# Noise correlations for faster and more robust learning

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

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

46
47 **Abstract:**
48

49 Distributed population codes are ubiquitous in the brain and pose a challenge to
50 downstream neurons that must learn an appropriate readout. Here we explore
51 the possibility that this learning problem is simplified through inductive biases
52 implemented by stimulus-independent noise correlations that constrain learning
53 to task-relevant dimensions. We test this idea in a set of neural networks that
54 learn to perform a perceptual discrimination task. Correlations among similarly
55 tuned units were manipulated independently of overall population signal-to-noise
56 ratio in order to test how the format of stored information affects learning. Higher
57 noise correlations among similarly tuned units led to faster and more robust
58 learning, favoring homogenous weights assigned to neurons within a functionally
59 similar pool, and could emerge through Hebbian learning. When multiple
60 discriminations were learned simultaneously, noise correlations across relevant
61 feature dimensions sped learning whereas those across irrelevant feature
62 dimensions slowed it. Our results complement existing theory on noise
63 correlations by demonstrating that when such correlations are produced without
64 degradation of signal-to-noise ratio, they can improve readout learning by
65 constraining it to appropriate dimensions.
66
67
68 **Introduction:**
69

70 The brain represents information using distributed population codes in which
71 particular feature values are encoded by large numbers of neurons. One
72 advantage of such codes is that a pooled readout across many neurons can
73 effectively reduce the impact of stimulus-independent variability (noise) in the
74 firing of individual neurons (Pouget et al., 2000). However, the extent to which
75 this benefit can be employed in practice is constrained by noise correlations, or
76 the degree to which stimulus-independent variability is shared across neurons in
77 the population (Averbeck et al., 2006). In particular, positive noise correlations
78 between neurons that share the same stimulus tuning can reduce the amount of
79 decodable information in the neural population (Averbeck et al, 2006; Moreno-
80 Bote et al., 2014; Hu et al., 2014). Despite their detrimental effect on encoding,
81 noise correlations of this type are reliably observed, even after years of training
82 on perceptual tasks (Cohen and Kohn, 2011). Furthermore, noise correlations
83 between neurons are dynamically enhanced under conditions where two neurons
84 provide evidence for the same response in a perceptual categorization task
85 (Cohen and Newsome, 2008), raising questions about whether they might serve
86 a function rather than simply reflecting a suboptimal encoding strategy.
87
88 At the same time, learning to effectively read out a distributed code also poses a
89 significant challenge. Learning the appropriate weights for potentially tens of

thousands of neurons in a low signal-to-noise regime is a difficult, high-dimensional problem, requiring a very large number of learning trials and entailing considerable risk of "over fitting" to specific patterns of noise across the neural populations encountered during learning trials. Nonetheless, people and animals can rapidly learn to perform perceptual discrimination tasks, albeit with performance that does not approach theoretically achievable levels (Hawkey et al., 2004; Stringer et al., 2019). In comparison, deep neural networks capable of achieving human level performance typically require a far greater number of learning trials than would be required by humans and other animals (Tsividis et al., 2017). This raises the question of how brains might implement inductive biases to enable efficient learning in high dimensional spaces.

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

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

3

133  correlations can make learning faster and more robust by controlling the
134  dimensions over which learning occurs.
135
136
137  **Methods:**

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

153  *Learning readout in perceptual learning task*

154  Simulations and analyses for a simple perceptual discrimination task were
155  performed with a simplified and statistically tractable two-layer feed-forward
156  neural network (figure 3A). The input layer consisted of two pools of 100 units
157  that were each "tuned" to one of two motion directions (left, right). On each trial
158  normalized firing rates for the neural population were drawn from a multivariate
159  normal distribution that was specified by a vector of stimulus-dependent mean
160  firing rates (signal: +1 for preferred stimulus, -1 for non-preferred stimulus) and a
161  covariance matrix. All elements of the covariance matrix corresponding to
162  covariance between units that were "tuned" to different stimuli were set to zero.
163  The key manipulation was to systematically vary the magnitude of diagonal
164  covariance components (eg. noise in the firing of individual units) and the "same
165  pool" covariance elements (eg. shared noise across identically tuned neurons)
166  while maintaining a fixed level of variance in the summed population response for
167  each pool:

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

168  Where $\sigma^2_{pool}$ is the variance on the sum of normalized firing rates from neurons
169  within a given pool, n is the number of units in the pool and the within pool
170  covariance ($Cov(within\ pool)$) specifies the covariance of pairs of units

4

171  belonging to the same pool. The signal-to-noise ratio (signal/$\sigma^2_{pool}$) for each pool
172  was fixed to one. Given this constraint, the fraction of noise that was shared
173  across neurons within the same pool was manipulated as follows:

174

$$\sigma^2_{unit} \; = \; \frac{\sigma^2_{pool}}{n \, + \, n(n-1)\phi} \quad Eq.\,2$$

175

176

$$Cov(within\,pool) \; = \; \phi\sigma^2_{unit} \quad Eq.\,3$$

177
178  Where $\phi$ reflects the fraction of noise that is correlated across units, which we
179  refer to in the text as noise correlations. Noise correlations ($\phi$) were manipulated
180  across values ranging from 0 to 0.2 for simulations. Note that, since $\phi$ appears in
181  the denominator of equation 2, adding noise correlations while sustaining a fixed
182  population signal-to-noise ratio leads to lower variance in the firing rates of single
183  neurons, differing from previous theoretical assumptions (compare figure 2a&b).
184
185  The input layer of the neural network was fully connected to an output layer
186  composed of two output units representing left and right responses. Output units
187  were activated on a given trial according to a weighted function of their inputs:
188
189

$$\boldsymbol{F_{output}} \; = \; \boldsymbol{wF_{input}} \quad Eq.\,4$$

190
191  Where $F_{output}$ is a vector of firing rates of output units, $F_{input}$ is a vector of firing
192  rates of the input units, and w is the weight matrix. Firing of an individual output
193  unit can also be written as a weighted sum over input unit activity:
194

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

195  where $F_j$ reflects the firing of the j[th] output unit, $F_i$ reflects the firing of the i[th] input
196  unit, and $w_{i,j}$ reflects the weight of the connection between the i[th] input unit and

5

197 the j$^{th}$ output unit. Actions were selected as a softmax function of output firing
198 rates:
199

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

200 where $\beta$ is an inverse temperature, which was set to a relatively deterministic
201 value (10000). Learning was implemented through reinforcement learning of
202 weights to the selected output neuron (subscripted j below):
203

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

204 Where $F_i$ is the normalized firing rate of the i$^{th}$ input neuron, $\delta$ is the reward
205 prediction error experienced on a given trial [+0.5 for correct trials and -0.5 for
206 error trials], and $\alpha$ is a learning rate (set to 0.0001 for simulations in figure 2). The
207 network was trained to correctly identify two stimuli (each of which was preferred
208 by a single pool of input neurons) over 100 trials (the last 20 trials of which were
209 considered testing). Simulations were repeated 1000 times for each level of $\phi$
210 and performance measures were averaged across all repetitions. Mean accuracy
211 per trial across all simulations was convolved with a Gaussian kernel (standard
212 deviation = 0.5 trials) for plotting in figure 2b. Mean accuracy across the final 20
213 trials was used as a measure of final accuracy (figure 2e). Statistics on model
214 performance were computed as Pearson correlations between noise correlations
215 $\phi$ and performance measures across all simulations and repetitions.
216

217 *Hebbian learning of noise correlations in three layer network*

218
219 We extended the two-layer feed-forward architecture described above to include
220 a third hidden layer in order to test whether Hebbian learning could facilitate
221 production of noise correlations among similarly tuned neurons (figure 4A). The
222 input layer was fully connected to the hidden layer, and each layer contained 200
223 neurons. In the input layer, neurons were tuned (100 leftward, 100 rightward) as
224 described above, with $\phi$ set to zero (eg. no noise correlations). Weights to the
225 hidden layer were initialized to favor one-to-one connections between input layer
226 units and hidden layer units by adding a small normal random weight perturbation
227 (mean=0, standard deviation = 0.01) to an identity matrix. During learning,
228 weights between the input and hidden layer were adjusted according to a
229 normalized Hebbian learning rule:
230

$$\Delta W = \alpha_{hebb} F'_1 F_2 \quad Eq.\,8$$

6

231

232  Where $F'_1$ is a normalized vector of firing rates corresponding to the input layer
233  and $F_2$ is a normalized vector of firing rates corresponding to the hidden layer
234  units. The learning rate for Hebbian plasticity ($\alpha_{hebb}$) was set to 0.00005 for
235  simulations in figure 4. The model was "trained" over 100 trials in the same
236  perceptual discrimination task described above and an additional 100 trials of the
237  task were completed to measure emergent noise correlations in the hidden layer.
238  Noise correlations were measured by regressing out variance attributable to the
239  stimulus on each trial, and then computing the Pearson correlation of residual
240  firing rate across each pair of neurons for the 100 testing trials (figure 4B&C).

241

242  *Learning readout in multiple discrimination task*

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

250

251  In order to model this task we extended our two-layer feed-forward network to
252  include 4 populations of input units, 4 output units, and 2 task units (figure 5A).
253  Each population of 100 input units encoded a conjunction of the movement
254  directions (up-right, up-left, down-right, down-left). On each trial, the mean firing
255  rate of each input unit population was determined according to their tuning
256  preferences:

257

258

$$\mu = V + H \quad Eq.9$$

259

260  Where V was +1/-1 for trials with the preferred/anti-preferred vertical motion
261  direction H was +1/-1 for trials with the preferred/anti-preferred horizontal motion
262  direction. Firing rates for individual neurons were sampled from a multivariate
263  Gaussian distribution with mean $\mu$ and a covariance matrix that depended on trial
264  type (vertical versus horizontal) and the level of same pool, relevant pool, and
265  irrelevant pool correlations.

266

267  In order to create a covariance matrix, we stipulated a desired standard error of
268  the mean for summed population activity (SEM=20 for simulations in figure 5)
269  and determined the summed population variance that would correspond to that
270  value ($\sigma^2_{pool}$). We then determined the variance on individual neurons that would
271  yield this population response under a given noise correlation profile as follows:

7

272

$$\sigma^2_{unit} \; = \; \frac{\sigma^2_{pool}}{n \, + \, n(n-1)\phi_{same} \, + \, n^2\phi_{relevant} \, - \, n^2\phi_{irrelevant}} \quad Eq.\,10$$

273

274 Where $\phi_{same}$ is the level of same pool correlations (range: 0-0.2 in our
275 simulations), $\phi_{relevant}$ is the level of relevant pool correlations (range: 0-0.2 in our
276 simulations), $\phi_{irrelevant}$ is the level of irrelevant pool correlations (range: 0-0.2 in
277 our simulations. Note that increasing the same pool or in pool correlations
278 reduces the overall variance in order to preserve the same level of variance on
279 the task relevant dimension in the population response, but that increasing
280 irrelevant pool correlations has the opposite effect. Covariance elements of the
281 covariance matrix were determined as follows:
282

$$Cov(same\ pool) = \; \phi_{same}\sigma^2_{unit} \quad Eq.\,11$$

$$Cov(relevant\ pool) \; = \; \phi_{relevant}\sigma^2_{unit} \quad Eq.\,12$$

$$Cov(irrelevant\ pool) \; = \; \phi_{irrelevant}\sigma^2_{unit} \quad Eq.\,13$$

283 Variance and covariance values above were used to construct a covariance
284 matrix for each trial type (vertical/horizontal) as depicted in figure 1.
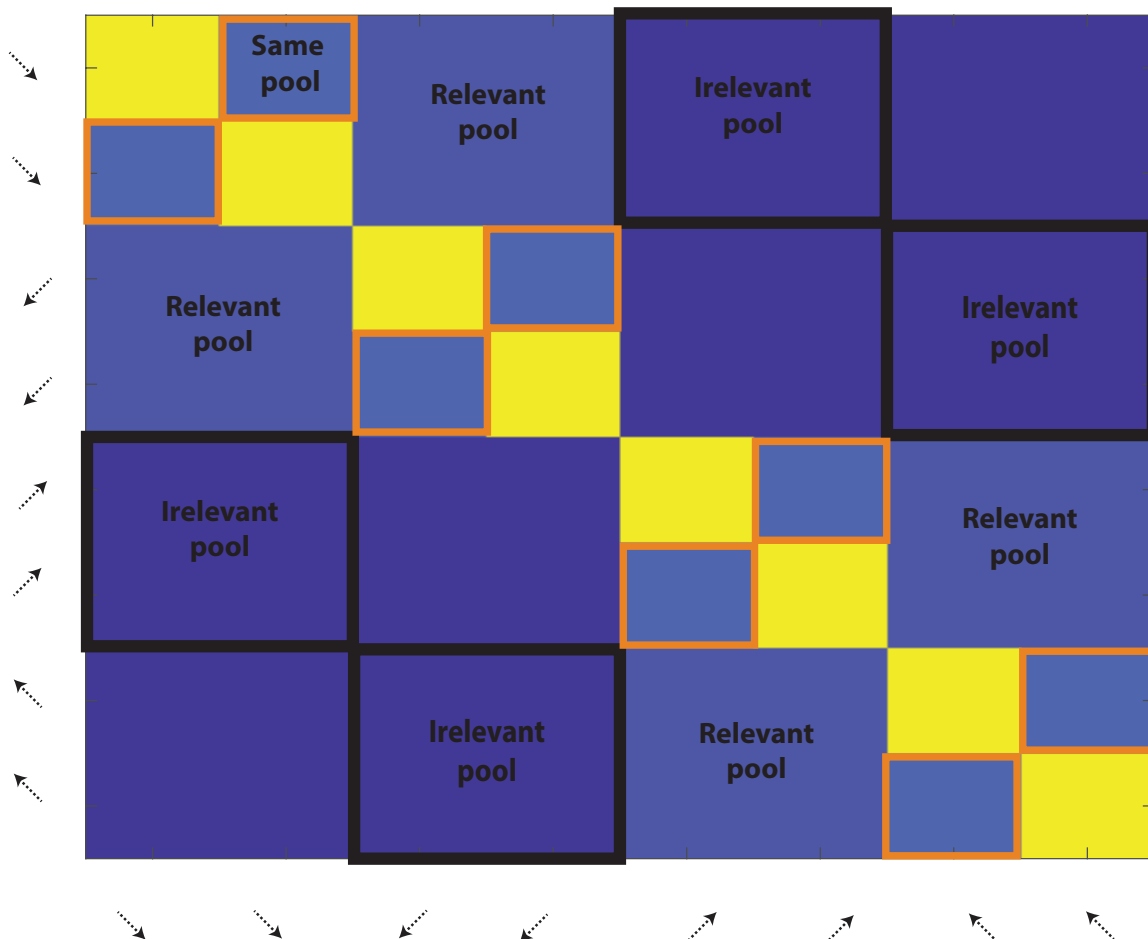
## Covariance matrix: vertical trials



**Figure 1: Schematic of covariance matrix for two-dimensional motion discrimination task.** Same pool correlations are controlled by covariance elements between neurons with identical tuning (orange boxes). Relevant pool correlations are controlled by covariance elements between neurons that are similarly tuned to the task-relevant feature. Task irrelevant correlations are controlled by covariance elements between neurons that are similarly tuned to the task-irrelevant feature. The covariance matrix shown here is for a vertical trial – on a horizontal trial the irrelevant pool and relevant pool locations would be reversed. Covariance elements for pairs of neurons that differed in tuning on both dimensions were set to zero. Each input population has been depicted as two units here for presentation purposes. Background color reflects the case where same pool correlations = 0.2 and relevant pool correlations = 0.1.

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

9

306
307
308 **Results:**
309
310 We examine how noise correlations affect learning in a simplified neural network
311 where the appropriate readout of hundreds of weakly tuned units is learned over
312 time through reinforcement. In order to isolate the effects of noise correlations on
313 learning, rather than their effects on other factors such as representational
314 capacity, we consider population encoding schemes at the input layer that can be
315 constrained to a fixed signal-to-noise ratio. This assumption differs from previous
316 work on noise correlations where the *variance* of the neural population is
317 assumed to be fixed and covariance is changed to produce noise correlations,
318 thereby affecting the representational capacity of the population (figure 2A;
319 (Averbeck et al., 2006; Moreno-Bote et al., 2014)). Under our assumptions, a
320 fixed signal-to-noise ratio can be achieved for any level of by scaling the variance
321 (figure 2B; equations 1-3), or, alternately scaling the magnitude of the signal (not
322 shown). While we do not discount the degree to which noise correlations affect
323 the encoding potential of neural populations, we believe that in many cases the
324 relevant information is limited by extrinsic factors (eg. the stimulus itself, or
325 upstream neural populations providing input (Beck et al., 2012; Kanitscheider et
326 al., 2015)). Under such conditions, reducing noise correlations can increase
327 information only until it saturates because all of the available incoming
328 information is encoded. Beyond that, increasing encoding potential is not
329 possible as it would be tantamount to the population "creating new information"
330 that was not communicated by inputs to the population. Therefore, our framework
331 can be thought of as testing how best to format limited available information in a
332 neural population in order to ensure that an acceptable readout can be rapidly
333 and robustly learned.
334
335 We propose that within this framework, noise correlations of the form that have
336 previously been shown to limit encoding are beneficial because they constrain
337 learning to occur over the most relevant dimensions. In general, a linear readout
338 can be thought of as hyperplane serving as a classification boundary in an N
339 dimensional space, where N reflects the number of neurons in a population.
340 Learning in such a framework involves adjustments of the hyperplane to minimize
341 classification errors. The most useful adjustments are in the dimension that best
342 discriminates signal from noise (central arrows in figure 2C&D), but adjustments
343 may also occur in dimensions orthogonal to the relevant one (such as "twisting"
344 of the hyperplane depicted by curved arrows in figure 2C&D) that could
345 potentially impair performance, or slow down learning. Our motivating hypothesis
346 is that by focusing population activity into the task relevant dimension, noise
347 correlations can increase the fraction of hyperplane adjustments that occur in the
348 task relevant dimension (figure 2D), thus reducing the effective dimensionality of
349 readout learning.
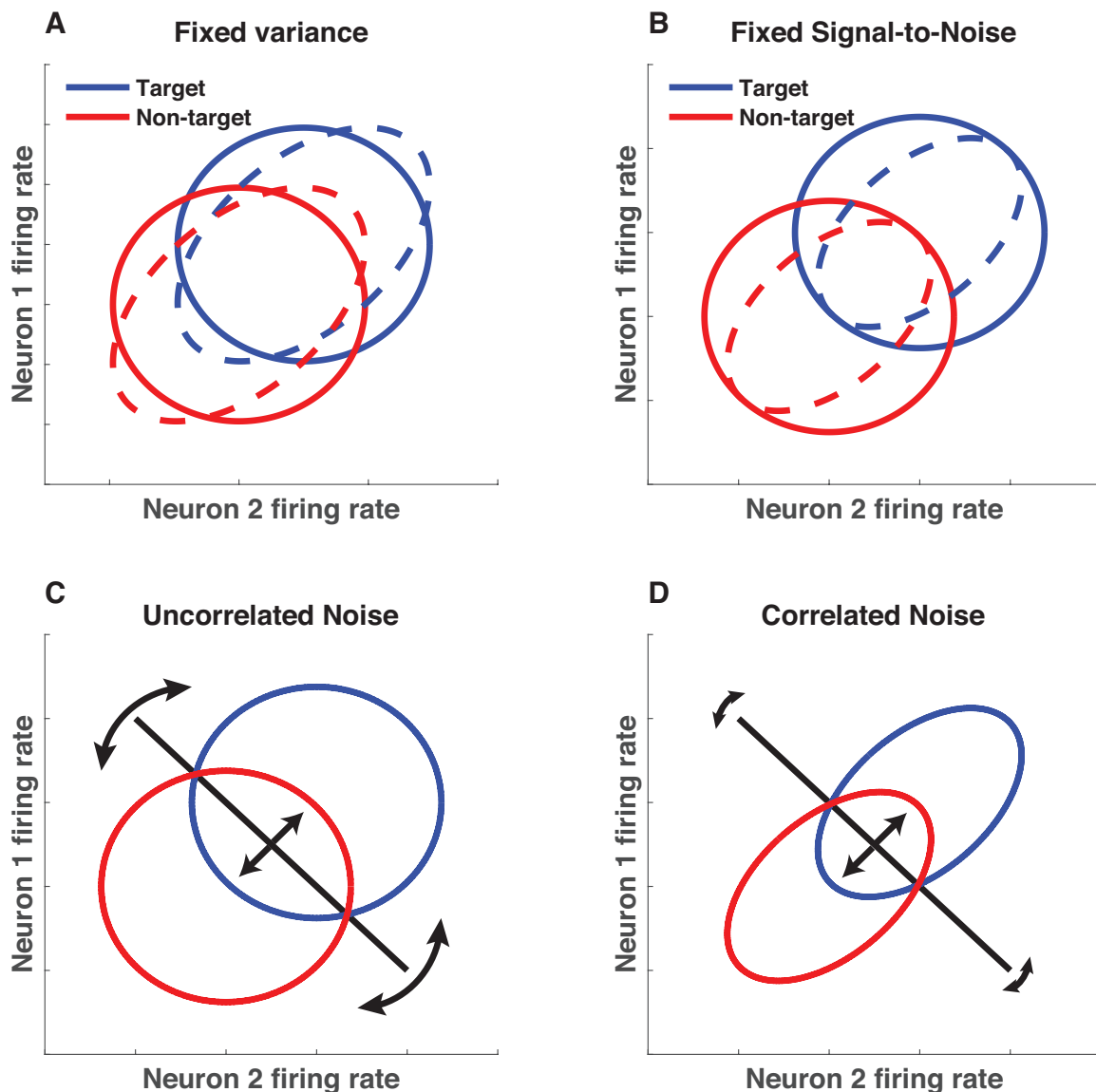
350
351
352
353



354
355
356 **Figure 2: Modeling noise correlations in under extrinsic constraint on signal-to-noise**
357 **ratio. A)** Previous work has modeled noise correlations by assuming that population variance is
358 fixed and that covariance is manipulated to produce noise correlations. Under such assumptions,
359 the firing rate of two similarly tuned neurons is plotted in the absence (solid) or presence (dotted)
360 of information-limiting noise correlations. **B)** Here we assume that the signal-to-noise ratio of the
361 neural population is limited to a fixed value such that noise correlations between similarly tuned
362 neurons do not affect theoretical performance. Thus, the percent overlap of blue (target) and red
363 (non-target) activity profiles does not differ in the presence (dotted) or absence (solid) of noise
364 correlations. **C&D)** Under this assumption, noise correlations among similarly tuned neurons
365 could compress the population activity to a plane orthogonal to the optimal decision boundary,

11

366 thereby minimizing boundary adjustments in irrelevant dimensions (**C**) and maximizing boundary
367 adjustments on relevant ones (**D**).
368
369
370 In order to test this hypothesis, we constructed a fully connected two-layer feed-
371 forward neural network in which input layer units responded to one of two
372 stimulus categories (pool 1 & pool 2) and each output unit produced a response
373 consistent with a category perception (left/right units in figure 3A). On each trial,
374 the network was presented with one stimulus at random, and input firing for each
375 pool was drawn from a multivariate Gaussian with a covariance that was
376 manipulated while preserving the population signal-to-noise ratio. Output units
377 were activated according to a weighted average of inputs and a response was
378 selected according to output unit activations. On each trial, weights to the
379 selected action were adjusted according to a reinforcement learning rule that
380 strengthened connections that facilitated a rewarded action and weakened
381 connections that facilitated an unrewarded action (Law and Gold, 2009).
382
383 Noise correlations led to faster and more robust learning of the appropriate
384 stimulus-response mapping. All neural networks learned to perform the requisite
385 discrimination, but neural networks that employed correlations among similarly
386 tuned neurons learned more rapidly (figure 3B). After learning, networks that
387 employed such noise correlations assigned more homogenous weights to input
388 units of a given pool than did networks that lacked noise correlations (compare
389 figure 3C&D). This led to better trained-task performance (figure 3E; Pearson
390 correlation between noise correlations and test performance: $R = 0.29$, $p < 10e$-
391 $50$) and greater robustness to adversarial noise profiles (figure 3F; $R = 0.81$, $p <$
392 $10e$-$50$) in the networks that employed noise correlations. Critically, these
393 learning advantages emerged despite the fact that optimal readout of all
394 networks achieved similar levels of performance and robustness (figure 3E&F,
395 compare optimal readout across conditions).
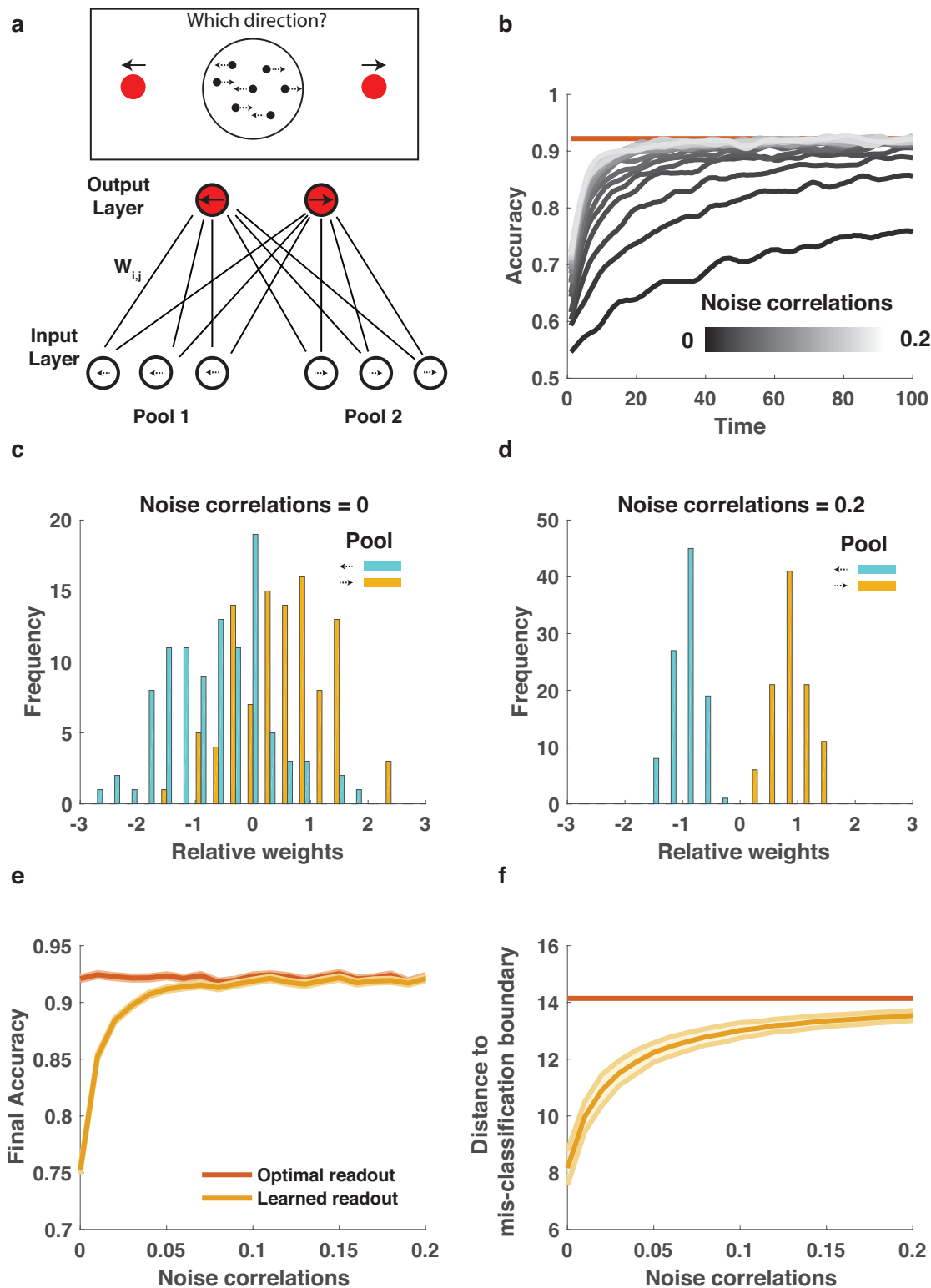396
397
398

**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

405    perceptual threshold. Input layer contains two pools of neurons that provide evidence for alternate
406    percepts (eg. leftward motion versus rightward motion) and output neurons encode alternate
407    courses of actions (eg. saccade left versus saccade right). Layers are fully connected with
408    weights randomized to small values and adjusted after each trial according to rewards (see
409    methods). **B)** Average learning curves for neural network models in which population signal-to-
410    noise ratio in pools 1 and 2 were fixed, but noise correlations (grayscale) were allowed to vary
411    from small (dark) to large (light) values. **C&D)** Weight differences (left output – right output) for
412    each input unit (color coded according to pool) after 100 timesteps of learning for low (**C**) and high
413    (**D**) noise correlations. **E)** Accuracy in the last 20 training trials is plotted as a function of noise
414    correlations for learned readouts (orange) and optimal readout (red). Lines/shading reflect
415    Mean/SEM. F) The shortest distance, in terms of neural activation, required to take the mean
416    input for a given category (eg. left or right) to the boundary that would result in misclassification is
417    plotted for the final learned (orange) and optimal (red) weights for each noise correlation condition
418    (abscissa). Lines/shading reflect Mean/SEM.
419
420
421

422    Given that noise correlations implemented in our previous simulation, like those

423    observed in the brain, depended on the tuning of individual units, we tested

424    whether such noise correlations might be produced via Hebbian plasticity.

425    Specifically, we considered an extension of our neural network in which an

426    additional intermediate layer is included between input and output neurons (figure

427    4a). Input units were again divided into two pools that differed in their encoding,

428    but variability was uncorrelated across neurons within a given pool. Connections

429    between the input layer and intermediate layer were initialized such that each

430    input unit strongly activated one intermediate layer unit, and shaped over time

431    using a Hebbian learning rule that strengthened connections between co-

432    activated neuron pairs. Despite the lack of noise correlations in the input layer of

433    this network (figure 4b; mean[std] in pool residual correlation = 0.0015[0.10]),

434    neurons in the intermediate layer developed tuning-specific noise correlations of

435    the form that were beneficial for learning in the previous simulations (figure 4c;

436    mean[std] in pool residual correlation = 0.55[0.07]; $t$-test on difference from input

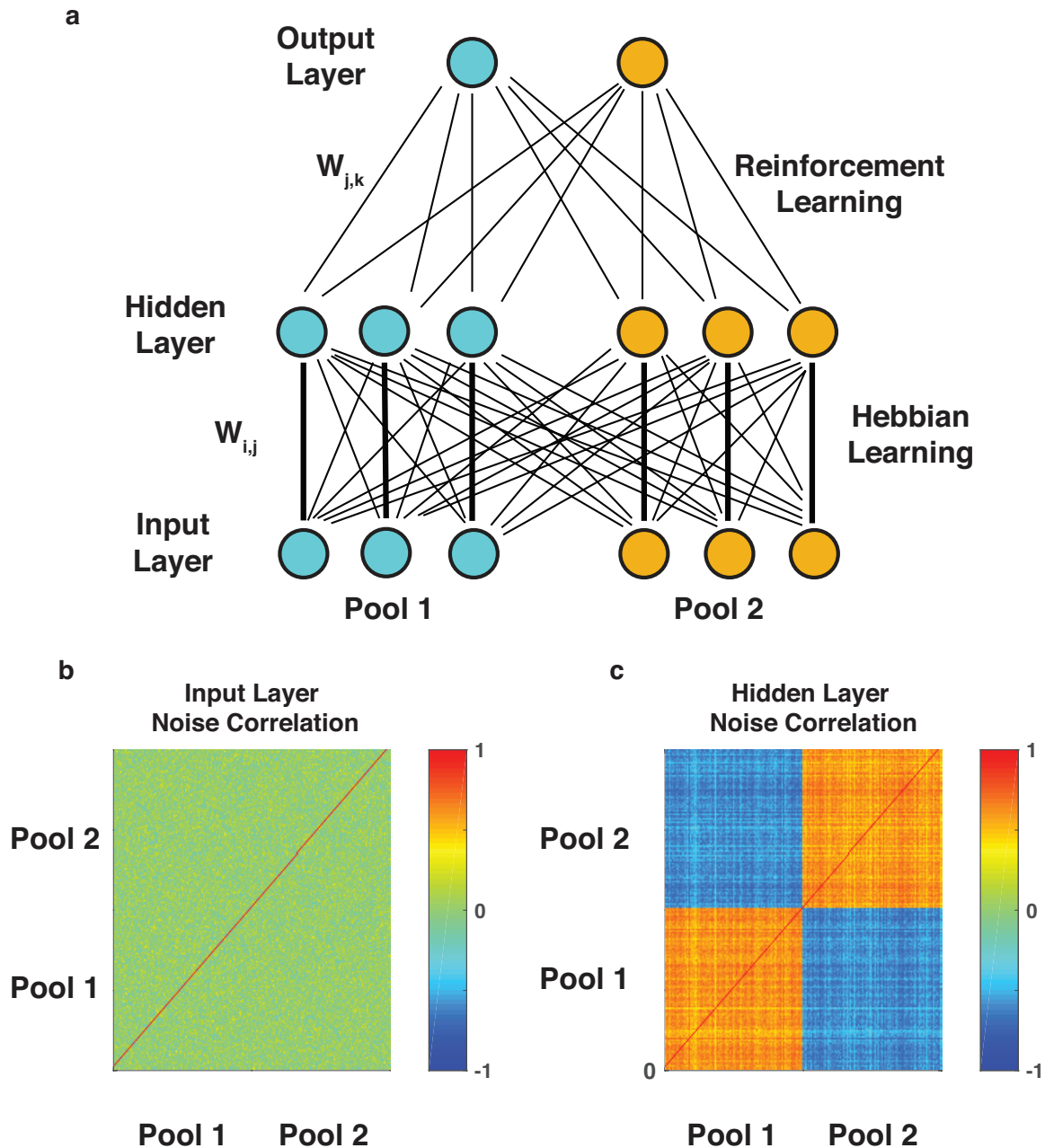437    layer correlations $t = 443$, $dof = 19800$, $p < 10e\text{-}50$).

438
439
440
441
442
443

444
445
446
447
448 **Figure 4: Hebbian learning produces correlations within similarly tuned populations in a**
449 **perceptual discrimination task. A)** Three-layer neural network architecture. Input layer feeds
450 forward to hidden layer, which is fully connected to an output layer. Input layer provides
451 uncorrelated inputs to hidden layer through projection weights that are adjusted according to a
452 Hebbian learning rule. **B&C)** Noise correlations observed in hidden layer units at the beginning
453 (**B**) and end (**C**) of training.
454
455
456
457

458

459    In order to understand how noise correlations might impact learning in mixed
460    encoding populations, we extended our perceptual discrimination task to include
461    two directions of motion discrimination (eg. up/down and left/right). On each trial,
462    a cue indicated which of two possible motion discriminations should be
463    performed (figure 5A, left; (Cohen and Newsome, 2008)). We extended our
464    neural network to include four populations of one hundred input units, each
465    population encoding a conjunction of motion directions (up-right, up-left, down-
466    right, down-left; figure 5A; input layer). Two additional inputs provided a perfectly
467    reliable "cue" regarding the relevant feature for the trial (figure 5A; task units).
468    Four output neurons encoded the four possible responses (up, left, down, right)
469    and were fully connected to the input layer (figure 5A; output layer). Task units
470    were hard wired to eliminate irrelevant task responses, but weights of input units
471    were learned over time as in our previous simulations.

472

473    Learning performance in the two-feature discrimination task depended not only
474    on the level of noise correlations, but also on the type. As in the previous
475    simulation, adding noise correlations to each individual population of identically
476    tuned units led to faster learning of the appropriate readout (Figure 5B&C,
477    compare blue and yellow; Figure 5D&E, vertical axis; mean[std] accuracy across
478    training: 0.53[0.05] and 0.614[0.08] for minimum (0) and maximum (0.2) in pool
479    correlations, t-test for difference in accuracy: $t = 95$, $dof = 19998$, $p < 10e\text{-}50$).

480

481    However, the more complex task design also allowed us to test whether dynamic
482    trial-to-trial correlations might further facilitate learning. Specifically, correlations
483    that increase shared variability among units that contribute evidence to the same
484    response have been observed previously (Cohen and Newsome, 2008), and
485    could in principle focus learning on relevant dimensions (figure 2C&D) even when
486    those dimensions change from trial to trial. Indeed, adding correlations among
487    separate pools that share the same encoding of the relevant feature (eg. UP on a
488    vertical trial) led to faster learning (figure 5B; mean[std] training accuracy for
489    model with relevant pool correlations: 0.64[0.09], $t$-test for difference from in pool
490    correlation only model: $t = 22$, $dof = 19998$, $p < 10e\text{-}50$) and weights that more
491    closely approached the optimal readout (figure 5E, horizontal axis). In contrast,
492    when positive noise correlations were introduced across separate encoding pools
493    that shared the same tuning for the irrelevant dimension on each trial (eg. UP on
494    a horizontal trial) learning was impaired dramatically (figure 5C; mean[std]
495    training accuracy for model with irrelevant pool correlations: 0.51[0.05], $t$-test for
496    difference from in pool correlation only model: $t = -112$, $dof = 19998$, $p < 10e\text{-}50$)
497    and learned weights diverged from the optimal readout (figure 5F, horizontal
498    axis). Model performance differences were completely attributable to learning the
499    readout, as all models performed similarly when using the optimal readout (figure
500    S1).

501

16

502  In order to test the idea that noise correlations might focus learning onto relevant
503  dimensions, we extracted weight updates from each trial and projected these
504  updates into a two-dimensional space where the first dimension captured the
505  relative sensitivity to leftward versus rightward motion and the second dimension
506  captured relative sensitivity to upward versus downward motion. In the model
507  where input units were only correlated within their identically tuned pool, weight
508  updates projected in all directions more or less uniformly (figure 5G), and did not
509  differ systematically across trial types (vertical versus horizontal). However,
510  dynamic noise correlations that shared variability across the relevant dimension
511  tended to push weight updates onto the appropriate dimension for a given trial
512  (figure 4F; $t$-test for difference in the magnitude of updating in up/down and
513  left/right dimensions across conditions [up/down – left/right]: $t = 3.4$, $dof$=98, $p =$
514  0.001). In contrast, dynamic noise correlations that shared variability across the
515  irrelevant dimension tended to push weight updates onto the wrong dimension
516  (figure 4H; t-test for difference in the magnitude of updating in up/down and
517  left/right dimensions across conditions [up/down – left/right]: $t = -9.5$, $dof$=98, $p =$
518  10e-14). Both of these trends were consistent across simulations, providing an
519  explanation for the performance improvements achieved by relevant noise
520  correlations (projection of learning onto an appropriate dimension) and
521  performance impairments produced by irrelevant noise correlations (projection of
522  learning onto an inappropriate dimension).
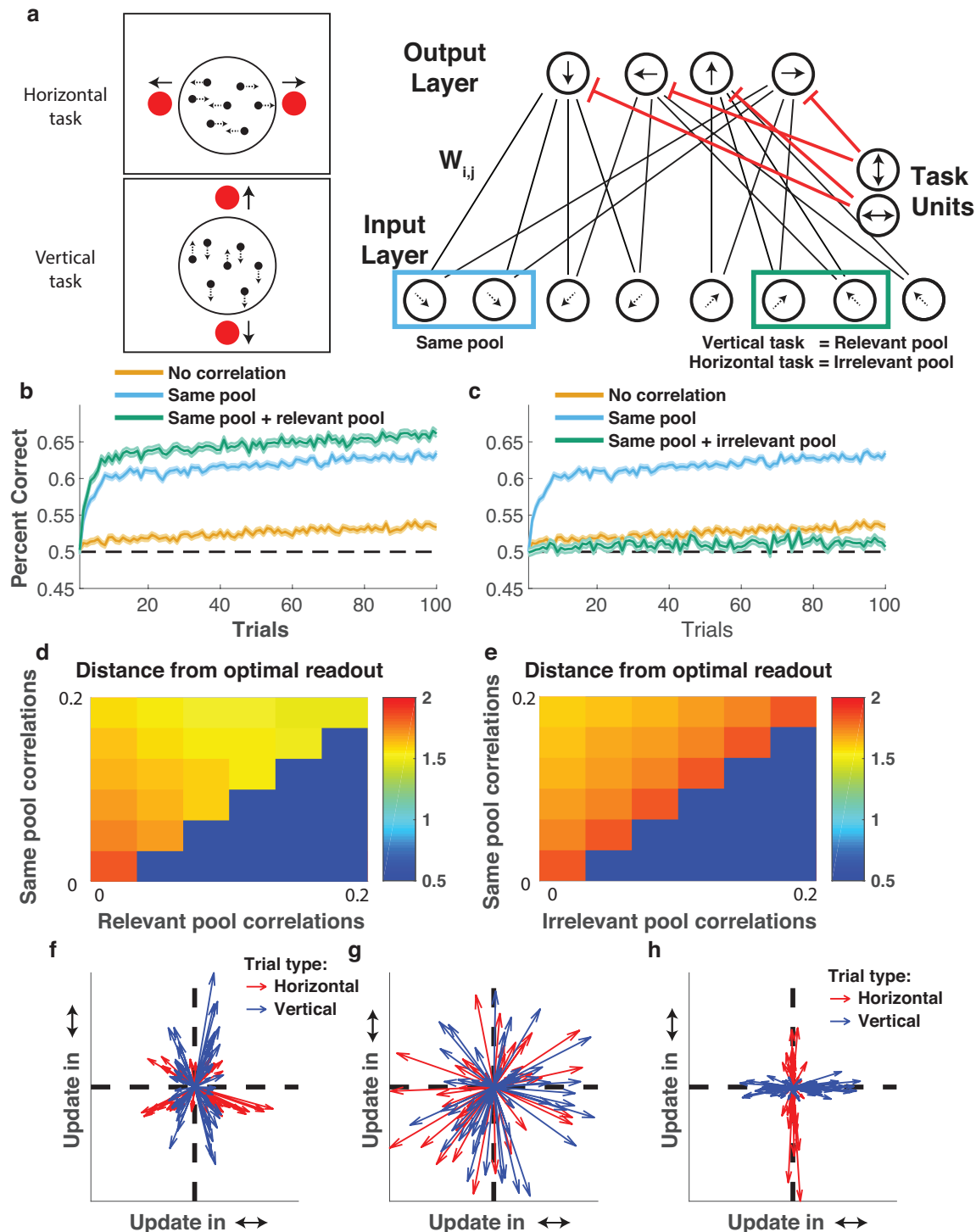523
524
525
526
527
528
529
530

**Figure 5: 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 population of input units represents two dimensions of motion (up versus down, and left versus right), and output units produce responses in favor of alternative actions (up, down, left, right). Two additional input units provide cue

540    information that biases output units to produce an output corresponding to the discrimination
541    appropriate on this trial (eg. horizontal or vertical). Noise correlations were manipulated among 1)
542    identically tuned neurons (blue rectangle; same pool), 2) neurons that have similar encoding of
543    the task relevant feature (green rectangle pair in vertical trials; relevant pool), and 3) neurons that
544    have similar encoding of the task irrelevant feature (green rectangle pair in horizontal trials;
545    irrelevant pool). **B&C**) Learning curves showing accuracy (ordinate) over trials (abscissa) for
546    models 1) lacking noise correlations (orange), 2) containing noise correlations that are limited to
547    neurons that have same tuning for both features (same pool; blue), 3) containing same pool noise
548    correlations along with correlations between neurons in different pools that have the same tuning
549    for the task-relevant feature (in pool+rel pool; green in **B**), and 4) containing in-pool noise
550    correlations along with correlations between neurons in different pools that have the same tuning
551    for the task irrelevant feature (in pool+irrel pool; green in **C**). **D&E**) Distance between learned
552    weights and the optimal readout (color) for models that differ in their level of "in pool" correlations
553    (ordinate, both plots), "relevant pool" correlations (abscissa, **D**), and "irrelevant pool" correlations
554    (abscissa, **E**). **F,G,H**) Weight updates for example learning sessions were projected into a two
555    dimensional space in which net updates to the relative contribution of vertical motion information
556    (eg. up versus down) is represented on the abscissa and updates to the relative contribution of
557    horizontal motion information (eg. left versus right) is represented on the ordinate. Arrows reflect
558    single trial weight updates and are colored according to the trial type (red = horizontal
559    discrimination, blue = vertical discrimination). Weight updates for a model with only "in pool"
560    correlations look similar across trial types (**G**), but weight updates for a model with "relevant pool"
561    correlations indicate more weight updating on the relevant feature (**F**), whereas the opposite was
562    observed in the case of "irrelevant pool" correlations (**H**).
563
564

565    **Discussion:**
566
567    Taken together, our results suggest that in settings where the population signal-
568    to-noise ratio is limited by external factors (eg. inputs) and relevant task
569    representations are low dimensional, noise correlations can make learning faster
570    and more robust by focusing learning on the most relevant dimensions. We
571    demonstrate this basic principle in a simple perceptual learning task (figure 3),
572    where beneficial noise correlations between similarly tuned units could be
573    produced through a simple Hebbian learning rule (figure 4). We extended our
574    framework to a contextual learning task to demonstrate that dynamic noise
575    correlations that bind task relevant feature representations facilitate faster
576    learning (figure 5b&d) by pushing learning onto task-relevant dimensions (figure
577    5f). Given the pervasiveness of noise correlations among similarly tuned sensory
578    neurons (Zohary et al., 1994; Maynard et al., 1999; Bair et al., 2001; Averbeck
579    and Lee, 2003; Cohen and Maunsell, 2009; Huang and Lisberger, 2009; Ecker et
580    al., 2010; Gu et al., 2011; Adibi et al., 2013), and that the noise correlations
581    dynamics beneficial for learning in our simulations are similar to those that have
582    been observed *in vivo* (Cohen and Newsome, 2008), we interpret our results as
583    suggesting that noise correlations between similarly tuned neurons are a feature
584    of neural coding architectures that ensures efficient readout learning, rather than
585    a bug that limits encoding potential.
586
587

588   This interpretation rests on several assumptions in our model. Of particular
589   importance is the assumption that signal-to-noise ratio of our populations is fixed,
590   meaning that our manipulation of noise correlations can focus variance on
591   specific dimensions without gaining or losing information. This assumption
592   reflects conditions in which information is limited at the level of the inputs to the
593   population, for instance due to noisy peripheral sensors (Beck et al., 2012;
594   Kanitscheider et al., 2015). In such conditions, even with optimal encoding,
595   population information saturates at an upper bound determined by the
596   information available in the inputs to the population. Therefore, fixing the signal-
597   to-noise ratio enabled us to examine the effect of noise correlations on
598   downstream processes that learn to read-out the population code in the absence
599   of any influence of noise correlations on the quantity of information contained
600   within that population code.
601
602   Previous theoretical work exploring the role of noise correlations in encoding has
603   typically assumed that single neurons have a fixed variance, such that tilting the
604   covariance of neural populations towards or away from the dimension of signal
605   encoding would have a large impact on the amount of information that can be
606   encoded by a population (figure 1a; (Averbeck et al., 2006; Moreno-Bote et al.,
607   2014)). Such assumptions lead to the idea that positive noise correlations among
608   similarly tuned neurons limit encoding potential, raising the question of why they
609   are so common in the brain (Cohen and Kohn, 2011). In considering the
610   implications of this framework, one important question is: if information encoded
611   by the population can be increased by changing the correlation structure among
612   neurons, where does this additional information come from? In some cases, the
613   neural population in question may indeed receive sufficient task relevant
614   information from upstream brain regions to reorganize its encoding in this way,
615   but in other cases it is likely that information is limited by the inputs to a neural
616   population (Kanitscheider et al., 2015; Kohn et al., 2016). In cases where
617   incoming information is limited, further increasing representational capacity is not
618   possible, and formatting information for efficient readout is essentially the best
619   that the population code could do. Here we show that the noise correlations that
620   have previously been described as "information limiting" are exactly the type of
621   correlations that format information most efficiently for readout learning under
622   such conditions.
623
624   Jointly considering these antagonistic perspectives on noise correlations provides
625   a more nuanced view of how neural representations are likely optimized for
626   learning. In order to optimize an objective function, a neural population can
627   reduce correlated noise in task relevant dimensions to increase its
628   representational capacity up to some level constrained by its inputs (Figure 6,
629   left). But once the population is fully representing all task relevant information that
630   has been provided to it, it can additionally optimize representations by pushing as
631   much variance onto task relevant dimensions as possible, thereby affording

20

632  efficient learning in downstream neural populations (Figure 6, right). In short,
633  optimization of a neural population code does not occur in a vacuum, and instead
634  depends critically on both upstream (eg. input constraints) and downstream (eg.
635  readout) neural populations (Figure 6). In this view, if a neural population is *not*
636  fully representing the decision relevant information made available to it, then
637  learning could improve the efficiency of representations by reducing rate limiting
638  noise correlations as has been observed in some paradigms (Figure 6, left; Gu et
639  al., 2011; Ni et al., 2018). In contrast, once available information is fully
640  represented, readout learning could be further optimized by reformatting
641  population codes such that variability is shared across neurons with similar
642  tuning for the relevant task feature, producing the sorts of dynamic noise
643  correlations that have been observed in well trained animals (Figure 6, right;
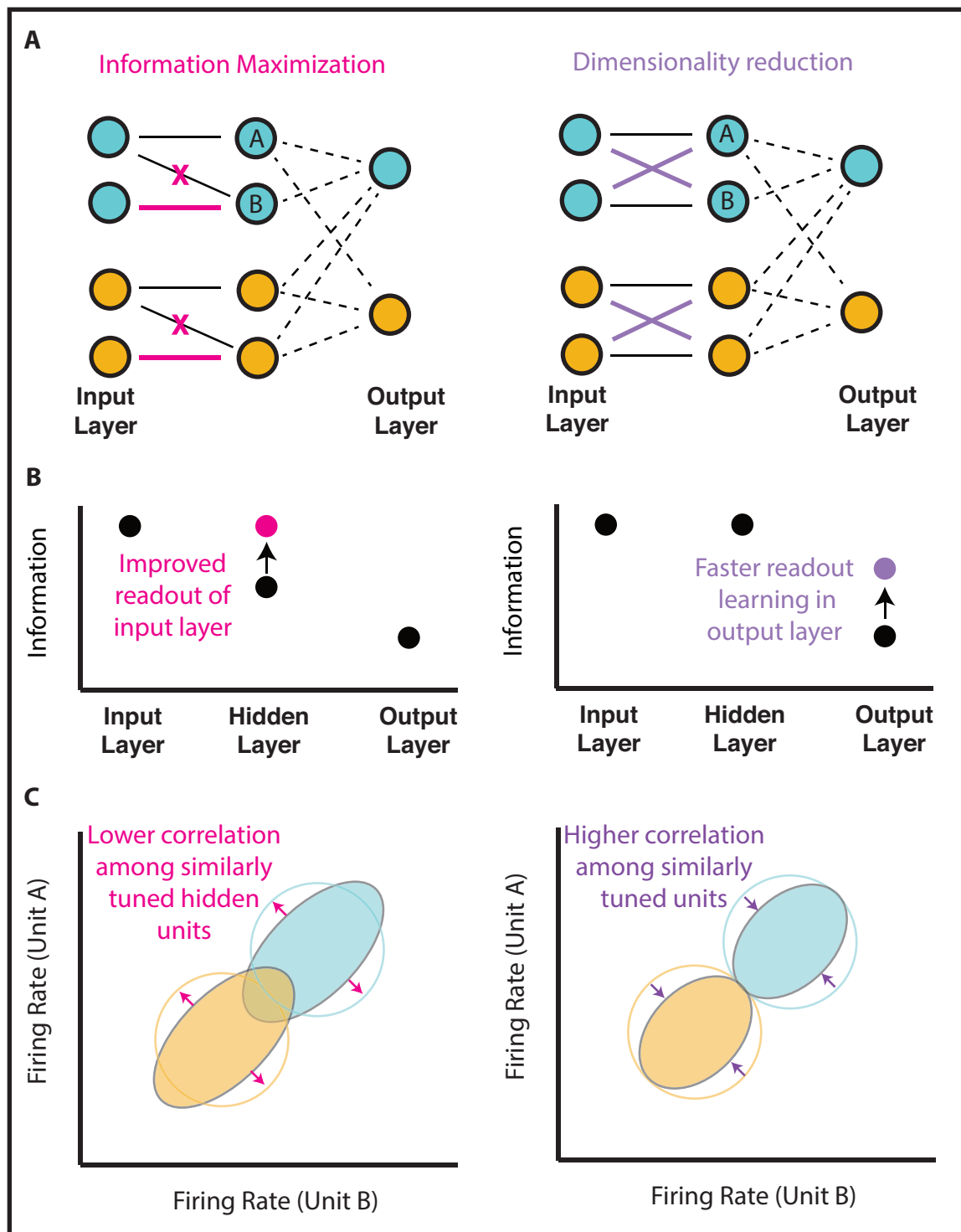644  Cohen and Newsome, 2008).

645
646 **Figure 6: Information maximization and dimensionality reduction can be useful for**
647 **learning under different situations and have opposite effects on noise correlations among**
648 **similarly tuned units. A)** A schematic representation of a three layer neural network in which
649 units provide evidence for one of two categorizations (blue/orange). In the left network, the hidden
650 layer initially has access to information from only one of two independent units in each pool, but
651 weights are subsequently adjusted to increase task-relevant information represented in the
652 hidden layer (pink). In the right network, the hidden layer initially has access to all task-relevant

653 information, but weights are subsequently adjusted to share signal and noise across similarly
654 tuned units to afford dimensionality reduction (purple). Note that the information maximizing
655 weight adjustments (left, pink) increase signal-to-noise ratio in the hidden layer but preserve the
656 variance in firing rate of individual neurons, whereas the dimensionality reducing weight
657 adjustments (right, purple) maintain a fixed signal-to-noise ratio in hidden units, but decrease the
658 variance of individual units by averaging across multiple similarly tuned inputs. Dashed lines to
659 output units reflect weights that need to be learned based on feedback. **B)** Task relevant
660 information (mutual information between unit activations and stimulus category; abscissa) is
661 depicted for each layer (ordinate). Weight adjustments affording information maximization (left)
662 increase task relevant information in the hidden layer (pink), whereas weight adjustments that
663 afford dimensionality reduction (right) do not affect task-relevant information in the hidden layer
664 itself but instead increase the rate of learning in the output layer, thereby leading to more task-
665 relevant information in the output layer (purple). **C)** Weight adjustments for information
666 maximization (pink in panel A) *decrease* correlations among hidden units A&B by removing
667 shared input from a single input unit and instead providing independent sources of input to each
668 unit (pink arrows). In contrast, weight adjustments for dimensionality reduction *increase* noise
669 correlations among hidden units A&B by providing them with the same mixture of information from
670 the two identically tuned input units. We propose that both of these processes play a critical role
671 in learning and that changes in noise correlations across learning will depend critically on which
672 process dominates. As shown in panel B, this will depend critically on whether the neural
673 population in question has already fully represented information available from its inputs. In
674 principle, these processes could occur serially, with early learning maximizing information
675 available in intermediate layers (left) and later learning compressing that information into a format
676 allowing rapid readout learning (right).
677
678
679 In addition to key assumptions about an external limitation on signal-to-noise, our
680 modeling included a number of simplifying assumptions that are unlikely to hold
681 up in real neural populations. For example, we consider discrete pools of
682 identically tuned neurons, rather than the heterogeneous populations observed in
683 sensory cortical regions of the brain. A primary goal of our work was to identify
684 the computational principles that control the speed at which readout can be
685 learned, and our simplified populations are considerably more tractable and
686 transparent than realistic neural populations. The principles that we identify here
687 are certainly at play in real neural populations, albeit with implications that are far
688 less transparent. We hope that our simplified results pave the way for future work
689 to assess nuances that can emerge in mixed heterogeneous populations, or in
690 more realistic architectures that go beyond the simple feed forward flow of
691 information considered here.
692
693 *Model predictions*
694
695 Our work shows that noise correlations can focus the gradient of learning onto
696 the most appropriate dimensions. Thus, our model predicts that the degree to
697 which similarly tuned neurons are correlated during a perceptual discrimination
698 should be positively related to performance improvements experienced on
699 subsequent discriminations. In contrast, our model predicts that the degree of
700 correlation between neurons that are similarly tuned to a task irrelevant feature
701 should control the degree of learning on irrelevant dimensions, and thus

702 negatively relate to performance improvements on subsequent discriminations.
703 These predictions are strongest for the earliest stages of learning where weight
704 adjustments are critical for subsequent performance, but they may also hold for
705 later stages of learning, when correlations on irrelevant dimensions, including
706 independent noise channels, could potentially lead to systematic deviations from
707 optimal readout (figure 2f, 4d&e). These predictions could be tested by recording
708 neural responses to a stimulus set that differs across multiple features to
709 characterize both signal-to-noise and correlated variability for each feature
710 discrimination. A strong prediction of our model is that correlated variability within
711 neurons tuned to a given feature should be a predictor of subsequent learning of
712 responses to that feature – above and beyond feature value discriminability.
713
714 One interesting special case involves tasks where the relevant dimension
715 changes in an unsignaled manner (Birrell and Brown, 2000). In such tasks, noise
716 correlations on the previously relevant dimension would, after such an
717 "extradimensional shift", force gradients into a task-irrelevant dimension and thus
718 impair learning performance. Interestingly, learning after extra-dimensional shifts
719 can be selectively improved by enhancing noradrenergic signaling (Devauges
720 and Sara, 1990; Lapiz and Morilak, 2006), which leads to increased arousal
721 (Joshi et al., 2016; Reimer et al., 2016) and decreased cortical pairwise noise
722 correlations in sensory and higher order cortex (Vinck et al., 2015; Joshi and
723 Gold, n.d.). While these observations have been made in different paradigms, our
724 model suggests that the reduction of noise correlations resulting from increased
725 sustained levels of norepinephrine after an extradimensional shift (Bouret and
726 Sara, 2005) could mediate faster learning by expanding the dimensionality of the
727 learning gradients (compare figure 5G to 5F) to consider features that have not
728 been task-relevant in the past.
729
730 *Relation to attentional effects on noise correlations*
731
732 In broad strokes, our finding that manipulation of noise correlations can focus
733 variance on specific dimensions is in line with specific models of attention. In
734 particular, noise reduction in task irrelevant dimensions might be considered in
735 the same light that is often cast on suppression of task irrelevant dimensions by
736 attentional mechanisms (Zanto and Gazzaley, 2009), in particular for purposes of
737 accurate credit assignment (Akaishi et al., 2016; Leong et al., 2017). One
738 possibility is that compressed low-dimensional task representations in higher-
739 order decision regions (Mack et al., 2019) may pass accumulated decision
740 related information back to sensory regions in order to approximate Bayesian
741 inference (Haefner et al., 2016; Bondy et al., 2018; Lange et al., 2018). As task
742 relevant features are learned, such a process would promote noise correlations
743 between neurons coding those relevant features. In other words, noise
744 correlations may reflect a chosen hypothesis about which feature is relevant for
745 predicting outcomes. Such a signal would be beneficial if it could persist (and

746 thus preserve correlations between neurons tuned to the same task relevant
747 feature value) until the time of feedback or reinforcement. Recent work showing
748 strengthened noise correlations between similarly tuned neurons during working
749 memory maintenance suggests that this might very well be the case (Merrikhi et
750 al., 2018).
751
752 One observation that seems at odds with this interpretation is that manipulations
753 of attention that cue a particular location or feature tend to decrease noise
754 correlations among neurons that encode that location or feature (Cohen and
755 Maunsell, 2009; Mitchell et al., 2009; Cohen and Maunsell, 2011; Herrero et al.,
756 2013; Doiron et al., 2016). The effects of attentional cuing on noise correlations
757 are dynamic in that cues change from one trial to the next, and contextual, in that
758 noise correlations are reduced most dramatically among neurons that contribute
759 evidence toward the same response in a manner consistent with increasing the
760 amount of task relevant information in the population code (Ruff and Cohen,
761 2014; Downer et al., 2015). Our model does not account for these attentional
762 effects, as we intentionally constrained the signal-to-noise ratio of our neural
763 populations, thereby eliminating any potential changes in information encoding
764 potential. However, we hope that our work motivates future studies to jointly
765 consider the impacts of noise correlations on both learning and immediate
766 performance in order to better understand the potentially competing imperatives
767 that the brain faces in dynamically controlling the correlation structure of its own
768 representations (see (Haimerl et al., 2019) for one attempt to do so).
769
770
771 *Origins of useful noise correlations*
772
773 One important question stemming from our work is how noise correlations
774 emerge in the brain. This question has been one of longstanding debate, largely
775 because there are so many potential mechanisms through which correlations
776 could emerge (Kanitscheider et al., 2015; Kohn et al., 2016). Noise correlations
777 could emerge from convergent and divergent feed forward wiring (Shadlen and
778 Newsome, 1998), local connectivity patterns within a neural population (Hansen
779 et al., 2012; Smith et al., 2013), or top down inputs provided separately to
780 different neural populations (Haefner et al., 2016). Here we show that static noise
781 correlations that are useful for perceptual learning emerge naturally from Hebbian
782 learning in a feed-forward network. While this certainly suggests that useful noise
783 correlations could emerge through feed forward wiring, it is also possible to
784 consider our Hebbian learning as occurring in a one-step recurrence of the input
785 units, and thus the same data support the possibility of noise correlations through
786 local recurrence. The context dependent noise correlations that speed learning
787 (figure 4), however, would not arise through simple Hebbian learning. Such
788 correlations could potentially be produced through selective top-down signals
789 from the choice neurons, as has been previously proposed (Wimmer et al., 2015;

790    Haefner et al., 2016; Bondy et al., 2018; Lange et al., 2018). Moreover, top-down
791    input may selectively target neuronal ensembles produced through Hebbian
792    learning (Collins and Frank, 2013). While previous work has suggested that such
793    a mechanism could be adaptive for accumulating information over the course of a
794    decision (Haefner et al., 2016), our work demonstrates that the same mechanism
795    could effectively be used to tag relevant neurons for weight updating between
796    trials, making efficient use of top-down circuitry. Haimerl et al. recently made a
797    similar point, showing that stochastic modulatory signals shared across task-
798    informative neurons can serve to tag them for a decoder (Haimerl et al., 2019).
799
800    *Noise correlations as inductive biases*
801
802    Artificial intelligence has undergone a revolution over the past decade leading to
803    human level performance in a wide range of tasks (Mnih et al., 2015). However, a
804    major issue for modern artificial intelligence systems, which build heavily on
805    neural network architectures, is that they require far more training examples than
806    a biological system would (Hassabis et al., 2017). This biological advantage
807    occurs despite the fact that the total number of synapses in the human brain,
808    which could be thought of as the free parameters in our learning architecture, is
809    much greater than the number of weights in even the most parameter-heavy
810    deep learning architectures. Our work provides some insight into why this occurs;
811    correlated variability across neurons in the brain constrain learning to specific
812    dimensions, thereby limiting the effective complexity of the learning problem
813    (figure 5F-G). We show that, for simple tasks, this can be achieved using
814    Hebbian learning rules (figure 4), but that contextual noise correlations, of the
815    form that might be produced through top-down signals (Haefner et al., 2016), are
816    critical for appropriately focusing learning in more complex circumstances. In
817    principle, algorithms that effectively learn and implement noise correlations might
818    reduce the amount of data needed to train AI systems by limiting degrees of
819    freedom to those dimensions that are most relevant. Furthermore, our work
820    suggests that large scale neural recordings in early stages of learning complex
821    tasks might serve as indicators of the inductive biases that constrain learning in
822    biological systems.
823
824    In summary, we show that under external constraints of task-relevant information,
825    noise correlations that have previously been called "rate limiting" can serve an
826    important role in constraining learning to task-relevant dimensions. In the context
827    of previous theory focusing on representation, our work suggests that neural
828    populations are subject to competing forces when optimizing covariance
829    structures; on one hand reducing correlations between pairs of similarly tuned
830    neurons can be helpful to fully represent available information, but increasing
831    correlations among similarly tuned neurons can be helpful for assigning credit to
832    task relevant features. We believe that this view of the learning process not only
833    provides insight to understanding the role of noise correlations in the brain, but

834 opens up the door to better understand the inductive biases that guide learning in
835 biological systems.
836
837
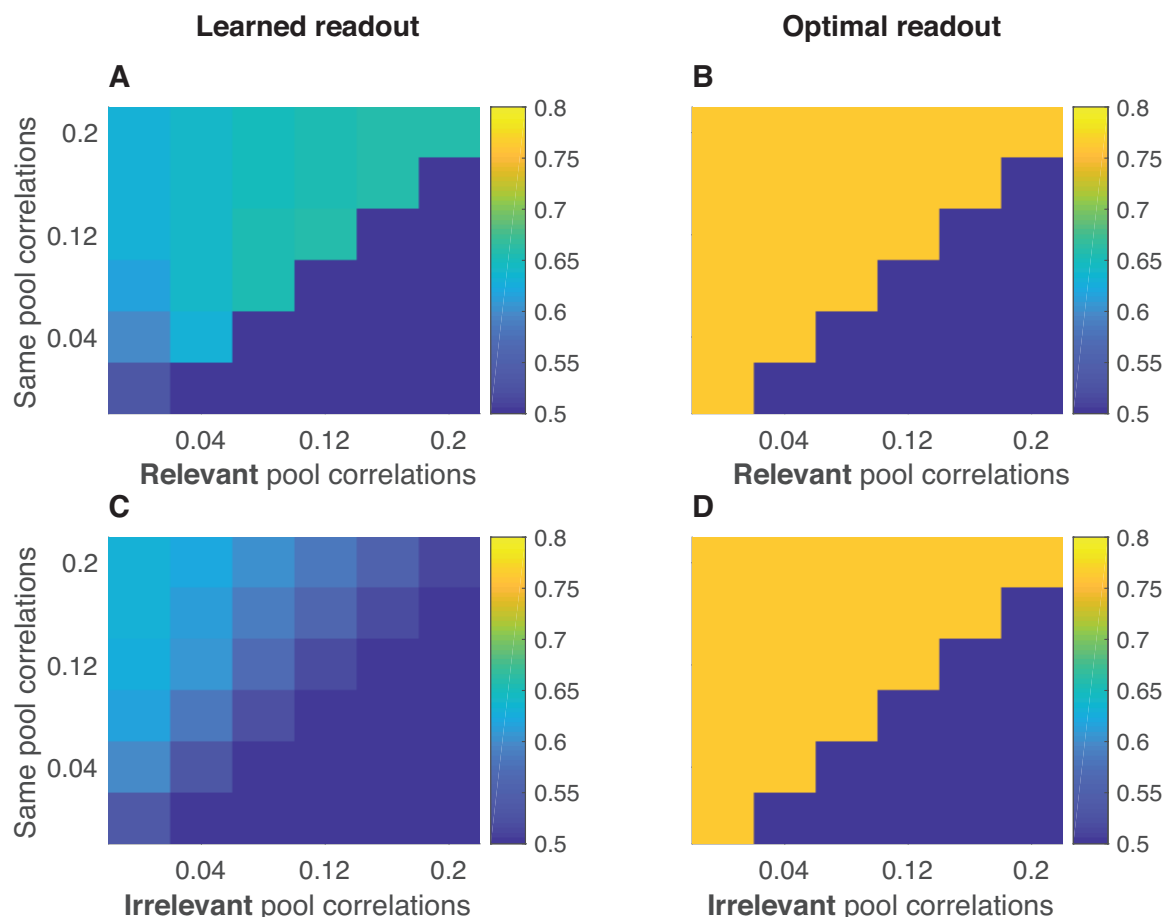838 Supplementary figures:
839
840
841



842
843
844
845 Figure S1: **Noise correlations affect speed of learning, but not performance using optimal**
846 **readout in multiple discrimination task**. **A)** Mean test accuracy (color) of all models spanning
847 the range of in pool correlations (abscissa) and relevant pool correlations (ordinate). **B)** Mean
848 accuracy of same models using optimal readout, rather than the learned readout. **C)** Mean test
849 accuracy (color) of all models spanning the range of in pool correlations (abscissa) and irrelevant
850 pool correlations (ordinate). **D)** Mean accuracy of same models using optimal readout, rather than
851 the learned readout. Note that performance of all models is identical when readout is optimal,
852 rather than learned.
853
854
855
856

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

860  Akaishi R, Kolling N, Brown JW, Rushworth M (2016) Neural Mechanisms of
861      Credit Assignment in a Multicue Environment. Journal of Neuroscience
862      36:1096–1112.

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

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

867  Bair W, Zohary E, Newsome WT (2001) Correlated firing in macaque visual area
868      MT: time scales and relationship to behavior. Journal of Neuroscience
869      21:1676–1697.

870  Beck JM, Ma WJ, Pitkow X, Latham PE, Pouget A (2012) Perspective. Neuron
871      74:30–39.

872  Birrell JM, Brown VJ (2000) Medial frontal cortex mediates perceptual attentional
873      set shifting in the rat. Journal of Neuroscience 20:4320–4324.

874  Bondy AG, Haefner RM, Cumming BG (2018) Feedback determines the structure
875      of correlated variability in primary visual cortex. Nature Publishing Group:1–
876      15.

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

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

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

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

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

887  Collins AGE, Frank MJ (2013) Cognitive control over learning: creating,
888      clustering, and generalizing task-set structure. Psychological Review
889      120:190–229.

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

892   Doiron B, Litwin-Kumar A, Rosenbaum R, Ocker GK, Josić K (2016) The
893        mechanics of state-dependent neural correlations. Nature Publishing Group
894        19:383–393.

895   Downer JD, Niwa M, Sutter ML (2015) Task engagement selectively modulates
896        neural correlations in primary auditory cortex. Journal of Neuroscience
897        35:7565–7574.

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

900   Gu Y, Liu S, Fetsch CR, Yang Y, Fok S, Sunkara A, DeAngelis GC, Angelaki DE
901        (2011) Perceptual learning reduces interneuronal correlations in macaque
902        visual cortex. Neuron 71:750–761.

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

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

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

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

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

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

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

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

922   Joshi S, Li Y, Kalwani RM, Gold JI (2016) Relationships between Pupil Diameter
923        and Neuronal Activity in the Locus Coeruleus, Colliculi, and Cingulate Cortex.

924        Neuron 89:221–234.

925    Kanitscheider I, Coen-Cagli R, Pouget A (2015) Origin of information-limiting
926        noise correlations. Proceedings of the National Academy of Sciences
927        112:E6973–E6982.

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

930    Krotov D, Hopfield JJ (2019) Unsupervised learning by competing hidden units.
931        Proceedings of the National Academy of Sciences 116:7723–7731.

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

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

938    Law C-T, Gold JI (2009) Reinforcement learning can account for associative and
939        perceptual learning on a visual-decision task. Nature Neuroscience 12:655–
940        663.

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

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

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

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

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

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

958     533.

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

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

964     Oja E (1982) Simplified neuron model as a principal component analyzer. Journal
965         of Mathematical Biology:1–7.

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

968     Reimer J, McGinley MJ, Liu Y, Rodenkirch C, Wang Q, McCormick DA, Tolias AS
969         (2016) Pupil fluctuations track rapid changes in adrenergic and cholinergic
970         activity in cortex. Nature Communications 7:13289.

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

973     Shadlen MN, Newsome WT (1998) The variable discharge of cortical neurons:
974         implications for connectivity, computation, and information coding. J Neurosci
975         18:3870–3896.

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

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

980     Tsividis P, Pouncy T, Xu JL, Tenenbaum JB, Gershman SJ (2017) Human
981         Learning in Atari. 2017 AAAI Spring Symposium Series, Science of
982         Intelligence: Computational Principles of Natural and Artificial Intelligence:1–
983         4.

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

987     Wimmer RD, Schmitt LI, Davidson TJ, Nakajima M, Deisseroth K, Halassa MM
988         (2015) Thalamic control of sensory selection in divided attention. Nature
989         526:705–709.

990     Zanto TP, Gazzaley A (2009) Neural Suppression of Irrelevant Information
991         Underlies Optimal Working Memory Performance. Journal of Neuroscience

992      29:3059–3066.

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

995