

426 A Notation

427 For $N, K \geq 2$, let \mathbb{R}^N denote N -dimensional Euclidean space equipped with the usual Euclidean
 428 norm $\|\cdot\|$ and let \mathbb{R}_+^N denote the non-negative orthant. Let $\mathbf{0} = [0, \dots, 0]^\top$ and $\mathbf{1} = [1, \dots, 1]^\top$
 429 respectively denote the vectors of zeros and ones, whose dimensions should be clear from context.

430 Let $\mathbb{R}^{N \times K}$ denote the set of $N \times K$ real-valued matrices. Let $\|\cdot\|_F$ denote the Frobenius norm and
 431 $\|\cdot\|_{\text{op}}$ denote the operator norm. Let $O(N)$ denote the set of $N \times N$ orthogonal matrices. Let \mathbb{S}^N
 432 (resp. $\mathbb{S}_+^N, \mathbb{S}_{++}^N$) denote the set of $N \times N$ symmetric (resp. positive semidefinite, positive definite)
 433 matrices. Let \mathbf{I}_N denote the $N \times N$ identity matrix.

434 Given vectors $\mathbf{u}, \mathbf{v} \in \mathbb{R}^N$, let $\mathbf{u} \circ \mathbf{v} = [u_1 v_1, \dots, u_N v_N]^\top \in \mathbb{R}^N$ denote the Hadamard (elementwise)
 435 product of \mathbf{u} and \mathbf{v} . Let $\text{diag}(\mathbf{u})$ denote the $N \times N$ diagonal matrix whose $(i, i)^{\text{th}}$ entry is u_i , for
 436 $i = 1, \dots, N$. Given a matrix $\mathbf{M} \in \mathbb{R}^{N \times N}$, we also let $\text{diag}(\mathbf{M})$ denote the N -dimensional vector
 437 whose i^{th} entry is M_{ii} .

438 B Separation of timescales for the gain and synaptic weight updates

439 In this section, we consider an algorithm where we directly optimize the objective in equation 3. In
 440 particular, for each context c , we first optimize over the gains \mathbf{g} and then take a gradient descent step
 441 with respect to \mathbf{W} .

442 We first compute the gains using the formula for the optimal gains derived in [25, equation 18]:

$$\mathbf{g} = [(\mathbf{W}^\top \mathbf{W})^{\circ 2}]^\dagger \text{diag} \left(\mathbf{W}^\top \mathbf{C}_{ss}^{1/2}(c) \mathbf{W} - \mathbf{W}^\top \mathbf{W} \right).$$

443 We then update the synaptic weights by taking the following gradient descent step:

$$\Delta \mathbf{w}_i = \mathbb{E}_{\mathbf{s} \sim p(\mathbf{s}|c)} [\eta_w (\mathbf{r} \mathbf{r}^\top \mathbf{w}_i g_i - \mathbf{w}_i g_i)] = \eta_w (\mathbf{M}(\mathbf{W}, \mathbf{g})^{-1} \mathbf{C}_{ss}(c) \mathbf{M}(\mathbf{W}, \mathbf{g})^{-1} - \mathbf{I}_N) \mathbf{w}_i g_i,$$

444 where $\mathbf{M}(\mathbf{W}, \mathbf{g}) := \alpha \mathbf{I}_N + \mathbf{W} \text{diag}(\mathbf{g}) \mathbf{W}^\top$. Combining these updates yields Algorithm 2, which
 445 takes context-dependent covariance matrices $\mathbf{C}_{ss}(c)$ as its input.

Algorithm 2: Adaptive whitening via synaptic plasticity and gain modulation

- 1: **Input:** Covariance matrices $\mathbf{C}_{ss}(1), \mathbf{C}_{ss}(2), \dots$
 - 2: **Initialize:** $\mathbf{W} \in \mathbb{R}^{N \times K}; \eta_w > 0$
 - 3: **for** $c = 1, 2, \dots$ **do**
 - 4: $\mathbf{g} \leftarrow [(\mathbf{W}^\top \mathbf{W})^{\circ 2}]^\dagger \text{diag} \left(\mathbf{W}^\top \mathbf{C}_{ss}^{1/2}(c) \mathbf{W} - \mathbf{W}^\top \mathbf{W} \right)$
 - 5: $\mathbf{G} \leftarrow \text{diag}(\mathbf{g})$
 - 6: $\mathbf{W} \leftarrow \mathbf{W} + \eta_w \left((\mathbf{W} \mathbf{G} \mathbf{W}^\top)^{-1} \mathbf{C}_{ss}(c) (\mathbf{W} \mathbf{G} \mathbf{W}^\top)^{-1} \mathbf{W} \mathbf{G} - \mathbf{W} \mathbf{G} \right)$
 - 7: **end for**
-

446 C Adaptive whitening of natural images

447 In this section, we elaborate on the converged structure of \mathbf{W}_T using natural image patches. To better
 448 visualize the relationship between the learned columns of \mathbf{W} and sinusoidal basis functions (e.g.
 449 DCT), we focus on 1-dimensional image patches (rows of pixels). The results are similar with 2D
 450 image patches.

451 It is well known that eigenvectors of natural images are well-approximated by sinusoidal basis
 452 functions [e.g. the DCT; 33; 34]. Using the same images from the main text [32], we generated 56
 453 contexts by sampling 16×1 pixel patches from separate images, with 2E4 samples each. We train
 454 Algorithm 2 with $K = N = 16$, $\eta_w = 5\text{E-}2$, and random $\mathbf{W}_0 \in O(16)$ on a training set of \mathbf{X} of the
 455 images, presented uniformly at random $T = 1\text{E}5$ times. Fig C.1A,B shows that \mathbf{W}_T approximates
 456 the principal components of the aggregated context-dependent covariance, $\mathbb{E}_{c \sim p(c)}[\mathbf{C}_{ss}(c)]$, which
 457 are closely aligned with the DCT. To show that this structure is inherent in the spatial statistics of
 458 natural images, we generated control contexts, $\mathbf{C}_{ss}(c)$, by forming covariance matrices with matching

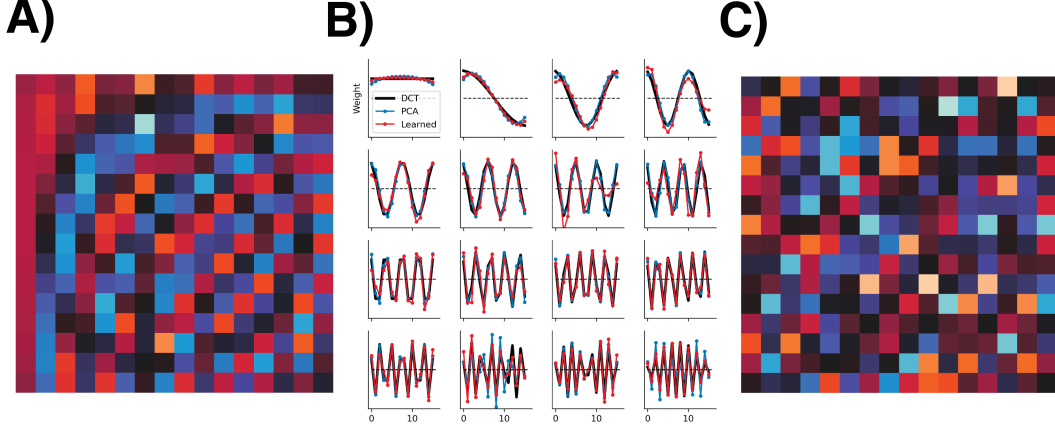


Figure C.1: Control experiment accompanying Sec. 5.2. **A)** \mathbf{W}_T learned from natural image patches. **B)** Basis vectors from **A** displayed as line plots, compared to the 1D DCT, and principal components of $\mathbb{E}_{c \sim p(c)}[C_{ss}(c)]$. **C)** Control condition. \mathbf{W}_T learned from spectrally-matched image patches with random eigenvectors.

eigenspectra, but each with *random* and distinct eigenvectors. This destroys the structure induced by natural image statistics. Consequently, the learned vectors in \mathbf{W}_T are no longer sinusoidal (Fig C.1C). As a result, whitening error with \mathbf{W}_T is much higher on the training set, with 0.3 ± 0.02 error (mean \pm standard error over 10 random initializations; Eq. 6) on natural image contexts and 2.7 ± 0.1 on the control contexts. While for the natural images, a basis approximating the DCT was sufficient to adaptively whiten all contexts in the ensemble, this is not the case for the generated control contexts.

Finally, we find that as K increases from $K = 1$ to $K = 16$, the basis vectors in \mathbf{W}_T progressively learn higher frequency components of the DCT (Fig. C.2). This is a sensible solution, due to the ℓ_2 reconstruction error of our objective, and the $1/f$ spectral content of natural image statistics. With more flexibility, as K increases past N (i.e. the overcomplete regime), the network continues to improve its whitening error (Fig. C.3A) by learning a basis, \mathbf{W}_T , that can account for within-context information that is insufficiently captured by the DCT (Fig. C.3B). Taken together, our model successfully learns a basis \mathbf{W}_T that exploits the spatial structure present in natural images.

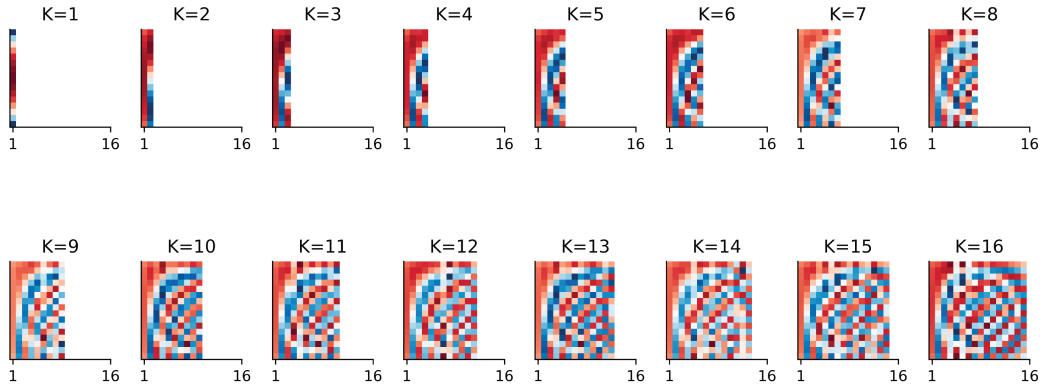


Figure C.2: As K increases, columns of \mathbf{W} progressively learn higher frequency components of the DCT.

472 D Modifications for increased biological realism

473 In this section, we modify Algorithm 1 to be more biologically realistic.

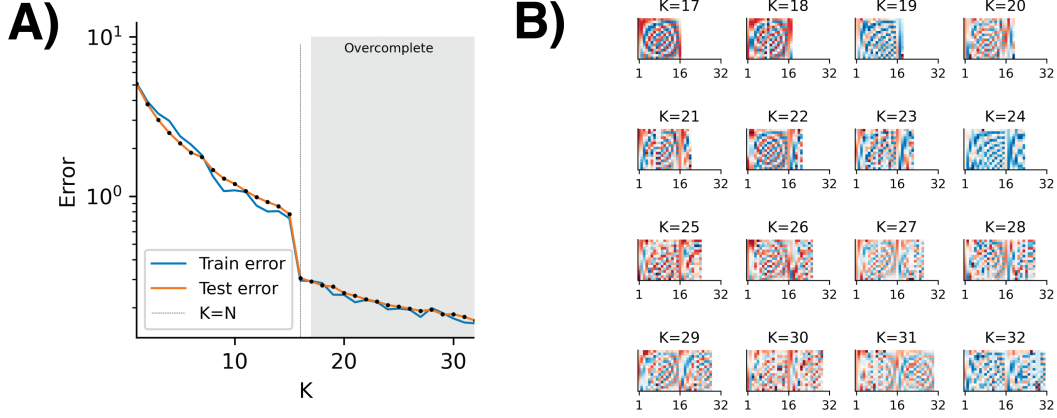


Figure C.3: **A)** Error on training and test set as a function of K . **B)** In the overcomplete regime, the network converges to a \mathbf{W}_T that helps to improve error compared to the $K \leq N$ regime.

474 D.1 Enforcing unit norm basis vectors

475 In our algorithm, there is no constraint on the magnitude of the column vectors of \mathbf{W} . We can enforce
 476 a unit norm (here measured using the Euclidean norm) constraint by adding Lagrange multipliers to
 477 the objective in equation 3:

$$\min_{\mathbf{W} \in \mathbb{R}^{N \times K}} \max_{\mathbf{m} \in \mathbb{R}^K} \mathbb{E}_{c \sim p(c)} \left[\min_{\mathbf{g} \in \mathbb{R}^K} \mathbb{E}_{\mathbf{s} \sim p(\mathbf{s}|c)} [g(\mathbf{W}, \mathbf{g}, \mathbf{r}, \mathbf{s})] \right], \quad (\text{D.1})$$

478 where

$$g(\mathbf{W}, \mathbf{g}, \mathbf{r}, \mathbf{s}) = \ell(\mathbf{W}, \mathbf{g}, \mathbf{r}, \mathbf{s}) + \sum_{i=1}^K m_i (\|\mathbf{w}_i\|^2 - 1).$$

479 Taking partial derivatives with respect to \mathbf{w}_i and \mathbf{m}_i results in the updates:

$$\begin{aligned} \Delta \mathbf{w}_i &= \eta_w (n_i \mathbf{r} - (g_i + m_i) \mathbf{w}_i) \\ \Delta m_i &= \|\mathbf{w}_i\|^2 - 1. \end{aligned}$$

480 Furthermore, since the weights are constrained to have unit norm, we can replace $\|\mathbf{w}_i\|^2$ with 1 in
 481 the gain update:

$$\Delta g_i = \eta_g (z_i^2 - 1).$$

482 D.2 Decoupling the feedforward and feedback weights

483 We replace the primary neuron-to-interneuron weight matrix \mathbf{W}^\top (resp. interneuron-to-primary
 484 neuron weight matrix $-\mathbf{W}$) with \mathbf{W}_{rn} (resp. $-\mathbf{W}_{nr}$). In this case, the update rules are

$$\begin{aligned} \mathbf{W}_{rn} &\leftarrow \mathbf{W}_{rn} + \eta_w (\mathbf{n}_t \mathbf{r}_t^\top - \text{diag}(\mathbf{g} + \mathbf{m}) \mathbf{W}_{rn}) \\ \mathbf{W}_{nr} &\leftarrow \mathbf{W}_{nr} + \eta_w (\mathbf{r}_t \mathbf{n}_t^\top - \mathbf{W}_{nr} \text{diag}(\mathbf{g} + \mathbf{m})). \end{aligned}$$

485 Let $\mathbf{W}_{rn,t}$ and $\mathbf{W}_{nr,t}$ denote the values of the weights \mathbf{W}_{rn} and \mathbf{W}_{nr} , respectively, after $t = 0, 1, \dots$
 486 iterates. Then for all $t = 0, 1, \dots$,

$$\mathbf{W}_{rn,t}^\top - \mathbf{W}_{nr,t} = (\mathbf{W}_{rn,0}^\top - \mathbf{W}_{nr,0}) (\mathbf{I}_N - \eta_w \text{diag}(\mathbf{g} + \mathbf{m}))^t.$$

487 Thus, if $g_i + m_i \in (0, 2\eta_w^{-1})$ for all i (e.g., by enforcing non-negative g_i, m_i and choosing $\eta_w > 0$
 488 sufficiently small), then the difference decays exponentially in t and the feedforward and feedback
 489 weights are asymptotically symmetric.

490 D.3 Sign-constraining the synaptic weights and gains

491 The synaptic weight matrix \mathbf{W} and gains vector \mathbf{g} are not sign-constrained in Algorithm 1, which is
 492 not consistent with biological evidence. We can modify the algorithm to enforce the sign constraints
 493 by rectifying the weights and gains at each step. Here $[\cdot]_+$ denote the elementwise rectification
 494 operation. This results in the updates

$$\begin{aligned}\mathbf{g} &\leftarrow [\mathbf{g} + \eta_g (\mathbf{z} \circ \mathbf{z} - \mathbf{1})]_+ \\ \mathbf{W}_{rn} &\leftarrow [\mathbf{W}_{rn} + \eta_w (\mathbf{n}_t \mathbf{r}_t^\top - \text{diag}(\mathbf{g} + \mathbf{m}) \mathbf{W}_{rn})]_+ \\ \mathbf{W}_{nr} &\leftarrow [\mathbf{W}_{nr} + \eta_w (\mathbf{r}_t \mathbf{n}_t^\top - \mathbf{W}_{nr} \text{diag}(\mathbf{g} + \mathbf{m}))]_+.\end{aligned}$$

495 D.4 Online algorithm with improved biological realism

496 Combining these modifications yields our more biologically realistic multi-timescale online algorithm,
 497 Algorithm 3.

Algorithm 3: Biologically realistic multi-timescale adaptive whitening

```

1: Input:  $\mathbf{s}_1, \mathbf{s}_2, \dots \in \mathbb{R}^N$ 
2: Initialize:  $\mathbf{W}_{nr} \in \mathbb{R}^{N \times K}$ ;  $\mathbf{W}_{rn} \in \mathbb{R}^{K \times N}$ ;  $\mathbf{m}, \mathbf{g} \in \mathbb{R}^K$ ;  $\eta_r, \eta_m > 0$ ;  $\eta_g \gg \eta_w > 0$ 
3: for  $t = 1, 2, \dots$  do
4:    $\mathbf{r}_t \leftarrow \mathbf{0}$ 
5:   while not converged do
6:      $\mathbf{z}_t \leftarrow \mathbf{W}_{rn} \mathbf{r}_t$ ; // interneuron inputs
7:      $\mathbf{n}_t \leftarrow \mathbf{g} \circ \mathbf{z}_t$ ; // gain-modulated interneuron outputs
8:      $\mathbf{r}_t \leftarrow \mathbf{r}_t + \eta_r (\mathbf{s}_t - \mathbf{W}_{nr} \mathbf{n}_t - \alpha \mathbf{r}_t)$ ; // recurrent neural dynamics
9:   end while
10:   $\mathbf{m} \leftarrow \mathbf{m} + \eta_m (\text{diag}(\mathbf{W}_{rn} \mathbf{W}_{nr}) - \mathbf{1})$ ; // weight normalization update
11:   $\mathbf{g} \leftarrow [\mathbf{g} + \eta_g (\mathbf{z}_t \circ \mathbf{z}_t - \mathbf{1})]_+$ ; // gains update
12:   $\mathbf{W}_{rn} \leftarrow [\mathbf{W}_{rn} + \eta_w (\mathbf{n}_t \mathbf{r}_t^\top - \text{diag}(\mathbf{g} + \mathbf{m}) \mathbf{W}_{rn})]_+$ ; // synaptic weights update
13:   $\mathbf{W}_{nr} \leftarrow [\mathbf{W}_{nr} + \eta_w (\mathbf{r}_t \mathbf{n}_t^\top - \mathbf{W}_{nr} \text{diag}(\mathbf{g} + \mathbf{m}))]_+$ 
14: end for
```
