases that guide learning in
1164  biological systems.

1167  References:

1170  Adibi M, McDonald JS, Clifford CWG, Arabzadeh E (2013) Adaptation improves
1171      neural coding efficiency despite increasing correlations in variability. Journal
1172      of Neuroscience 33:2108–2120.

1173  Akaishi R, Kolling N, Brown JW, Rushworth M (2016) Neural Mechanisms of
1174      Credit Assignment in a Multicue Environment. Journal of Neuroscience
1175      36:1096–1112.

1176  Averbeck BB, Latham PE, Pouget A (2006) Neural correlations, population
1177      coding and computation. Nature Reviews Neuroscience 7:358–366.

1178  Averbeck BB, Lee D (2003) Neural noise and movement-related codes in the
1179      macaque supplementary motor area. Journal of Neuroscience 23:7630–
1180      7641.

1181  Bair W, Zohary E, Newsome WT (2001) Correlated firing in macaque visual area
1182      MT: time scales and relationship to behavior. Journal of Neuroscience
1183      21:1676–1697.

1184  Beck JM, Ma WJ, Pitkow X, Latham PE, Pouget A (2012) Perspective. Neuron
1185      74:30–39.

1186  Birrell JM, Brown VJ (2000) Medial frontal cortex mediates perceptual attentional

1187    set shifting in the rat. Journal of Neuroscience 20:4320–4324.

1188    Bondy AG, Haefner RM, Cumming BG (2018) Feedback determines the structure
1189        of correlated variability in primary visual cortex. Nature Publishing Group:1–
1190        15.

1191    Bouret S, Sara SJ (2005) Network reset: a simplified overarching theory of locus
1192        coeruleus noradrenaline function. Trends in Neurosciences 28:574–582.

1193    Chelaru MI, Dragoi V (2008) Efficient coding in heterogeneous neuronal
1194        populations. Proceedings of the National Academy of Sciences 105:16344–
1195        16349.

1196    Cohen MR, Kohn A (2011) Measuring and interpreting neuronal correlations.
1197        Nature Publishing Group 14:811–819.

1198    Cohen MR, Maunsell JHR (2009) Attention improves performance primarily by
1199        reducing interneuronal correlations. Nature Publishing Group 12:1594–1600.

1200    Cohen MR, Maunsell JHR (2011) Using neuronal populations to study the
1201        mechanisms underlying spatial and feature attention. Neuron 70:1192–1204.

1202    Cohen MR, Newsome WT (2008) Context-Dependent Changes in Functional
1203        Circuitry in Visual Area MT. Neuron 60:162–173.

1204    Collins AGE, Frank MJ (2013) Cognitive control over learning: creating,
1205        clustering, and generalizing task-set structure. Psychological Review
1206        120:190–229.

1207    Devauges V, Sara SJ (1990) Activation of the noradrenergic system facilitates an
1208        attentional shift in the rat. Behavioural Brain Research 39:19–28.

1209    Doiron B, Litwin-Kumar A, Rosenbaum R, Ocker GK, Josić K (2016) The
1210        mechanics of state-dependent neural correlations. Nature Publishing Group
1211        19:383–393.

1212    Downer JD, Niwa M, Sutter ML (2015) Task engagement selectively modulates
1213        neural correlations in primary auditory cortex. Journal of Neuroscience
1214        35:7565–7574.

1215    Ecker AS, Berens P, Keliris GA, Bethge M, Logothetis NK, Tolias AS (2010)
1216        Decorrelated neuronal firing in cortical microcircuits. Science 327:584–587.

1217    Ecker AS, Berens P, Tolias AS, Bethge M (2011) The effect of noise correlations
1218        in populations of diversely tuned neurons. Journal of Neuroscience
1219        31:14272–14283.

1220  Gu Y, Liu S, Fetsch CR, Yang Y, Fok S, Sunkara A, DeAngelis GC, Angelaki DE
1221       (2011) Perceptual learning reduces interneuronal correlations in macaque
1222       visual cortex. Neuron 71:750–761.

1223  Haefner RM, Pietro Berkes, Fiser J (2016) Perceptual Decision-Making as
1224       Probabilistic Inference by Neural Sampling. Neuron 90:649–660.

1225  Haimerl C, Savin C, Simoncelli EP (2019) Flexible and accurate decoding of
1226       neural populations through stochastic comodulation. Biorxiv 21:598.

1227  Hansen BJ, Chelaru MI, Dragoi V (2012) Correlated variability in laminar cortical
1228       circuits. Neuron 76:590–602.

1229  Hassabis D, Kumaran D, Summerfield C, Botvinick M (2017) Neuroscience-
1230       Inspired Artificial Intelligence. Neuron 95:245–258.

1231  Hawkey DJC, Amitay S, Moore DR (2004) Early and rapid perceptual learning.
1232       Nature Publishing Group 7:1055–1056.

1233  Herrero JL, Gieselmann MA, Sanayei M, Thiele A (2013) Attention-induced
1234       variance and noise correlation reduction in macaque V1 is mediated by
1235       NMDA receptors. Neuron 78:729–739.

1236  Hu, Y., Zylberberg, J., & Shea-Brown, E. (2014). The sign rule and beyond: boundary
1237       effects, flexibility, and noise correlations in neural population codes. PLoS Comput
1238       Biol, 10(2), e1003469.

1239  Huang X, Lisberger SG (2009) Noise correlations in cortical area MT and their
1240       potential impact on trial-by-trial variation in the direction and speed of
1241       smooth-pursuit eye movements. Journal of Neurophysiology 101:3012–3030.

1242  Joshi S, Gold JI (n.d.) Context-Dependent Relationships between Locus
1243       Coeruleus Firing Patterns and Coordinated Neural Activity in the Anterior
1244       Cingulate Cortex. Biorxiv.

1245  Joshi S, Li Y, Kalwani RM, Gold JI (2016) Relationships between Pupil Diameter
1246       and Neuronal Activity in the Locus Coeruleus, Colliculi, and Cingulate Cortex.
1247       Neuron 89:221–234.

1248  Kanitscheider I, Coen-Cagli R, Pouget A (2015) Origin of information-limiting
1249       noise correlations. Proceedings of the National Academy of Sciences
1250       112:E6973–E6982.

1251  Kohn A, Coen-Cagli R, Kanitscheider I, Pouget A (2016) Correlations and
1252       Neuronal Population Information. Annu Rev Neurosci 39:237–256.

1253  Krotov D, Hopfield JJ (2019) Unsupervised learning by competing hidden units.

1254    Proceedings of the National Academy of Sciences 116:7723–7731.

1255    Lange RD, Chattoraj A, Beck JM, Yates JL, Haefner RM (2018) A confirmation
1256        bias in perceptual decision-making due to hierarchical approximate inference.
1257        Biorxiv.

1258    Lapiz MDS, Morilak DA (2006) Noradrenergic modulation of cognitive function in
1259        rat medial prefrontal cortex as measured by attentional set shifting capability.
1260        Neuroscience 137:1039–1049.

1261    Law C-T, Gold JI (2009) Reinforcement learning can account for associative and
1262        perceptual learning on a visual-decision task. Nature Neuroscience 12:655–
1263        663.

1264    Leong YC, Radulescu A, Daniel R, DeWoskin V, Niv Y (2017) Dynamic
1265        Interaction between Reinforcement Learning and Attention in
1266        Multidimensional Environments. Neuron 93:451–463.

1267    Mack ML, Preston AR, Love BC (2019) Ventromedial prefrontal cortex
1268        compression during concept learning. Nature Communications:1–11.

1269    Maynard EM, Hatsopoulos NG, Ojakangas CL, Acuna BD, Sanes JN, Normann
1270        RA, Donoghue JP (1999) Neuronal interactions improve cortical population
1271        coding of movement direction. Journal of Neuroscience 19:8083–8093.

1272    Merrikhi Y, Clark K, Noudoost B (2018) Concurrent influence of top-down and
1273        bottom-up inputs on correlated activity of Macaque extrastriate neurons.
1274        Nature Communications 9:5393.

1275    Mitchell JF, Sundberg KA, Reynolds JH (2009) Spatial attention decorrelates
1276        intrinsic activity fluctuations in macaque area V4. Neuron 63:879–888.

1277    Mnih V, Kavukcuoglu K, Silver D, Rusu AA, Veness J, Bellemare MG, Graves A,
1278        Riedmiller M, Fidjeland AK, Ostrovski G, Petersen S, Beattie C, Sadik A,
1279        Antonoglou I, King H, Kumaran D, Wierstra D, Legg S, Hassabis D (2015)
1280        Human-level control through deep reinforcement learning. Nature 518:529–
1281        533.

1282    Moreno-Bote R, Beck J, Kanitscheider I, Pitkow X, Latham P, Pouget A (2014)
1283        Information-limiting correlations. Nature Publishing Group 17:1410–1417.

1284    Ni AM, Ruff DA, Alberts JJ, Symmonds J, Cohen MR (2018) Learning and
1285        attention reveal a general relationship between population activity and
1286        behavior. Science 359:463–465.

1287    Oja E (1982) Simplified neuron model as a principal component analyzer. Journal

1288      of Mathematical Biology:1–7.

1289    Pouget A, Dayan P, Zemel R (2000) Information processing with population
1290        codes. Nature Reviews Neuroscience 1:125–132.

1291    Rabinowitz NC, Goris RL, Cohen M, Simoncelli EP (2015) Attention stabilizes the
1292        shared gain of V4 populations. eLife 4:e08998.

1293    Reimer J, McGinley MJ, Liu Y, Rodenkirch C, Wang Q, McCormick DA, Tolias
1294        AS (2016) Pupil fluctuations track rapid changes in adrenergic and
1295        cholinergic activity in cortex. Nature Communications 7:13289.

1296    Ruff DA, Cohen MR (2014) Attention can either increase or decrease spike count
1297        correlations in visual cortex. Nature Publishing Group 17:1591–1597.

1298    Ruff DA, Cohen MR (2016) Stimulus Dependence of Correlated Variability across
1299        Cortical Areas. Journal of Neuroscience 36:7546–7556.

1300    Shadlen MN, Newsome WT (1998) The variable discharge of cortical neurons:
1301        implications for connectivity, computation, and information coding. J Neurosci
1302        18:3870–3896.

1303    Shamir M, Sompolinsky H (2004) Nonlinear population codes. Neural Comput
1304        16:1105–1136.

1305    Shamir M, Sompolinsky H (2006) Implications of neuronal diversity on population
1306        coding. Neural Comput 18:1951–1986.

1307    Smith MA, Jia X, Zandvakili A, Kohn A (2013) Laminar dependence of neuronal
1308        correlations in visual cortex. Journal of Neurophysiology 109:940–947.

1309    Stringer C, Michaelos M, Pachitariu M (2019) High precision coding in mouse
1310        visual cortex. Biorxiv.

1311    Tsividis P, Pouncy T, Xu JL, Tenenbaum JB, Gershman SJ (2017) Human
1312        Learning in Atari. 2017 AAAI Spring Symposium Series, Science of
1313        Intelligence: Computational Principles of Natural and Artificial Intelligence:1–
1314        4.

1315    Vinck M, Batista-Brito R, Knoblich U, Cardin JA (2015) Arousal and Locomotion
1316        Make Distinct Contributions to Cortical Activity Patterns and Visual Encoding.
1317        Neuron 86:740–754.

1318    Wimmer RD, Schmitt LI, Davidson TJ, Nakajima M, Deisseroth K, Halassa MM
1319        (2015) Thalamic control of sensory selection in divided attention. Nature
1320        526:705–709.

1321    Zanto TP, Gazzaley A (2009) Neural Suppression of Irrelevant Information
1322        Underlies Optimal Working Memory Performance. Journal of Neuroscience
1323        29:3059–3066.

1324    Zohary E, Shadlen MN, Newsome WT (1994) Correlated neuronal discharge rate
1325        and its implications for psychophysical performance. Nature 370:140–143.
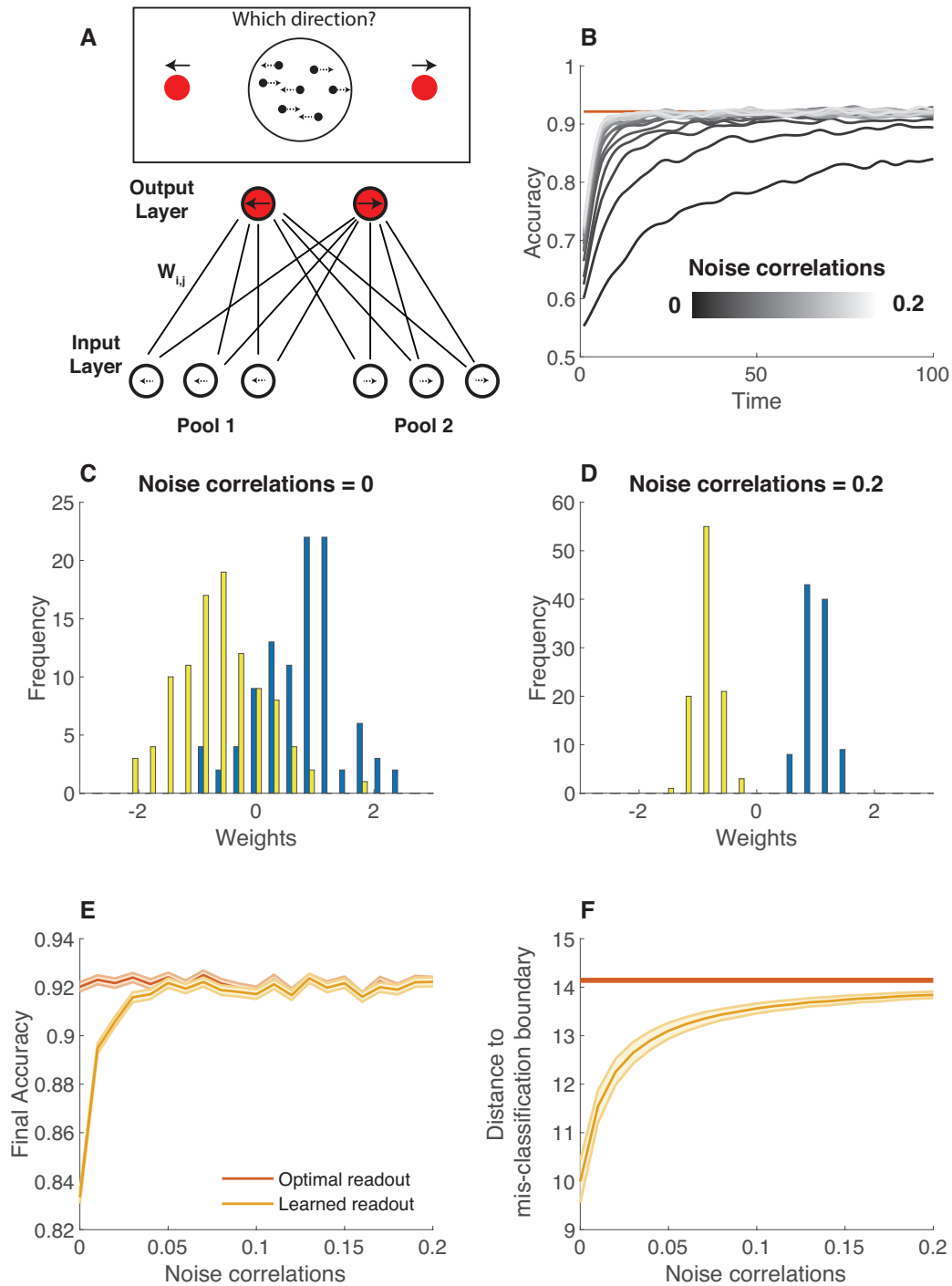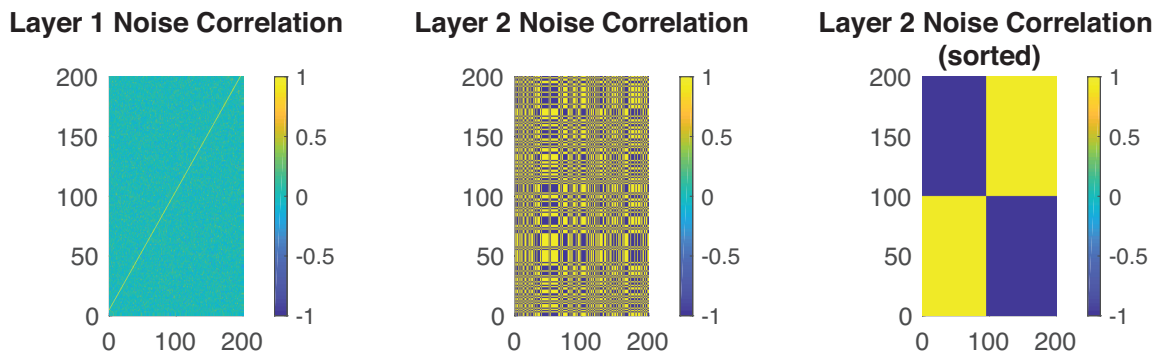
1326

1327
1328
1329
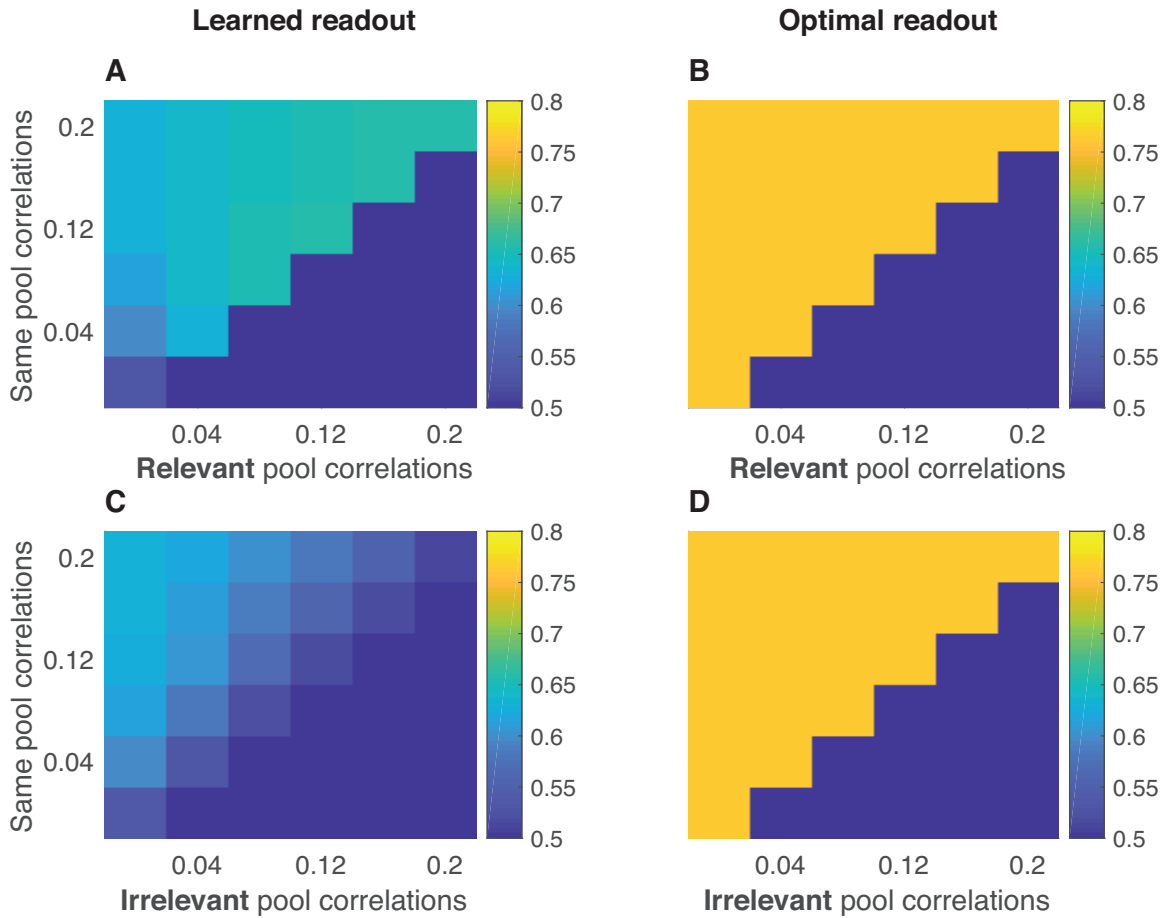1330    Extended data:
1331
1332
1333

1334
1335 **Extended data figure 3-1: Noise correlations that maintain signal-to-noise ratio by scaling**
1336 **signal lead to faster and more robust learning of a perceptual discrimination.** This figure is
1337 a replication of results reported in figure 3 of the main text, except that noise correlations are
1338 produced using equation 25 such that each unit maintains the same fixed variance across noise
1339 correlation conditions, and signal is scaled to maintain a fixed signal-to-noise ratio. **A)** A two-layer
1340 feed-forward neural network was designed to solve a two alternative forced choice motion
1341 discrimination task at or near perceptual threshold. Input layer contains two pools of neurons that
1342 provide evidence for alternate percepts (eg. leftward motion versus rightward motion) and output

41

1343  neurons encode alternate courses of actions (eg. saccade left versus saccade right). Layers are
1344  fully connected with weights randomized to small values and adjusted after each trial according to
1345  rewards (see methods). **B)** Average learning curves for neural network models in which
1346  population signal-to-noise ratio in pools 1 and 2 were fixed, but noise correlations (grayscale)
1347  were allowed to vary from small (dark) to large (light) values. **C&D)** Weight differences (left output
1348  – right output) for each input unit (color coded according to pool) after 100 timesteps of learning
1349  for low (**C**) and high (**D**) noise correlations. **E)** Accuracy in the last 20 training trials is plotted as a
1350  function of noise correlations for learned readouts (orange) and optimal readout (red).
1351  Lines/shading reflect Mean/SEM. F) The shortest distance, in terms of neural activation, required
1352  to take the mean input for a given category (eg. left or right) to the boundary that would result in
1353  misclassification is plotted for the final learned (orange) and optimal (red) weights for each noise
1354  correlation condition (abscissa). Lines/shading reflect Mean/SEM.

1355
1356
1357
1358
1359
1360
1361

**Layer 1 Noise Correlation**　　**Layer 2 Noise Correlation**　　**Layer 2 Noise Correlation (sorted)**



1362
1363  **Extended data figure 6-1**: **Emergence of noise correlations from Hebbian learning does not**
1364  **depend on weight initialization**. In order to test whether beneficial noise correlations might have
1365  emerged in our Hebbian learning simulations due to our initialization biasing one-to-one
1366  connectivity between the input and hidden layers (see figure 4), we repeated these simulations
1367  with in a network that was initialized with random normal weight projects from layer 1 to layer 2.
1368  Simulations included performance of 200 trials with noise correlations measured in the final 100
1369  trials of the simulation. Activity of layer 1 units was defined by a multivariate Gaussian with zero
1370  covariance elements, and thus it is not surprising that pairwise noise correlations measured in the
1371  activity of that layer were near zero (Left). Activity of layer 2 units was sculpted through Hebbian
1372  learning that shaped connectivity between layer 1 and layer 2. This learning led many pairs of
1373  neurons in layer 2 to become highly correlated and many other pairs to become anti-correlated
1374  (middle, blue and yellow elements, respectively). In order to understand the structure defining
1375  these correlations, we sorted the layer 2 units according to their relative projections to the two
1376  output units (weight to output 1 – weight to output 2), and recomputed the pairwise correlations
1377  (right). This sorting reveals that layer 2 units projecting to the same output unit are positively
1378  correlated with one another, whereas they negatively correlate with the layer 2 units that project
1379  to the opposing output neuron (note block diagonal structure).

1380
1381
1382
1383
1384

1385
1386



1387
1388
1389
1390 **Extended data figure 7-1**: **Noise correlations affect speed of learning, but not performance**
1391 **using optimal readout in multiple discrimination task**. **A)** Mean test accuracy (color) of all
1392 models spanning the range of in pool correlations (abscissa) and relevant pool correlations
1393 (ordinate). **B)** Mean accuracy of same models using optimal readout, rather than the learned
1394 readout.  **C)** Mean test accuracy (color) of all models spanning the range of in pool correlations
1395 (abscissa) and irrelevant pool correlations (ordinate). **D)** Mean accuracy of same models using
1396 optimal readout, rather than the learned readout. Note that performance of all models is identical
1397 when readout is optimal, rather than learned.
1398
1399

1400