

A Single Model Explains both Visual and Auditory Precortical Coding

Honghao Shan,¹ Matthew H. Tong,¹ Garrison W. Cottrell^{1*}

¹Department of Computer Science and Engineering, University of California, San Diego
La Jolla, CA 92093-0404

*To whom correspondence should be addressed; E-mail: gary@ucsd.edu.

¹HS and GWC designed the experiments, HS implemented the models, HS, MHT, and GWC wrote the manuscript.

March 7, 2024

1 Precortical neural systems encode information collected by the senses,
2 but the driving principles of the encoding used have remained a subject of
3 debate. We present a model of retinal coding that is based on three con-
4 straints: information preservation, minimization of the neural wiring, and
5 response equalization. The resulting novel version of sparse principal com-
6 ponents analysis successfully captures a number of known characteristics of
7 the retinal coding system, such as center-surround receptive fields, color op-
8 ponency channels, and spatiotemporal responses that correspond to magno-
9 cellular and parvocellular pathways. Furthermore, when trained on auditory
10 data, the same model learns receptive fields well fit by gammatone filters,
11 commonly used to model precortical auditory coding. This suggests that
12 efficient coding may be a unifying principle of precortical encoding across
13 modalities.

14 Introduction

15 Sensory information goes through various forms of processing before it reaches
16 the cerebral cortex. Visual information is transformed into neural signals at
17 the retina, where it passes through retinal ganglion cells that are charac-
18 terized by their center-surround shaped receptive fields (Enroth-Cugell &
19 Robson, 1966); auditory information, on the other hand, is passed to the
20 brain through the auditory nerve fibers whose filtering properties can be well
21 described by gammatone filters (Kiang, Watanabe, Thomas, & Clark, 1965).
22 Since such peripheral processing prepares the data that the subsequent cor-
23 tical processing relies on, its functional role has attracted a great deal of
24 attention in the past several decades (Srinivasan, Laughlin, & Dubs, 1982;
25 Field, 1987; Atick & Redlich, 1992; Lewicki, 2002; Vincent, Baddeley, Tros-
26 cianko, & Gilchrist, 2005; Graham, Chandler, & Field, 2006; Doi & Lewicki,
27 2007).

28 Despite intensive research, there are still mysteries concerning the func-
29 tional role of pre-cortical processing. For example, do different sensory
30 modalities (visual, auditory, somatosensory, etc.) adopt the same compu-
31 tational principles in their pre-cortical stages? Although it is tantalizing to
32 assume so, recent studies suggest otherwise. For example, (Lewicki, 2002)
33 learned gammatone filters from natural sound using independent compo-
34 nent analysis (ICA), which was previously applied to natural image patches
35 to learn edge/bar shaped filters resembling the V1 simple cells' receptive
36 fields (Olshausen & Field, 1996; Lewicki & Sejnowski, 2000). Since gamma-
37 tones model pre-cortical auditory nerve fibers while V1 is a region of cortex,
38 this gives rise to a puzzle: Why would the brain use the same strategy for pre-
39 processing at a *pre-cortical* stage in the auditory pathway and early *cortical*
40 processing in the visual pathway (Olshausen & O'Connor, 2002)?

41 Questions remain even for peripheral processing in a single modality. Re-
42 cently, Graham et al. proposed that decorrelation, response equalization,
43 and sparseness form the minimum constraints that must be considered to ac-
44 count for the known linear properties of retinal coding (Graham et al., 2006).
45 This hypothesis is the combination of several previous theories. The response
46 equalization theory hypothesizes that retinal coding seeks a representation
47 that ensures that each neuron has approximately the same average activity
48 level when the animal is presented with natural scenes (Field, 1987; Brady &
49 Field, 1995; Field & Brady, 1997; Brady & Field, 2000). The output decor-
50 relation theory follows the efficient coding principle (Attneave, 1954; Barlow,

1961; Atick & Redlich, 1992). It hypothesizes that retinal coding represents the most efficient coding of the information in the visual domain by capturing the second-order statistical structure of the visual inputs and making the signals from these neurons less correlated. Both of these theories are derivatives of whitening theory, which hypothesizes that retina coding produces a flattened response spectrum for natural visual inputs from a specific range of spatial frequencies (Srinivasan et al., 1982). This whitening theory links the properties of retinal coding with the statistics of natural scenes and is now part of the prevailing view of retinal processing. A third theory suggests that the system is trying to minimize energy usage or wiring cost (Vincent & Baddeley, 2003). Vincent et al. argued that systems that try to minimize energy usage by minimizing wiring give center surround receptive fields. Given these various objectives, it is still not clear what constraints are actually operating in the specification of the retinal coding system. It would be desirable to build a retinal coding model that integrates the different ideas behind these theories and explains the origins of the observed center-surround receptive fields. Ideally, this model should also be able to explain the pre-cortical processing of other modalities.

Our model takes into account the following considerations. The retina compresses the approximately 100 million photoreceptor responses into a million ganglion cell responses. Hence the first consideration is that we would like the ganglion cells to retain the maximum amount of information about the photoreceptor responses. If we make the simplifying assumption that ganglion cells respond linearly, then the optimal linear compression technique in terms of reconstruction error is principal components analysis (PCA). One can map PCA into a neural network as in Figure 1(a) (Cottrell, Munro, & Zipser, 1989; Baldi & Hornik, 1989). The weight vectors of each hidden unit in this network each correspond to one eigenvector of the covariance matrix of the data. In standard PCA, there is an ordering to the hidden units, such that the first hidden unit has very high response variance and the last hidden unit has practically no variance, which means the first hidden unit is doing orders of magnitude more work than the last one. The second consideration, then, is that we would like to spread the work evenly among the hidden units. Hence we impose a threshold on the average squared output of the hidden units. As we will see from the simulations, in order to preserve the maximum information, the units all hit this threshold, which equalizes the work. The third consideration is that PCA is profligate with connections - every ganglion cell would have non-zero connections to every

89 photoreceptor. Hence we also impose a constraint on the connectivity in the
 90 network. In this latter constraint we were inspired by the earlier work of
 91 (Vincent et al., 2005). They proposed a model of retinal and early cortical
 92 processing based on energy minimization and showed that it could create
 93 center-surround shaped receptive fields for grayscale images. However, their
 94 system sometimes led to cells with two center-surround fields, and the opti-
 95 mization itself was unstable.

96 These considerations lead to our objective function. Images can be rep-
 97 resented as high-dimensional real-valued data; if L photoreceptors are repre-
 98 senting the input image, the observed image can be represented as $x \in R^L$.
 99 Given input vectors $\mathbf{x} \in R^L$, we seek to find the output responses $\mathbf{s} \in R^M$
 100 (the signal from the retinal ganglion cells) and basis functions $A \in R^{L \times M}$
 101 (the connections from the photoreceptors to the ganglion cells) such that the
 102 following objective function is minimized:

$$E = \left\langle \frac{\|\mathbf{x} - \mathbf{A}\mathbf{s}\|_2^2}{2} \right\rangle + \lambda \|\mathbf{A}\|_1 \quad (1)$$

103 subject to:

$$\langle s_i^2 \rangle \leq 1 \quad \forall i \quad (2)$$

104 where $\langle \cdot \rangle$ denotes taking average over all the input samples. The first term in
 105 Equation 1 minimizes the reconstruction error and maximizes the informa-
 106 tion maintained by the encoding. When the sparsity weight λ is small and
 107 $L > M$ (i.e., the encoding compresses the information), the reconstruction
 108 error reduces to a term that only involves the correlation matrix: $\mathbf{C} = \langle \mathbf{x}\mathbf{x}^t \rangle$
 109 (see supplementary materials). This concurs with the idea that the system
 110 is only sensitive to second order statistics. The second constraint minimizes
 111 the connections from the photoreceptors to the ganglion cells, incorporating
 112 sparsity and an economy of elementary features. The constraint on the aver-
 113 age energy of the ganglion cells equalizes the work across the ganglion cells.
 114 The system will in fact push this term to the threshold of 1 in order to main-
 115 tain the maximum information. Thus this objective function integrates three
 116 major theories of retinal coding: efficient coding (Attneave, 1954), response
 117 equalization (Field, 1987), and the economy of elementary features (Vincent
 118 et al., 2005). While this does not directly embody the decorrelation the-
 119 ory (Atick & Redlich, 1992), we have not found that assumption necessary
 120 to obtain our results.

Put another way, the three terms in the objective function determine different aspects of the basis functions (i.e., the columns of \mathbf{A}): the reconstruction error determines the subspace that the basis functions span; the output constraint specifies the lengths of the basis functions; and the sparsity penalty rotates the basis functions within the subspace determined by the reconstruction error. In this sense, all three terms, and hence all three theories that they embody, are necessary to fully characterize the retinal coding model. This observation partially agrees with the prediction by (Graham et al., 2006): “we conclude that a minimum of three constraints must be considered to account for the known linear properties – decorrelation, response equalization, and size/sparseness.” Our three constraints are efficient coding of information, response equalization, and sparseness of connections.

An important feature of this objective function is that it allows us to derive an efficient algorithm to learn the model parameters, because the revised model turns out to be a particular variation of Sparse PCA (Zou, Hastie, & Tibshirani, 2006) that is reducible to sparse coding (Olshausen & Field, 1996). We can therefore efficiently estimate the parameters of the model and apply it to a larger range of data than has typically been used in the past. A critical insight provided by the mapping to sparse coding is that this model *is* exactly sparse coding, applied to the *transpose* of the data matrix rather than the data matrix itself. This means we can also use our model for dimensionality *expansion* (overcomplete representations) as well as dimensionality reduction, although when doing expansion, we can no longer use the efficient approximation derived in the Supplementary Materials.

In what follows, we show that this simple objective function is able to account for both retinal ganglion cell receptive fields and gammatone filters that have been used to characterize the signals in the auditory nerve.

Results

Going forward, it is important to understand the distinction between *features* and *filters*. The features are the rows of \mathbf{A}^t , and represent the connections between the photoreceptors and the hidden units in Figure 1(b). The filters, on the other hand, correspond to the rows of $\mathbf{W} = (\mathbf{A}^t \mathbf{A})^{-1} \mathbf{A}^t$, the pseudoinverse of \mathbf{A} , which correspond to the receptive fields of the ganglion cells that would result from reverse correlation. We visualize the computation as in Figure 1(b): the network receives input from the photoreceptors, and then

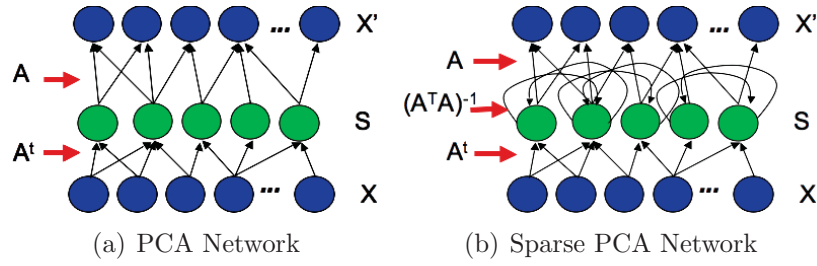


Figure 1: Two neural networks that can implement PCA and Sparse PCA. The left hand panel represents a network that performs PCA. The weights (rows of \mathbf{A}^t) to each hidden unit (ganglion cell) from the pixels (photoreceptors) represent the coordinates of a unit-length eigenvector of the covariance matrix of the data, so that the activations are the projections onto the eigenvectors. The same weights (transposed) can be used to reconstruct the data (these are not part of the model). The right panel represents Sparse PCA. It is the connections from the pixels that are sparse (\mathbf{A}^t). This is followed by recurrent connections that give the center-surround shape. This would be reflected in any recordings of the hidden units, hence the receptive field of a hidden unit is represented by a row of $\mathbf{W} = (\mathbf{A}^t \mathbf{A})^{-1} \mathbf{A}^t$.

156 there is inhibition of hidden units with similar receptive fields (represented
 157 by the recurrent connections $(\mathbf{A}^t \mathbf{A})^{-1}$).

158 **Grayscale Images** We applied the model to four grayscale image datasets.
 159 Results were qualitatively similar across the sets; the results we describe
 160 here are from the subset of the Van Hateren natural image set described
 161 in (Karklin & Lewicki, 2009). For this simulation, we used 20×20 patches
 162 of pixels and reduced them to 100 dimensions. Our SPCA model captures
 163 99.23% of the variance that is captured by standard PCA with 100 eigen-
 164 vectors retained, while 96.31% of the connection weights (the rows of \mathbf{A}^t)
 165 are absolute zero; in contrast, none of the connection weights in standard
 166 PCA are zero. Figure 8(a) and 8(b) plots the distribution of the connection
 167 weights in \mathbf{A} learned by our model versus standard PCA. Each ganglion cell
 168 is directly connected to only 3.69% of the input neurons in the 20×20 patch
 169 on average. Clearly, this sparsity would be advantageous for a biological
 170 system.

171 The learned elementary features (i.e., the columns of \mathbf{A}) are blobs of sim-
 172 ilar size that tile the 20×20 image patch. The top panel in Figure 2(a)
 173 displays 10 features randomly selected from all 100 features. We fit all the
 174 features with 2D Gaussians, and plot them as circles in Figure 2(b). The cen-
 175 ter and the radius of each circle represent the center and twice the standard
 176 deviation of the fitted Gaussian. To visualize how well the Gaussians fit the
 177 features, we display the first feature in Figure 2(a) and highlight its fitted
 178 Gaussian. As shown in the figure, the Gaussians provide a mosaic coverage
 179 of the image patch. If we reduce the number of hidden units to 32, the blobs
 180 enlarge to cover the image, as shown in Figure 2(d).

181 The optimal filters (i.e., the rows of \mathbf{W}) are center-surround shaped, as
 182 shown in the bottom panel of Figure 2(a). The first filter, for example, recov-
 183 ers the weight assigned to the first feature in the top panel. It is tantalizing to
 184 think that some of the filters are ON-centered while others are OFF-centered.
 185 However, we can switch the signs of the features (and hence the signs of the
 186 optimal filters) without changing the model’s objective function. Hence our
 187 model does not provide insight into the difference between the ON-centered
 188 and the OFF-centered cells (Chichilnisky & Kalmar, 2002) beyond the usual
 189 explanation that neurons cannot fire both positively and negatively.

190 It is interesting to see why a population of Gaussian blob shaped features
 191 should give rise to center-surround shaped filters. As shown in Figure 2(c),

each filter is a weighted sum of all the elementary features and can be viewed as the result of a sequence of efforts to recover the contribution of its corresponding feature. Each feature is first applied as a template filter on the image patch to estimate its contribution. However, this estimation is inaccurate because this feature overlaps with its neighboring features. To get a more accurate estimation the contribution from neighboring features must be subtracted. This potentially overcompensates, so get an even more accurate estimation, one must add back the contribution from the features neighboring the features that surround the first feature. This process repeats, moving ever outward. However, the weight reduces quickly for features removed from the first feature, which makes the resulting filter effectively localized and keeps the filter center-surround shaped (for low lambda, some additional ripples can be observed - however, these additional ripples have also been observed in ganglion receptive fields (Dearworth Jr. & Granda, 2002)).

Which aspects of natural scene images give rise to the learned features we observe? To answer this question, we apply our algorithm to white noise images, which contain no statistical structure, and pink noise images, which follow the same $1/f$ power law as natural scene images (Field, 1987) but otherwise contain no structure. Figure 3 displays two example images and some of the learned filters. On white noise images, the learned features are one-pixel image templates; the corresponding filters also only contain one non-zero pixel. That is, the model simply keeps 64 pixels and ignores the other pixels. That's the best it can do with 64 features, because white noise images contain no structure. On pink noise images, we learn essentially the same elementary features (and hence the same filters) as those learned from natural scene images. This result supports the hypothesis that the center-surround shaped filters come from the $1/f$ power spectrum of natural scene images, which agrees with the classic whitening theory (Srinivasan et al., 1982; Field, 1987; Atick & Redlich, 1992).

The sparseness level λ plays an important role in shaping the learned features and filters. As λ increases, the model puts more emphasis on sparse connections at the cost of keeping less information about the inputs. In a biological system, this may occur when the system is on a strict energy budget. Here we check how the learned filters change with larger λ values.

By analyzing the filters in Fourier space, we can plot amplitude at various frequencies, giving a contrast sensitivity function. As shown in Figure 4, with larger λ value, the model becomes less sensitive to low frequency information, but more sensitive to high frequency information. This change matches with

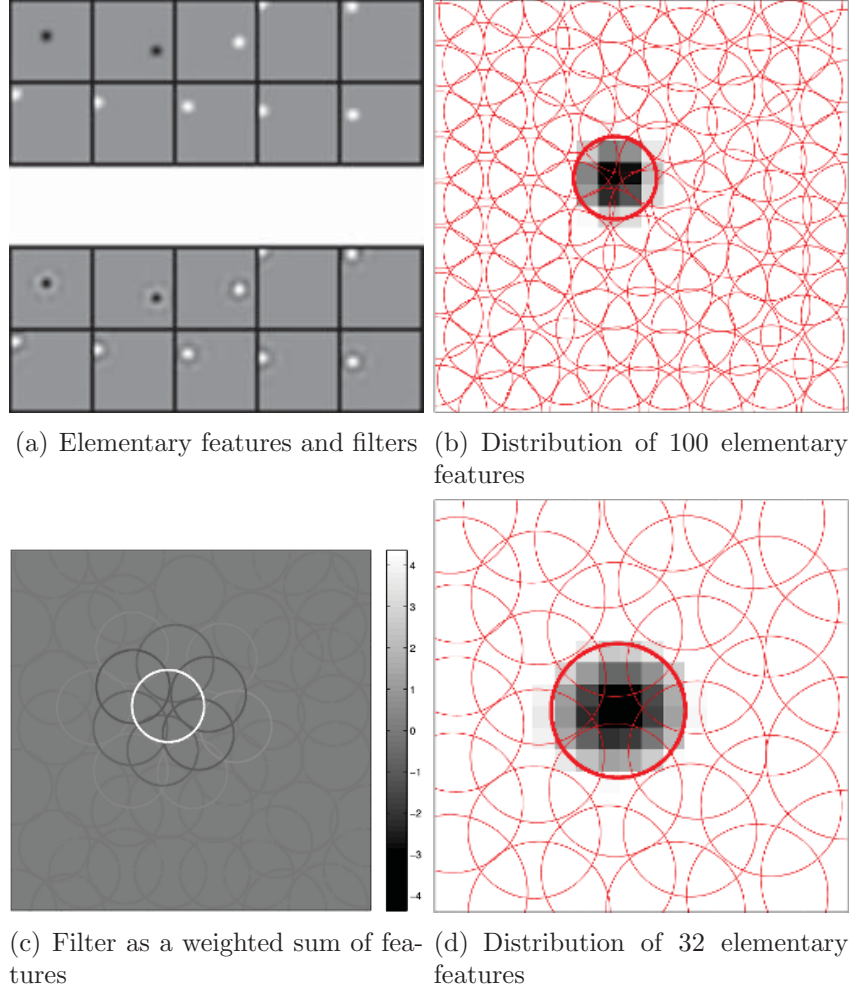


Figure 2: Elementary features and their corresponding optimal filters, learned from 20×20 grayscale image patches. In Figure 2(a), the top panel displays 10 features randomly selected from all the 100 features; the bottom panel displays their optimal filters. We fit the features with Gaussian blobs and plot them as circles in Figure 2(b). The radius of each circle represents twice the standard deviation of the fitted Gaussian blob. To help visualize how well the Gaussian blobs fit with the features, we display the first feature in Figure 2(a) and highlight its fitted Gaussian. These Gaussian will become bigger if we use a smaller number of features to “construct” the image patches, as shown in Figure 2(d). Figure 2(c) displays the first filter in Figure 2(a) as a weighted sum of all the 100 features. Each feature is plotted as a circle, as in Figure 2(b) and 2(d). The color of each circle represents the weight assigned to this feature (read the main text for details).

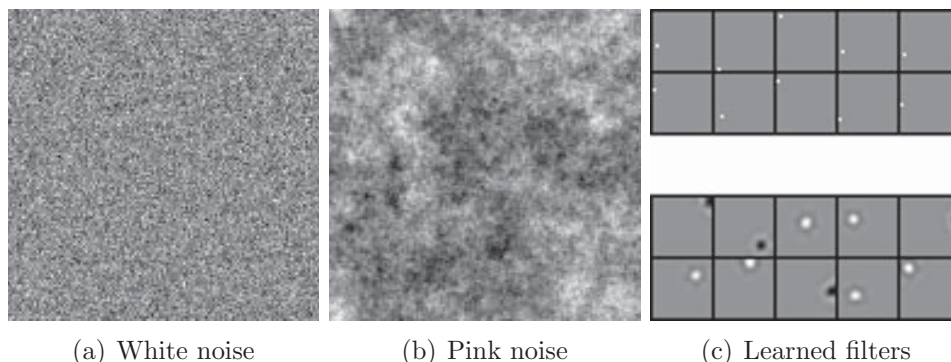


Figure 3: Experiments on white noise and pink noise images. Figure 3(a) and 3(b) display images containing white noise and pink noise. The top panel in Figure 3(c) displays 10 filters learned from white noise images; the bottom panel plays the filters learned from pink noise images.

230 psychophysical studies of contrast sensitivity in children with chronic malnu-
 231 trition. Compared with normal children, malnourished children are reported
 232 to be less sensitive to low spatial frequencies, but slightly more sensitive to
 233 high spatial frequencies (dos Santos & Alencar, 2010). This shift of acuity
 234 towards high frequencies, as suggested by our result, might due to the effort
 235 of the neural system to capture more visual information with a limited neural
 236 wiring budget.

237 **Chromatic Images** We applied our algorithm to four chromatic image
 238 datasets and again found that we learn qualitatively similar features with
 239 each. Here, we report the features learned from Kyoto image dataset. Reti-
 240 nal L, M, and S cones are estimated and given as input to the model. The re-
 241 sulting model captures 99.75% of the variance that is captured by an optimal
 242 linear model (PCA) with 256 output neurons, with 96.11% of its connections
 243 being absolute zero.

244 Figure 5(a) displays 6 representative features as well as their correspond-
 245 ing filters, learned from chromatic image patches. We visualize the con-
 246 nection strength from three of the filters to the L/M/S channels in Fig-
 247 ure 5(b). Among all the 256 learned features, 193 are black/white blobs, 48
 248 are blue/yellow blobs, 15 are red/green blobs. Figure 5(c), 5(d), and 5(e)
 249 plot the spatial layout of learned features.

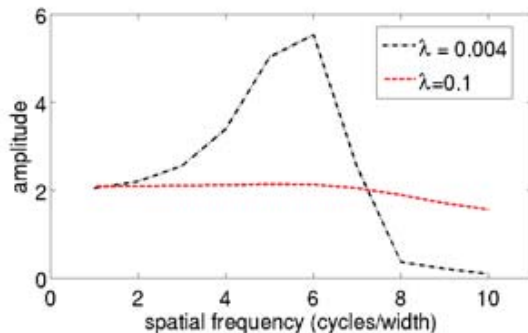


Figure 4: Experiments with increased sparseness level λ . We plot the amplitudes of different frequency component of the filters learned with $\lambda = 0.004$ and $\lambda = 0.1$. As shown in the figure, with an increased λ value, the filter becomes less sensitive to low frequencies, but more sensitive to high frequencies.

250 The above result replicates the segregation of the spatial channel and the
 251 color channel at the retina stage (Calkins & Sterling, 1999). This segregation
 252 was explored in previous research that applied information-theoretic meth-
 253 ods to natural color spectra/images, such as PCA (Buchsbaum & Gottschalk,
 254 1983; Derrico & Buchsbaum, 1991) and ICA (Tailor, Finkel, & Buchsbaum,
 255 2000; Wachtler, Lee, & Sejnowski, 2001; Lee, Wachtler, & Sejnowski, 2002;
 256 Doi, Inui, Lee, Wachtler, & Sejnowski, 2003; Caywood, Willmore, & Tolhurst,
 257 2004). One common observation in these studies is that the learned visual fea-
 258 tures (eigenvectors or independent components) segregate into black/white,
 259 blue/yellow, and red/green opponent structures, either shaped as Fourier ba-
 260 sis functions or Gabor kernel functions (Wachtler et al., 2001). These results
 261 are similar to what is obtained with ZCA (Zero-component analysis) (Doi et
 262 al., 2003), although ZCA has small connections to all of the inputs, since it
 263 does not inherently try to minimize connections.

264 **Grayscale Videos** To explore the spatio-temporal structure of natural
 265 videos, we collected a video dataset of 27 clips from nature documentaries.

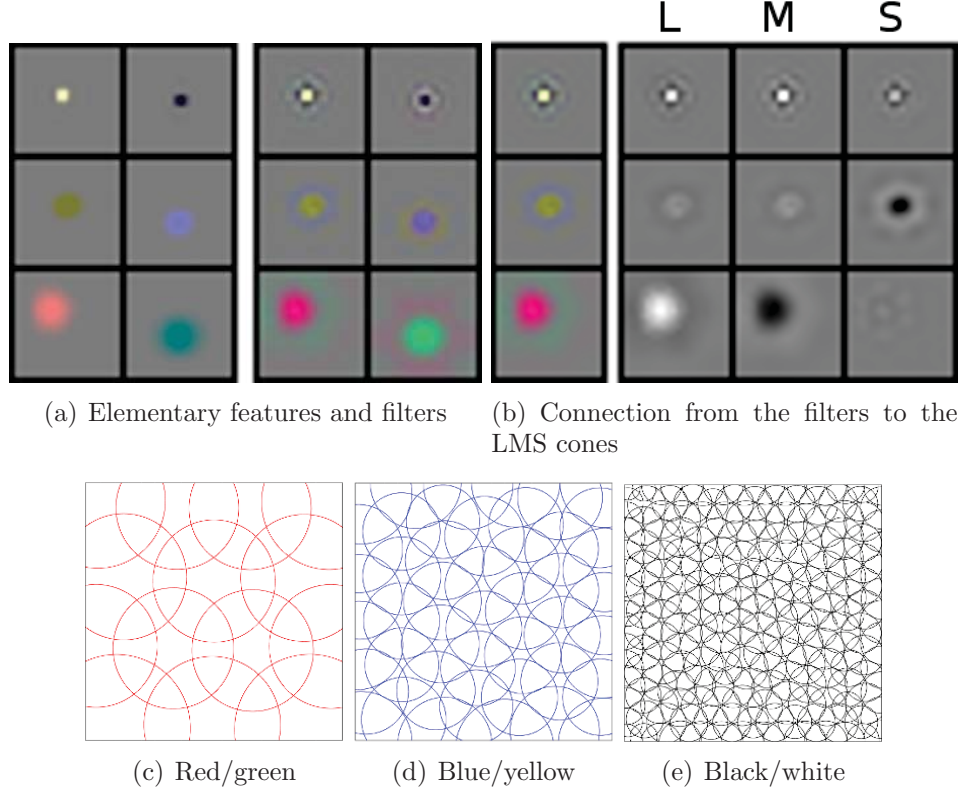


Figure 5: Elementary features and the corresponding filters learned from 20×20 chromatic image patches with $\lambda = 0.001$. In Figure 5(a), the left panel displays 6 representative features; the right panel displays their corresponding filters. The features belong to three categories: black/white blobs, blue/yellow blobs, and red/green blobs. The corresponding optimal filters are center-surround shaped, with black/white, blue/yellow, or red/green antagonism. Figure 5(b) plots the connection strength from the filters to the L, M, S cones. Figure 5(c), 5(d), and 5(e) plot the spatial layout of learned features, as we did in Figure 2(b).

266 Just as a two dimensional image patch can be flattened into a vector of
267 input responses, a three dimensional spatiotemporal patch of video can also
268 be tranformed into a vector. These vectors can then be given to the model
269 as input.

270 The learned features are black/white blobs whose contrast changes over
271 time. As shown in Figure 6(a), the features can be well fitted to spatio-
272 temporal Gaussians. The corresponding filters are spatially center-surround
273 shaped. Their temporal profile seems to provide an “edge” detector along
274 the temporal axis, which is similar to the temporal profile describing retinal
275 ganglion cells (Pillow et al., 2008). To see the animation file of the learned
276 filters, see http://cseweb.ucsd.edu/~gary/video_W.gif.

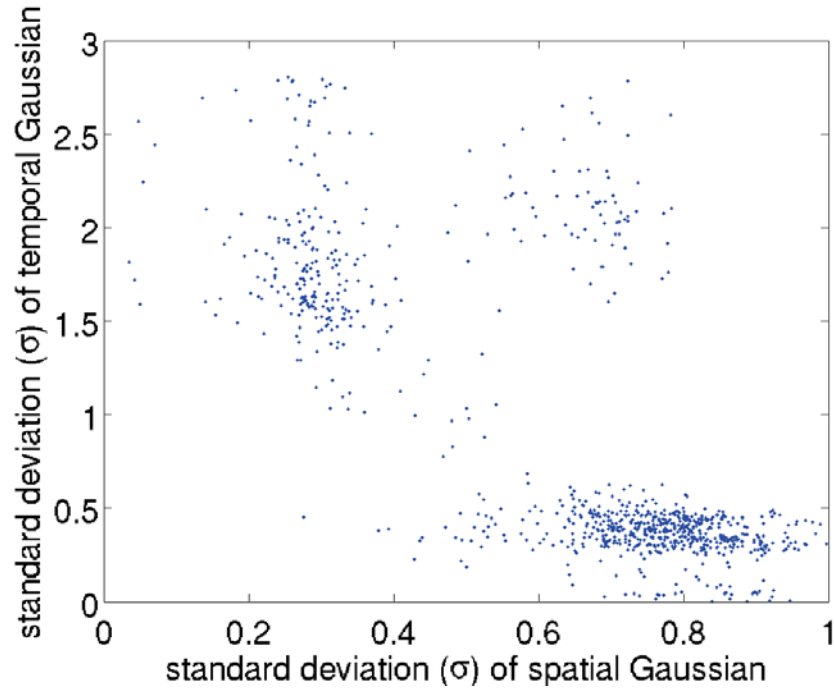
277 Another interesting observation is that most of the features segregate into
278 two groups: those with low-spatial and high-temporal frequencies (centered
279 around (0.3,1.75) in the figure), and those with high-spatial and low-temporal
280 frequencies (centered around (.75,0.4) in the figure), as plotted in Figure 6(b).
281 This suggests that the division of ganglion cells into the magno-pathway and
282 the parvo-pathway represents an efficient encoding of the visual environment.
283 This segregation appears to reveal statistical properties of natural videos,
284 instead of coming from our specific algorithm. In fact, we found that this
285 segregation exists even for features learned with standard PCA. As in PCA
286 of static images, however, the features are not biologically plausible.

287 **Sound** We applied our algorithm to three sound datasets and get quali-
288 tatively similar results on these three datasets. As shown in Figure 7, the
289 learned filters can be well fitted to gammatone filters. Gammatone filters
290 resemble the filtering properties of auditory nerve fibers estimated using the
291 reverse correlation technique from animal such as cats (Carney, 1990) and
292 chinchillas (Recio-Spinoso, Temchin, van Dijk, Fan, & Ruggero, 2005).

293 Note that this is (to the best of our knowledge) the first time a non-ICA
294 algorithm has learned gammatone-like filters from sound. Also, since we used
295 the same algorithm for both visual and auditory modalities, this provides an
296 answer to the question posed by Olshausen and O’Connor (2002): “Perhaps
297 an even deeper question is why ICA accounts for neural response properties
298 at the very earliest stage of analysis in the auditory system, whereas in the
299 visual system ICA accounts for the response properties of cortical neurons,
300 which are many synapses removed from photoreceptors.” Our model suggests
301 that it is not necessary to use ICA to obtain gammatone filters from sound;

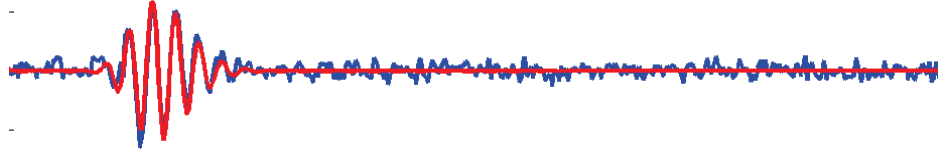


(a) two video filters learned by Sparse PCA (time from left to right)

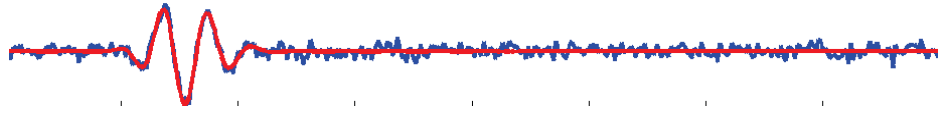


(b) distribution of video features learned by Sparse PCA

Figure 6: Video features segregate into two groups: small, persistent features, and large, brief features.



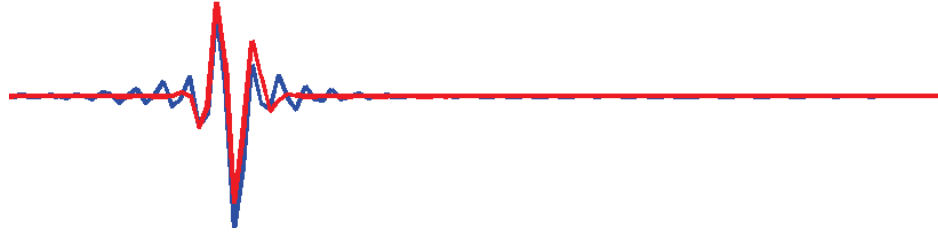
(a) revcor filter from cat



(b) revcor filter from cat



(c) filter learned from the TIMIT dataset



(d) filter learned from the Pittsburgh dataset

Figure 7: Figure 7(a) and 7(b) plot the revcor filters estimated from cat's auditory nerve fibers using the linear reverse correlation technique, as well as the fitted gammatone filter. Figure 7(c) and 7(d) plot filters learned from the TIMIT speech dataset and the Pittsburgh environmental sound dataset respectively. The blue line plots the estimated filter; the red line plots the fitted gammatone filter.

rather, Sparse PCA can account for the receptive fields of neurons at the very earliest stages of analysis in both auditory and visual modalities.

Discussion

We have suggested three principles that can be used to explain precortical encoding: information preservation, minimization of the neural wiring, and response equalization. Each of these principles can be independently justified via evolutionary and energy minimization arguments. Clearly, an organism should try to extract as much relevant information as possible from its environment. As organisms evolve to survive in more enriched environments, which information is relevant becomes more difficult to encode in the genome. As Barlow and Attneave have argued, redundancy reduction for efficient coding is a reasonable response to environmental complexity. Minimizing energy usage suggests constructing the minimal architecture possible, via minimizing wiring, which is a win in terms of both development and daily energy budgets. Finally, equalizing the work (response normalization) results in no single component being crucial to the organism.

The model is closely related to previous theories, but differs in crucial respects. We were inspired by Vincent et al.’s model (Vincent et al., 2005), which also attempted to derive center surround receptive fields by information preservation and minimizing wiring. By including the response normalization constraint, we were able to obtain a more stable algorithm with none of the occasional double receptive fields generated by their model. Unlike many previous models, explicitly decorrelating the outputs of our model was not necessary in order to obtain our results. Hence the model integrates two of the three components suggested as being necessary for any retinal coding model by Graham, Chandler, and Field (Graham et al., 2006) - response equalization and sparseness, but replaces decorrelation with minimization of the neural wiring. The Sparse PCA model has an interesting link to Olshausen’s sparse coding model; both models try to minimize the reconstruction error with response equalization, but our model imposes sparseness on the dictionary, while the sparse coding model imposes sparseness on the output. Finally, the model is closely related to that described in (Doi & Lewicki, 2007), where they proposed a model in which the retinal output is a linear transform of the input $\mathbf{s} = \mathbf{W}\mathbf{x}$, ignoring optical blur. The objective of their model is to minimize the difference between the input \mathbf{x} and its

337 reconstruction $\mathbf{A}\mathbf{W}\mathbf{x}$. Hence their model can be seen as the approximation
 338 we take when λ is small and $L > M$ (described in the Supplementary mate-
 339 rials). They also include two extra terms during learning, to regularize the
 340 average amplitude of the outputs, and to impose sparseness on \mathbf{W} (instead
 341 of \mathbf{A}). Our version, with sparseness on \mathbf{A} , leads to a simple interpretation
 342 of locally-connected ganglion cells with an inhibitory surround. Also, the
 343 convexity properties of our model lead to excellent convergence properties.

344 We derived an efficient algorithm to learn the model parameters by trans-
 345 forming it into a sparse coding problem. Our approximate algorithm uses
 346 only the covariance matrix \mathbf{C} , and runs orders of magnitude faster than the
 347 exact algorithm, while obtaining results that are less than one percent differ-
 348 ent in the objective function and learned basis functions. Our approximation
 349 works well under the assumption that λ is small and $L > M$. While fast
 350 approximation is not a necessary part of a successful model, the speed of
 351 computation allows it to be used on more rich data, such as video.

352 We applied our algorithm to grayscale images, color images, grayscale
 353 videos, human speech, and environmental sound, and learned visual and au-
 354 ditory filters that resemble the filtering properties of retinal ganglion cells
 355 and auditory nerve fibers. Some of the learned filters are novel. For exam-
 356 ple, it learns the magno and parvo segregation pathways; and it learns the
 357 gammatone filters from natural sound.

358 Finally, as noted above, our model suggests an answer the question of
 359 why ICA gives features corresponding to cortical receptive fields in vision,
 360 but seems necessary to obtain gammatone-like filters for sound, which is a
 361 pre-cortical level of processing. Our suggestion is that ICA is not necessary
 362 to obtain gammatone filters; rather, a PCA algorithm with sparsity and
 363 response equalization constraints can result in gammatone filters for sound,
 364 while also producing receptive fields similar to peripheral neurons in vision.

365 Materials and Methods

366 **Model** As previously discussed, our model seeks to find the output vectors
 367 $\mathbf{s} \in R^M$ (the signal from the retinal ganglion cells) and basis functions $\mathbf{A} \in$
 368 $R^{L \times M}$ such that the following objective function is minimized:

$$E = \left\langle \frac{\|\mathbf{x} - \mathbf{A}\mathbf{s}\|_2^2}{2} \right\rangle + \lambda \|\mathbf{A}\|_1 \quad (3)$$

369 (where $\langle \cdot \rangle$ denotes taking average over all the input samples) subject to the
 370 constraint that the average output of each cell $\langle s_i^2 \rangle \leq 1$. Upon convergence,
 371 the model will satisfy $\langle s_i^2 \rangle = 1$ as otherwise the objective function could
 372 easily be further reduced.

373 Vincent et al. interpret \mathbf{A} as the synaptic strength between the input
 374 and output neurons; and hence they interpret the sparsity penalty on \mathbf{A} as
 375 the desire to minimize the neural wiring cost. From a generative point of
 376 view, the columns of \mathbf{A} are the elementary features that the model uses to
 377 “construct” the observed inputs. Thus the overall objective function can be
 378 understood as capturing most information of the inputs using an economic
 379 dictionary of elementary features. However, when \mathbf{A} is fixed and full rank,
 380 the optimal output is in fact given by $\mathbf{s}^* = \mathbf{W}\mathbf{x}$, where \mathbf{W} is the pseudoin-
 381 verse of \mathbf{A} (if we ignore our constraint on the average activation of a cell
 382 - an assumption we experimentally verified as yielding a suitable approxi-
 383 mate solution). If we apply the linear reverse correlation technique to the
 384 model neurons (Chichilnisky, 2001), which recovers the filters transforming
 385 the inputs to the outputs, we will get the rows of \mathbf{W} as the model neurons’
 386 receptive fields.

387 The columns of \mathbf{A} and the rows of \mathbf{W} describe different aspects of the
 388 model. The columns of \mathbf{A} describe the elementary *features* that the model
 389 uses to construct the observed inputs; this forms a kind of visual dictionary.
 390 The rows of \mathbf{W} , on the other hand, represent the best linear *filters* to recover
 391 the weights assigned to the elementary features when generating the observed
 392 inputs. Since \mathbf{W} better reflects the properties of cells’ receptive fields, we
 393 focus primarily on \mathbf{W} throughout the paper.

394 Our objective function allows us to derive an efficient algorithm to learn
 395 the model parameters. Our objective function can be written as:

$$E = \frac{\|\mathbf{X} - \mathbf{A}\mathbf{S}\|_F^2}{2n} + \lambda \|\mathbf{A}\|_1 \quad (4)$$

396 where the columns of \mathbf{X} and \mathbf{S} store the input and output vectors; $\|\cdot\|_F$
 397 denotes taking the Frobenius norm of a matrix (i.e., the square root of the
 398 sum of squared entries); n denotes the number of samples. Our constraint
 399 now is that each row of \mathbf{S} is constrained to have L2 norm less than or equal
 400 to \sqrt{n} .

401 As shown by (Mairal, Bach, Ponce, & Sapiro, 2010), the Sparse Coding
 402 problem (Olshausen & Field, 1996) can be expressed in the matrix factoriza-

tion form as:

$$E = \frac{\|\mathbf{X} - \mathbf{A}\mathbf{S}\|_F^2}{2} + \lambda\|\mathbf{S}\|_1 \quad (5)$$

Each column of \mathbf{A} is constrained to have L2 norm less than or equal to q . Sparse PCA can then utilize Sparse Coding algorithms because the objective function in Eq (4) can be re-written as:

$$E = \frac{\|\mathbf{X}^t - \mathbf{S}^t\mathbf{A}^t\|_F^2}{2n} + \lambda\|\mathbf{A}^t\|_1 \quad (6)$$

Now we see that Eq (6) can be symbolically mapped to Eq (5), if we replace \mathbf{X}^t with \mathbf{X} , \mathbf{S}^t with \mathbf{A} , \mathbf{A}^t with \mathbf{S} , with extra care taken to deal with n and \sqrt{n} .

The above discussion suggests another interpretation of the retinal coding model: it can be interpreted as removing redundancy between input samples, because Sparse Coding is usually interpreted as removing redundancy between input dimensions (Lewicki & Olshausen, 1999; Lewicki & Sejnowski, 2000). In this sense, Architecture-1 ICA (Bartlett, Lades, & Sejnowski, 1998), which applies ICA to \mathbf{X}^t instead of \mathbf{X} , was perhaps the first Sparse PCA algorithm ever proposed.

We then only need to pick some efficient Sparse Coding algorithms to optimize our model parameters. Typically, optimizing the objective function in Eq (5) is factored into two sub-problems: optimizing \mathbf{A} while fixing \mathbf{S} , and optimizing \mathbf{S} while fixing \mathbf{A} . Both sub-problems are convex optimization problems (this is one reason we impose $\langle s_i^2 \rangle \leq 1$ instead of $\langle s_i^2 \rangle = 1$; otherwise optimizing \mathbf{S} is no longer a convex optimization problem). Recently, it was shown that the coordinate descent algorithm is considerably faster than competing methods for both sub-problems (Friedman, Hastie, & Tibshirani, 2010; Mairal et al., 2010). These algorithms are described in Appendix A.

The above algorithm has a computational complexity that depends on the number of samples n . Here we give an approximate algorithm whose complexity only relies on L and M , the input and output dimensionalities. In our experiments on grayscale images, this reduces the computation time from 37 minutes to 10 seconds, with the learned parameters very close to those learned without approximation. The derivation of the algorithm also helps us to see the connection between this retinal coding model and the output decorrelation theory (Atick & Redlich, 1992).

Our algorithm utilizes two approximations, both of which rely on the condition that λ is small and $L > M$. Under such a condition, the first approximation we use is to calculate the optimal output \mathbf{s} using (See Appendix B for detailed derivations):

$$\mathbf{s}^* \approx \mathbf{W}\mathbf{x} \quad (\mathbf{W} = (\mathbf{A}^t \mathbf{A})^{-1} \mathbf{A}^t) \quad (7)$$

Replacing the above approximation into the original objective function, we get

$$E \approx \frac{\text{Tr}((\mathbf{I} - \mathbf{A}\mathbf{W})\mathbf{C}(\mathbf{I} - \mathbf{A}\mathbf{W})^t)}{2} + \lambda \|\mathbf{A}\|_1 \quad (8)$$

where $\mathbf{C} = \langle \mathbf{x}\mathbf{x}^t \rangle$, \mathbf{I} denotes the identity matrix, and Tr denotes taking the trace of a matrix. The constraint $\langle s_i^2 \rangle \leq 1$ can be expressed as $\text{diag}(\mathbf{W}\mathbf{C}\mathbf{W}^t) \leq \mathbf{1}$.

Since $\mathbf{C} = \langle \mathbf{x}\mathbf{x}^t \rangle$ is positive semidefinite, we can factor it using the eigenvalue decomposition $\mathbf{C} = \mathbf{U}\mathbf{V}\mathbf{U}^t$, where \mathbf{U} is a unitary matrix (i.e., $\mathbf{U}\mathbf{U}^t = \mathbf{I}$) containing the eigenvectors as its columns; \mathbf{V} is a diagonal matrix with the eigenvalues on its diagonal. Let $\mathbf{B} = \mathbf{U}\mathbf{V}^{1/2}$, we get $\mathbf{C} = \mathbf{B}\mathbf{B}^t$. Replacing $\mathbf{C} = \mathbf{B}\mathbf{B}^t$ into Eq (8), we get

$$E \approx \frac{\|\mathbf{B}^t - \mathbf{Z}^t \mathbf{A}^t\|_F^2}{2} + \lambda \|\mathbf{A}^t\|_1 \quad (9)$$

where $\mathbf{Z} = \mathbf{W}\mathbf{B}$. The second approximation we use is to relax \mathbf{Z} to a free variable instead of constraining it to $\mathbf{Z} = \mathbf{W}\mathbf{B}$. The constraint becomes that each column of \mathbf{Z}^t should have L2 norm less than or equal to 1.

Now we see that Eq (9) can also be symbolically mapped to Eq (5), if we replace \mathbf{B}^t with \mathbf{X} , \mathbf{Z}^t with \mathbf{A} , and \mathbf{A}^t with \mathbf{S} . Hence the objective function in Eq (9) can also be optimized by efficient Sparse Coding algorithms, such as the coordinate descent algorithms in Appendix A. Since Sparse Coding reduces the redundancy between inputs, the above derivation also brings up another interpretation of the retinal coding model: it can be seen as removing the redundancy between the eigenvectors of \mathbf{C} when λ is small and $L > M$.

Our algorithm efficiently minimizes the objective function. We test two optimization methods: directly optimize the objective function in Eq (6), or first initialize \mathbf{A} by optimizing Eq (9) and then optimize Eq (6). We implement the experiments in single precision on a computer server with Intel Core i7 processors. When we reduce the dimensionality from 400 to

100 and set $\lambda = 0.004$, it takes 37 minutes to directly optimize Eq (6). On the other hand, it takes less than 10 seconds to initialize \mathbf{A} by optimizing Eq (9) and another 40 seconds to optimize Eq (6). As shown in Figure 10(a), our approximate method efficiently minimizes the objective function. In fact, after \mathbf{A} is initialized with the approximate method, further optimizing the objective function with the direct method only causes a less than 0.1% change of the objective function, on average the weights in \mathbf{A} are only changed by 0.4%, and no visible difference in the features and filters can be observed. The closeness of the approximation means that in later experiments (color images, video, and sound) the approximation provided by Eq (9) is used without additional optimization.

Grayscale Images The four image sets used were van Hateren natural images (van Hateren & van der Schaaf, 1998), Kyoto natural images (Doi et al., 2003), Berkeley segmentation dataset (Martin, Fowlkes, Tal, & Malik, 2001), and Caltech-256 object category dataset (Griffin, Holub, & Perona, 2007). We didn’t observe any qualitative difference between the features learned from these datasets. The features reported are obtained using van Hateren natural images. We use a subset of it selected by (Karklin & Lewicki, 2009), which contains 110 1536×1024 grayscale images.

For each image, we discard two pixels off the image borders, normalize the pixel values to $[0, 1]$, then apply a nonlinear function that simulates the cone processing (Karklin & Lewicki, 2009):

$$x = 1 - \exp(-k \cdot x) \quad (10)$$

where x denotes the pixel value, and k is selected for each image such that its average pixel value equals 0.5 after the nonlinearity. This nonlinearity does not seem to alter the features learned by the retinal coding model, but might help to expose higher-order statistical structure (Karklin & Lewicki, 2009). Then we randomly sample 1000 20×20 image patches from each image. Figure 10 compares the distributions of connection weights between PCA and SPCA.

Chromatic Images We applied our algorithm to four chromatic image datasets: the Kyoto natural image dataset (Doi et al., 2003), the McGill color image dataset (Olmos & Kingdom, 2004), the Berkeley segmentation dataset (Martin et al., 2001), and the Caltech-256 object category

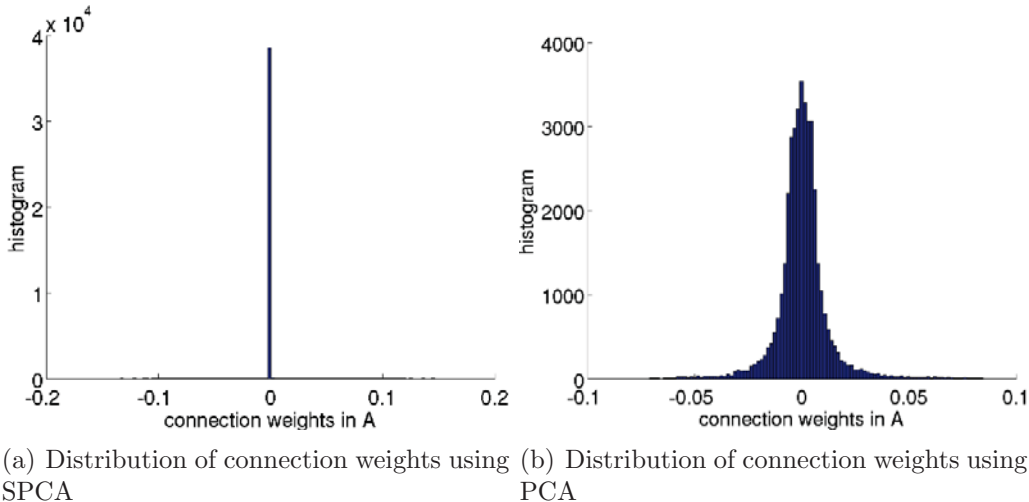


Figure 8: Comparison of SPCA and PCA connection weights on grayscale image patches. Figure 8(a) and 8(b) plot the distribution of the connection weights in \mathbf{A} learned by our model versus those learned by standard PCA.

dataset (Griffin et al., 2007). From these datasets we learn qualitatively similar features. Below we report the features learned from Kyoto image dataset, which contains 62 1000×1280 pixel chromatic images of natural scenes.

First, we estimate the retinal L, M, S cones’ responses to those images. The original images are stored in sRGB color representation (Stokes, Anderson, Chandrasekar, & Motta, 1996). We normalize the pixel values of each image to $[0, 1]$, then transform the image to the CIE XYZ color space (T. Smith & Guild, 1931), from which we estimate the LMS cone responses following the CIECAM02 color appearance model (Fairchild, 2001). In this manner, we can roughly estimate how the retinal L, M, and S cones would respond when presented with the image content.

After that, we apply the cone nonlinearity in Eq (10) to the estimated LMS cone responses. Then we extract all 20 by 20 image patches, and estimate the matrix \mathbf{C} . We apply the retinal coding model using the approximation in Eq (9) to reduce the dimensionality from 1200 to 256 with $\lambda = 0.002$. The resulting model captures 99.75% of the variance captured by an optimal linear model (i.e., standard PCA) with 256 output neurons, with 96.11% of its connections being absolute zero.

515 **Grayscale Videos** To explore the spatio-temporal structure of natural
516 videos, we collected a video dataset of 27 clips from YouTube (see online
517 supporting material for their URLs). All the videos are from natural his-
518 tory shows from the BBC World Wide channel ([www.youtube.com/user/](http://www.youtube.com/user/BBCWorldwide)
519 [BBCWorldwide](http://www.youtube.com/user/BBCWorldwide)). In order to have a realistic sample of natural videos, we
520 eliminate those that are obviously unnatural, such as walking dinosaurs,
521 cavemen, or pigeon-mounted cameras. We calculated the power spectrum
522 of the videos in order to eliminate interlaced-format videos, which have an
523 ellipse-shaped power spectrum elongated in the horizontal direction. We ini-
524 tially applied our algorithm to van Hateren video dataset (van Hateren &
525 Ruderman, 1998). The learned features and filters are qualitatively similar
526 to those reported below, except that some of the learned features are ellipse
527 or even bar shaped and are elongated along the horizontal direction. These
528 features most likely result from the fact that the original videos are inter-
529 laced (i.e., recording only the odd numbered lines in one frame, and the even
530 numbered lines in the next frame). Although the videos were de-interlaced
531 by block averaging with $2 \times 2 \times 2$ (van Hateren, personal communication),
532 their spatial power spectrum is still ellipse-shaped with more energy along
533 the horizontal direction.

534 The original YouTube videos are in color. We transform them to grayscale
535 videos using the method described by the Matlab function `rgb2gray`, normal-
536 ize the pixel values between $[0, 1]$, and apply the cone nonlinearity in Eq 10.
537 Then we estimate the correlation matrix $\mathbf{C} = \mathbf{xx}^t$ for all the $12 \times 12 \times 12$
538 video cubes (with the local mean removed - keeping in the local mean did
539 not substantially change the qualitative results, but it did seem to produce
540 slightly noisier filters), and apply our algorithm to \mathbf{C} using the approximate
541 method.

542 **Sound** We applied our algorithm to three sound datasets: the Pittsburgh
543 natural sounds dataset (E. C. Smith & Lewicki, 2006), the TIMIT speech
544 dataset (Lamel, Kassel, & Seneff, 1986), and rainforest mammal vocalization
545 dataset (Emmons, Whitney, & Ross, n.d.). The Pittsburgh natural sounds
546 dataset contains 48 recordings of natural sound recorded around the Pitts-
547 burgh region, including ambient sounds (such as rain, wind, and streams) and
548 quick acoustic events (such as snapping twigs, breaking wood, and rock im-
549 pacts). The TIMIT speech dataset contains English speech from 630 speak-
550 ers, with each person speaking 10 sentences. The rainforest mammal vocal-

551 ization dataset contains the characteristic sounds of 109 species of rainforest
552 mammals, such as primates, anteaters, bats, jaguars, and manatees.

553 For each dataset, each recording is re-sampled at 16 kHz. We normalize
554 the maximum amplitude of each recording to 1, take all the segments formed
555 using a sliding window of 128 sample points (about a 8 millisecond window),
556 then estimate the matrix \mathbf{C} for the sound segments. We then apply the
557 approximate method to learn the features and filters.

558 References

- 559 Atick, J. J., & Redlich, A. N. (1992). What does the retina know about
560 natural scenes? *Neural Computation*, 4(2), 196–210.
- 561 Attneave, F. (1954). Some informational aspects of visual perception. *Psy-*
562 *chological Review*, 61(3), 183–193.
- 563 Baldi, P., & Hornik, K. (1989). Neural networks and principal component
564 analysis: Learning from examples without local minima. *Neural Net-*
565 *works*, 2, 53–58.
- 566 Barlow, H. B. (1961). Possible principles underlying the transformation of
567 sensory messages. In W. A. Rosenblith (Ed.), *Sensory communication*
568 (pp. 217–234). Cambridge, MA, USA: MIT Press.
- 569 Bartlett, M. S., Lades, M., & Sejnowski, T. J. (1998, January). Independent
570 component representations for face recognition. *Proceedings of the SPIE*
571 *Symposium on Electronic Imaging: Science and Technology; Conference*
572 *on Human Vision and Electronic Imaging III*, 3299.
- 573 Brady, N., & Field, D. J. (1995). What’s constant in contrast constancy?
574 the effects of scaling on the perceived contrast of bandpass patterns.
575 *Vision Research*, 35(6), 739–756.
- 576 Brady, N., & Field, D. J. (2000). local contrast in natural images: normali-
577 sation and coding efficiency. *Perception*, 29(9), 1041–1055.
- 578 Buchsbaum, G., & Gottschalk, A. (1983). Trichromacy, opponent colours
579 coding and optimum colour information transmission in the retina. *Pro-*
580 *ceedings of the Royal Society B: Biological Sciences*, 220(1218), 89–
581 113.
- 582 Calkins, D. J., & Sterling, P. (1999, October). Evidence that circuits for
583 spatial and color vision segregate at the first retinal synapse. *Neuron*,
584 24, 313–321.

- 585 Carney, L. H. (1990). Sensitivities of cells in anteroventral cochlear nucleus
586 of cat to spatiotemporal discharge patterns across primary afferents.
587 *Journal of neurophysiology*, 64(2), 437–456.
- 588 Caywood, M. S., Willmore, B., & Tolhurst, D. J. (2004, June). Independent
589 components of color natural scenes resemble V1 neurons in their spatial
590 and color tuning. *Journal of Neurophysiology*, 91(6), 2859–2873.
- 591 Chichilnisky, E. J. (2001). A simple white noise analysis of neuronal light
592 responses. *Network: Computation in Neural Systems*, 12(2), 199–213.
- 593 Chichilnisky, E. J., & Kalmar, R. S. (2002). Functional asymmetries in ON
594 and OFF ganglion cells of primate retina. *Journal of Neuroscience*,
595 22(7), 2737–2747.
- 596 Cottrell, G. W., Munro, P., & Zipser, D. (1989). Image compression
597 by back propagation: An example of extensional programming. In
598 N. E. Sharkey (Ed.), *Models of cognition: A review of cognitive science*
599 (pp. 208–240). Norwood, New Jersey: Ablex.
- 600 Dearworth Jr., J. R., & Granda, A. M. (2002). Multiplied functions unify
601 shapes of ganglion-cell receptive fields in retina of turtle. *Journal of*
602 *Vision*, 2(3).
- 603 Derrico, J. B., & Buchsbaum, G. (1991). A computational model of spati-
604 ochromatic image coding in early vision. *Journal of Visual Communi-*
605 *cation and Image Representation*, 2(1), 31–38.
- 606 Doi, E., Inui, T., Lee, T.-W., Wachtler, T., & Sejnowski, T. J. (2003).
607 Spatiochromatic receptive field properties derived from information-
608 theoretic analyses of cone mosaic responses to natural scenes. *Neural*
609 *Computation*, 15(2), 397–417.
- 610 Doi, E., & Lewicki, M. S. (2007). A theory of retinal population coding. In
611 *Advances in neural information processing systems* (Vol. 19, pp. 353–
612 360). Cambridge, MA, USA: MIT Press.
- 613 dos Santos, N. A., & Alencar, C. C. G. (2010, August). Early malnutrition
614 diffusely affects children contrast sensitivity to sine-wave gratings of
615 different spatial frequencies. *Nutritional Neuroscience*, 13(4), 189–194.
- 616 Emmons, L. H., Whitney, B. M., & Ross, D. L. (n.d.). *Sounds of neotrop-*
617 *ical rainforest mammals [audio cd] (library of natural sounds, cornell*
618 *laboratory of ornithology, ithaca, new york, 1997).*
- 619 Enroth-Cugell, C., & Robson, J. G. (1966). The contrast sensitivity of retinal
620 ganglion cells of the cat. *Journal of Physiology*, 187(3), 517–552.
- 621 Fairchild, M. D. (2001). A revision of CIECAM97s for practical applications.
622 *Color Research & Application*, 26(6), 418–427.

- 623 Field, D. J. (1987). Relations between the statistics of natural images and
624 the response properties of cortical cells. *Journal of the Optical Society*
625 *of American, A*, 4(12), 2379–2394.
- 626 Field, D. J., & Brady, N. (1997). Visual sensitivity, blur and the sources of
627 variability in the amplitude spectra of natural scenes. *Vision Research*,
628 37(23), 3367–3383.
- 629 Friedman, J., Hastie, T., & Tibshirani, R. (2010). Regularization paths for
630 generalized linear models via coordinate descent. *Journal of Statistical*
631 *Software*, 33, 1–22.
- 632 Graham, D. J., Chandler, D. M., & Field, D. J. (2006). Can the theory of
633 “whitening” explain the center-surround properties of retinal ganglion
634 cell receptive fields? *Vision Research*, 46(18), 2901–2913.
- 635 Griffin, G., Holub, A., & Perona, P. (2007). *Caltech-256 object category*
636 *dataset* (Tech. Rep. No. 7694). California Institute of Technology. Re-
637 trieved from <http://authors.library.caltech.edu/7694>
- 638 Jolliffe, I. T. (2002). *principal component analysis*. Springer-Verlag.
- 639 Karklin, Y., & Lewicki, M. S. (2009). Emergence of complex cell properties
640 by learning to generalize in natural scenes. *Nature*, 457, 83–86.
- 641 Kiang, N. Y.-S., Watanabe, T., Thomas, E. C., & Clark, L. F. (1965).
642 *Discharge patterns of single fibers in the cat’s auditory nerve*. MIT
643 Press Cambridge.
- 644 Lamel, L. F., Kassel, R. H., & Seneff, S. (1986). Speech database de-
645 velopment: design and analysis of the acoustic-phonetic corpus. In
646 *Proceedings of the DARPA speech recognition workshop* (pp. 100–109).
- 647 Lee, T.-W., Wachtler, T., & Sejnowski, T. J. (2002). Color opponency con-
648 stitutes a sparse representation for the chromatic structure of natural
649 scenes. *Vision Research*, 42(17), 2095–2103.
- 650 Lewicki, M. S. (2002). Efficient coding of natural sounds. *Nature Neuro-*
651 *science*, 5(4), 356–363.
- 652 Lewicki, M. S., & Olshausen, B. A. (1999). A probabilistic framework for
653 the adaptation and comparison of image codes. *Journal of the Optical*
654 *Society of America A*, 16(7), 1587–1601.
- 655 Lewicki, M. S., & Sejnowski, T. J. (2000). Learning overcomplete represen-
656 tations. *Neural Computation*, 12(2), 337–365.
- 657 Mairal, J., Bach, F., Ponce, J., & Sapiro, G. (2010). Online learning for
658 matrix factorization and sparse coding. *Journal of Machine Learning*
659 *Research*, 11, 10–60.
- 660 Martin, D., Fowlkes, C., Tal, D., & Malik, J. (2001, July). A database

661 of human segmented natural images and its application to evaluating
662 segmentation algorithms and measuring ecological statistics. In *Pro-*
663 *ceedings of the 8th international conference on computer vision* (Vol. 2,
664 pp. 416–423).

665 Olmos, A. A., & Kingdom, F. A. A. (2004). *Mcgill calibrated colour image*
666 *database*. <http://tabby.vision.mcgill.ca>.

667 Olshausen, B. A., & Field, D. J. (1996). Emergence of simple-cell receptive
668 field properties by learning a sparse code for natural images. *Nature*,
669 *381*, 607–609.

670 Olshausen, B. A., & O’Connor, K. N. (2002). A new window on sound.
671 *Nature*, *5*(4), 292–294.

672 Pillow, J. W., Shlens, J., Paninski, L., Sher, A., Litke, A. M., Chichilnisky,
673 E. J., & Simoncelli, E. P. (2008). Spatio-temporal correlations and
674 visual signalling in a complete neuronal population. *Nature*, *454*(7207),
675 995–999.

676 Recio-Spinoso, A., Temchin, A. N., van Dijk, P., Fan, Y.-H., & Ruggero,
677 M. A. (2005). Wiener-kernel analysis of responses to noise of chinchilla
678 auditory-nerve fibers. *Journal of neurophysiology*, *93*(6), 3615–3634.

679 Smith, E. C., & Lewicki, M. S. (2006). Efficient auditory coding. *Nature*,
680 *439*(7079), 800–805.

681 Smith, T., & Guild, J. (1931). The CIE colorimetric standards and their
682 use. *Transactions of the Optical Society*, *33*, 73–134.

683 Srinivasan, M. V., Laughlin, S. B., & Dubs, A. (1982). Predictive coding: a
684 fresh view of inhibition in the retina. *Proceedings of the Royal Society*
685 *of London. Series B, Biological Sciences*, *216*(1205), 427–459.

686 Stokes, M., Anderson, M., Chandrasekar, S., & Motta, R. (1996). A standard
687 default color space for the Internet–sRGB. *Microsoft and Hewlett-*
688 *Packard Joint Report*.

689 Tailor, D. R., Finkel, L. H., & Buchsbaum, G. (2000). Color-opponent
690 receptive fields derived from independent component analysis of natural
691 images. *Vision Research*, *40*(19), 2671–2676.

692 van Hateren, J. H., & Ruderman, D. L. (1998). Independent component
693 analysis of natural image sequences yields spatio-temporal filters similar
694 to simple cells in primary visual cortex. *Proceedings of the Royal Society*
695 *B: Biological Sciences*, *265*(1412), 2315–2315.

696 van Hateren, J. H., & van der Schaaf, A. (1998). Independent compo-
697 nent filters of natural images compared with simple cells in primary
698 visual cortex. *Proceedings of the Royal Society B: Biological Sciences*,

- 699 265(1394), 359–366.
- 700 Vincent, B. T., & Baddeley, R. J. (2003). Synaptic energy efficiency in retinal
701 processing. *Vision Research*, 43(11), 1283–1290.
- 702 Vincent, B. T., Baddeley, R. J., Troscianko, T., & Gilchrist, I. D. (2005).
703 Is the early visual system optimised to be energy efficient? *Network:
704 Computation in Neural Systems, special issue on Sensory Coding and
705 the Natural Environment*, 16(2/3), 1283–1290.
- 706 Wachtler, T., Lee, T.-W., & Sejnowski, T. J. (2001). Chromatic structure
707 of natural scenes. *Journal of the Optical Society of America A*, 18(1),
708 65–77.
- 709 Zou, H., Hastie, T., & Tibshirani, R. (2006). Sparse principal component
710 analysis. *Journal of computational and graphical statistics*, 15(2), 265–
711 286.

712 Acknowledgements

713 We thank Lingyun Zhang, Wensong Xu, and Chris Kanan for helpful dis-
714 cussions, Ben Vincent for sharing the source code of his model, Yan Karklin
715 for sharing his image preprocessing code, Eizaburo Doi for sharing his code
716 to calculate LMS cone representations, Vivienne Ming and Mike Lewicki for
717 providing the auditory data, Malcolm Slaney for suggesting the experiments
718 on noise images, and Terry Sejnowski for sharing computational resources.

719 **A Coordinate Descent Algorithms**

720 Suppose we want to optimize an objective function $f(\mathbf{x})$, where \mathbf{x} is a high
721 dimensional vector. To do so, the coordinate descent algorithm cyclically
722 optimizes each dimension of \mathbf{x} . Each time, it optimizes one dimension of \mathbf{x}
723 while fixing the other dimensions. Once this dimension is optimized, the algo-
724 rithm optimizes another dimension. This process repeats until the objective
725 function can no longer be optimized.

726 Here we list the coordinate descent algorithm for optimizing the Sparse
727 Coding objective function:

$$E = \left\langle \frac{1}{2} \|\mathbf{x} - \mathbf{A}\mathbf{s}\|_2^2 + \lambda \|\mathbf{s}\|_1 \right\rangle \quad (11)$$

728 Each column of \mathbf{A} is constrained to have L2 norm less than or equal to q .

729 **Optimizing \mathbf{s} While Fixing \mathbf{A}** When we try to optimize the i -th coordi-
730 nate of \mathbf{s} , with \mathbf{A} and all the other coordinates of \mathbf{s} being fixed, the optimal
731 s_i is given by (Friedman et al., 2010):

$$y = \mathbf{a}_i^t \mathbf{x} - \sum_{j \neq i} \mathbf{a}_i^t \mathbf{a}_j s_j \quad (12)$$

$$s_i^* = \begin{cases} (y - \lambda) / \|\mathbf{a}_i\|_2^2, & \text{if } y > \lambda; \\ (y + \lambda) / \|\mathbf{a}_i\|_2^2, & \text{if } y < -\lambda; \\ 0, & \text{otherwise.} \end{cases} \quad (13)$$

732 where \mathbf{a}_i denotes the i -th column of \mathbf{A} .

733 **Optimizing \mathbf{A} While Fixing \mathbf{s}** When we try to optimize the i -th column
734 of \mathbf{A} , with \mathbf{s} and all the other columns of \mathbf{A} being fixed, the optimal \mathbf{a}_i is
735 given by (Mairal et al., 2010):

$$\mathbf{u} = \langle s_i \mathbf{x} \rangle - \sum_{j \neq i} \mathbf{a}_j \langle s_j s_i \rangle \quad (14)$$

$$\mathbf{a}_i^* = \frac{\mathbf{u}}{\max(\langle s_i^2 \rangle, \|\mathbf{u}\|_2 / q)} \quad (15)$$

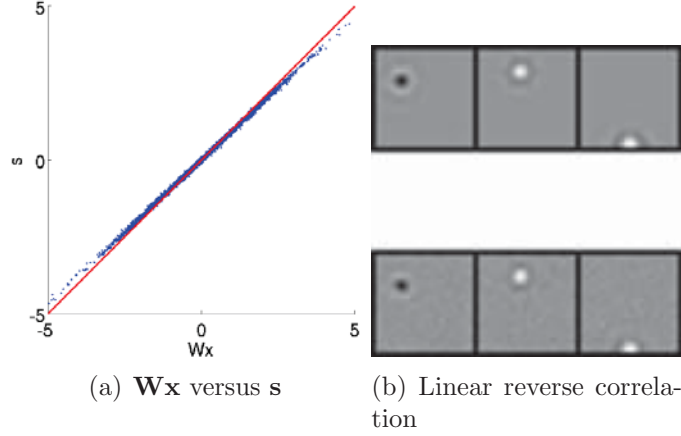


Figure 9: \mathbf{W} approximates the model neurons’ receptive fields. In Figure 9(a), we compare $\mathbf{W}\mathbf{x}$ with the optimal \mathbf{s} inferred using the direct method during the experiments on grayscale image patches. If $\mathbf{s} = \mathbf{W}\mathbf{x}$ holds exactly, all the points should lie on the red solid line. Once we have learned the optimal \mathbf{A} on grayscale image patches, we feed the model with white noise inputs and infer their optimal outputs, then use $\langle \mathbf{x}\mathbf{s}^t \rangle$ to estimate the model neurons’ receptive fields. The top panel in Figure 9(b) displays three rows of \mathbf{W} ; the bottom panel displays the estimated receptive fields for the corresponding neurons.

736 B Derivation of the approximate algorithm

737 The revised retinal coding model aims to minimize the following objective
738 function:

$$E = \left\langle \frac{\|\mathbf{x} - \mathbf{A}\mathbf{s}\|_2^2}{2} \right\rangle + \lambda \|\mathbf{A}\|_1 \quad (16)$$

739 subject to the constraint that

$$\langle s_i^2 \rangle \leq 1 \quad \text{for every } i \quad (17)$$

740 Below we show how to approximate the above objective function when λ
741 is small and $L > M$ (i.e., when we are reducing the data dimensionality).

742 **The First Approximation** Let’s first check the functional roles of the
743 three terms in the objective function: the reconstruction error, the sparsity

744 penalty, and the constraint.

745 If the objective function only contains the reconstruction error term, the
 746 model is reduced to the standard PCA problem (Jolliffe, 2002). The optimal
 747 basis functions (i.e., the columns of \mathbf{A}) should span the subspace spanned
 748 by the eigenvectors of $\langle \mathbf{x}\mathbf{x}^t \rangle$ with top eigenvalues. The optimal outputs are
 749 given by $\mathbf{s}^* = \mathbf{W}\mathbf{x}$, where \mathbf{W} is the pseudo-inverse of \mathbf{A} . Changing the
 750 basis functions' individual directions and lengths within the subspace won't
 751 change the reconstruction error, because for any full-rank matrix $\mathbf{G} \in R^{M \times M}$
 752 we have

$$\mathbf{A}\mathbf{s} = (\mathbf{A}\mathbf{G})(\mathbf{G}^{-1}\mathbf{s}) \quad (18)$$

753 That is, for any $\mathbf{A}_{new} = \mathbf{A}\mathbf{G}$ which spans the same subspace as \mathbf{A} but with
 754 different individual lengths and directions, we can find $\mathbf{s}_{new} = \mathbf{G}^{-1}\mathbf{s}$ such
 755 that the reconstruction error remains to be the minimum.

756 The objective function is not changed by adding the constraint $\langle s_i^2 \rangle \leq 1$.
 757 For any value of \mathbf{A} and \mathbf{s} , we can always divide s_i (the i -th coordinate of
 758 \mathbf{s}) and multiply \mathbf{a}_i (the i -th column of \mathbf{A}) with some value α to satisfy the
 759 constraint without changing the reconstruction error. In other words, the
 760 constraint term only specifies the lengths of the basis functions.

761 The sparsity penalty $\|\mathbf{A}\|_1$ will shrink the basis functions' lengths and
 762 rotate their directions. However, when $\lambda > 0$ is sufficiently small such that
 763 the subspace that the basis functions span is mainly determined by the re-
 764 construction error, the sparsity penalty will only serve to rotate the basis
 765 functions within the subspace determined by the reconstruction error.

766 Our first approximation utilizes the above analysis. When λ is small and
 767 $L > M$, we use $\mathbf{s}^* \approx \mathbf{W}\mathbf{x}$, where \mathbf{W} is the pseudo-inverse of \mathbf{A} . Substituting
 768 it into the objective function in Eq (3), we get:

$$E \approx \left\langle \frac{\|\mathbf{x} - \mathbf{A}(\mathbf{W}\mathbf{x})\|_2^2}{2} \right\rangle + \lambda \|\mathbf{A}\|_1 \quad (19)$$

$$= \left\langle \frac{\text{Tr}((\mathbf{x} - \mathbf{A}\mathbf{W}\mathbf{x})(\mathbf{x} - \mathbf{A}\mathbf{W}\mathbf{x})^t)}{2} \right\rangle + \lambda \|\mathbf{A}\|_1 \quad (20)$$

$$= \frac{\text{Tr}((\mathbf{I} - \mathbf{A}\mathbf{W})\langle \mathbf{x}\mathbf{x}^t \rangle(\mathbf{I} - \mathbf{A}\mathbf{W})^t)}{2} + \lambda \|\mathbf{A}\|_1 \quad (21)$$

$$= \frac{\text{Tr}((\mathbf{I} - \mathbf{A}\mathbf{W})\mathbf{C}(\mathbf{I} - \mathbf{A}\mathbf{W})^t)}{2} + \lambda \|\mathbf{A}\|_1 \quad (\text{let } \mathbf{C} = \langle \mathbf{x}\mathbf{x}^t \rangle) \quad (22)$$

769 where \mathbf{I} denotes the identity matrix, Tr denotes taking the trace of a matrix.

770 The constraint can now be approximated as:

$$\langle s_i^2 \rangle \approx \langle (\mathbf{w}_i \mathbf{x})(\mathbf{w}_i \mathbf{x})^t \rangle = \mathbf{w}_i \langle \mathbf{x} \mathbf{x}^t \rangle \mathbf{w}_i^t \leq 1, \quad \text{or} \quad \text{diag}(\mathbf{W} \mathbf{C} \mathbf{W}^t) \leq \mathbf{1} \quad (23)$$

771 where \mathbf{w}_i denotes the i -th row of \mathbf{W} , and $\mathbf{1}$ denotes a vector of 1's. Hence,
 772 when λ is small and $L > M$, the model mainly serves to capture the second-
 773 order statistical structure of the inputs.

774 **The Second Approximation** Since $\mathbf{C} = \langle \mathbf{x} \mathbf{x}^t \rangle$ is positive semidefinite,
 775 we can factor it using the eigenvalue decomposition $\mathbf{C} = \mathbf{U} \mathbf{V} \mathbf{U}^t$, where \mathbf{U} is
 776 a unitary matrix (i.e., $\mathbf{U} \mathbf{U}^t = \mathbf{I}$) containing the eigenvectors as its columns;
 777 \mathbf{V} is a diagonal matrix with the eigenvalues on its diagonal. Let $\mathbf{B} = \mathbf{U} \mathbf{V}^{1/2}$,
 778 we get $\mathbf{C} = \mathbf{B} \mathbf{B}^t$. Substituting this into the objective function in Eq (22),
 779 yields

$$E = \frac{\text{Tr}((\mathbf{I} - \mathbf{A} \mathbf{W}) \mathbf{B} \mathbf{B}^t (\mathbf{I} - \mathbf{A} \mathbf{W})^t)}{2} + \lambda \|\mathbf{A}\|_1 \quad (24)$$

$$= \frac{\text{Tr}((\mathbf{B} - \mathbf{A} \mathbf{W} \mathbf{B})(\mathbf{B} - \mathbf{A} \mathbf{W} \mathbf{B})^t)}{2} + \lambda \|\mathbf{A}\|_1 \quad (25)$$

$$= \frac{\|\mathbf{B} - \mathbf{A} \mathbf{W} \mathbf{B}\|_F^2}{2} + \lambda \|\mathbf{A}\|_1 \quad (26)$$

$$= \frac{\|\mathbf{B} - \mathbf{A} \mathbf{Z}\|_F^2}{2} + \lambda \|\mathbf{A}\|_1 \quad (\text{let } \mathbf{Z} = \mathbf{W} \mathbf{B}) \quad (27)$$

780 The constraint becomes that each row of \mathbf{Z} should have L2 norm less than
 781 or equal to 1:

$$\text{diag}(\mathbf{W} \mathbf{C} \mathbf{W}^t) = \text{diag}(\mathbf{W} \mathbf{B} \mathbf{B}^t \mathbf{W}^t) = \text{diag}(\mathbf{Z} \mathbf{Z}^t) \leq \mathbf{1} \quad (28)$$

782 Our second approximation is to relax \mathbf{Z} to a free variable instead of con-
 783 straining it to $\mathbf{Z} = \mathbf{W} \mathbf{B}$ because when λ is small and $L > M$, this free
 784 variable will converge to $\mathbf{Z} \approx \mathbf{W} \mathbf{B}$ following the same analysis in our first
 785 approximation. Now the objective function can be written as:

$$E = \frac{\|\mathbf{B}^t - \mathbf{Z}^t \mathbf{A}^t\|_F^2}{2} + \lambda \|\mathbf{A}^t\|_1 \quad (29)$$

786 Each column of \mathbf{Z}^t should have L2 norm less than or equal to 1. We see that
 787 this objective function can also be symbolically mapped to Eq (5). Hence its
 788 parameters can be optimized by efficient Sparse Coding algorithms.

789 In the grayscale image experiment, we verified the efficacy of this ap-
 790 proximation. We tested two optimization methods: directly optimizing the
 791 objective function in Eq (6), or first initialize \mathbf{A} by optimizing Eq (9) and
 792 then optimize Eq (6). We implemented the experiments in single precision on
 793 a computer server with Intel Core i7 processors. When we reduce the dimen-
 794 sionality from 400 to 100 and set $\lambda = 0.004$, it takes 37 minutes to directly
 795 optimize Eq (6). On the other hand, it takes less than 10 seconds to initialize
 796 \mathbf{A} by optimizing Eq (9) and another 40 seconds to optimize Eq (6). As shown
 797 in Figure 10(a), our approximate method efficiently minimizes the objective
 798 function. In fact, after \mathbf{A} is initialized with the approximate method, further
 799 optimizing the objective function with the direct method only causes a less
 800 than 0.1% change of the objective function, and on average, the weights in
 801 \mathbf{A} are only changed by 0.4%. No visible difference in the features and filters
 802 can be observed. Because of the closeness of the approximation, in later
 803 experiments (color images, video, and sound), we used the approximation
 804 provided by Eq (9) without additional optimization.

805 C URLs of Video Clips

806 www.youtube.com/watch?v=cMIRwCNvI94 www.youtube.com/watch?v=J7eRGHVx3p0
 807 www.youtube.com/watch?v=8R1g0t00vGM www.youtube.com/watch?v=K61FRGpvmfE
 808 www.youtube.com/watch?v=M-nXN5SGmhw
 809 www.youtube.com/watch?v=tOn2RhH36Mc www.youtube.com/watch?v=gc9jFmkjizQ
 810 www.youtube.com/watch?v=ZiW96Uci624 www.youtube.com/watch?v=1YQrLPW5DdY
 811 www.youtube.com/watch?v=xKksJ3fvB1Q
 812 www.youtube.com/watch?v=43id_NRajDo www.youtube.com/watch?v=NRWehNVSA1A
 813 www.youtube.com/watch?v=oJ-KzdRsQC4 www.youtube.com/watch?v=VuMRDZbrdXc
 814 www.youtube.com/watch?v=2rlZVtKKWnk
 815 www.youtube.com/watch?v=aIQB0NFcFog www.youtube.com/watch?v=B71T_GpA2AM
 816 www.youtube.com/watch?v=XB-8Grn6sRo www.youtube.com/watch?v=JxrIWShNPko
 817 www.youtube.com/watch?v=gVjqL-9Fh3E
 818 www.youtube.com/watch?v=yKKabd3W904 www.youtube.com/watch?v=4ZFoqh8PQ88
 819 www.youtube.com/watch?v=NR3Z4p5hspI www.youtube.com/watch?v=RB9uzMjiYSQ
 820 www.youtube.com/watch?v=a7XuXi3mqYM
 821 www.youtube.com/watch?v=PBrStxu0Jbs www.youtube.com/watch?v=u6ouWOGJk5E
 822

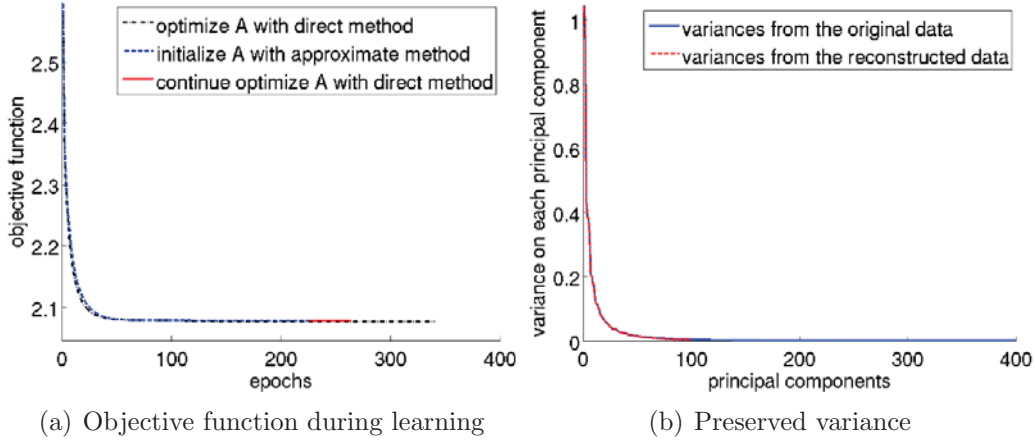


Figure 10: Experiments on grayscale image patches. Figure 10(a) plots how the objective function changes during learning. The black dash-dotted line plots the objective function using the direct optimization method; the blue dashed line plots the objective function when we use the approximate method to find \mathbf{A} ; the red solid line plots the objective function when we further fine tune the parameters using the direct method after initializing \mathbf{A} with the approximate method. Figure 10(b) plots the eigenvalues of \mathbf{C} from the original data (the blue solid line) versus those from the reconstructed data (the red dashed line). Our model captures 99.23% of the variance that could be captured by an optimal linear model with 100 output neurons.

References

- Atick, J. J., & Redlich, A. N. (1992). What does the retina know about natural scenes? *Neural Computation*, 4(2), 196–210.
- Attneave, F. (1954). Some informational aspects of visual perception. *Psychological Review*, 61(3), 183–193.
- Baldi, P., & Hornik, K. (1989). Neural networks and principal component analysis: Learning from examples without local minima. *Neural Networks*, 2, 53–58.
- Barlow, H. B. (1961). Possible principles underlying the transformation of sensory messages. In W. A. Rosenblith (Ed.), *Sensory communication* (pp. 217–234). Cambridge, MA, USA: MIT Press.
- Bartlett, M. S., Lades, M., & Sejnowski, T. J. (1998, January). Independent component representations for face recognition. *Proceedings of the SPIE Symposium on Electronic Imaging: Science and Technology; Conference on Human Vision and Electronic Imaging III*, 3299.
- Brady, N., & Field, D. J. (1995). What’s constant in contrast constancy? the effects of scaling on the perceived contrast of bandpass patterns. *Vision Research*, 35(6), 739–756.
- Brady, N., & Field, D. J. (2000). local contrast in natural images: normalisation and coding efficiency. *Perception*, 29(9), 1041–1055.
- Buchsbaum, G., & Gottschalk, A. (1983). Trichromacy, opponent colours coding and optimum colour information transmission in the retina. *Proceedings of the Royal Society B: Biological Sciences*, 220(1218), 89–113.
- Calkins, D. J., & Sterling, P. (1999, October). Evidence that circuits for spatial and color vision segregate at the first retinal synapse. *Neuron*, 24, 313–321.
- Carney, L. H. (1990). Sensitivities of cells in anteroventral cochlear nucleus of cat to spatiotemporal discharge patterns across primary afferents. *Journal of neurophysiology*, 64(2), 437–456.
- Caywood, M. S., Willmore, B., & Tolhurst, D. J. (2004, June). Independent components of color natural scenes resemble V1 neurons in their spatial and color tuning. *Journal of Neurophysiology*, 91(6), 2859–2873.
- Chichilnisky, E. J. (2001). A simple white noise analysis of neuronal light responses. *Network: Computation in Neural Systems*, 12(2), 199–213.
- Chichilnisky, E. J., & Kalmar, R. S. (2002). Functional asymmetries in ON and OFF ganglion cells of primate retina. *Journal of Neuroscience*,

- 860 22(7), 2737–2747.
- 861 Cottrell, G. W., Munro, P., & Zipser, D. (1989). Image compression
862 by back propagation: An example of extensional programming. In
863 N. E. Sharkey (Ed.), *Models of cognition: A review of cognitive science*
864 (pp. 208–240). Norwood, New Jersey: Ablex.
- 865 Dearworth Jr., J. R., & Granda, A. M. (2002). Multiplied functions unify
866 shapes of ganglion-cell receptive fields in retina of turtle. *Journal of*
867 *Vision*, 2(3).
- 868 Derrico, J. B., & Buchsbaum, G. (1991). A computational model of spati-
869 ochromatic image coding in early vision. *Journal of Visual Communi-*
870 *cation and Image Representation*, 2(1), 31–38.
- 871 Doi, E., Inui, T., Lee, T.-W., Wachtler, T., & Sejnowski, T. J. (2003).
872 Spatiochromatic receptive field properties derived from information-
873 theoretic analyses of cone mosaic responses to natural scenes. *Neural*
874 *Computation*, 15(2), 397–417.
- 875 Doi, E., & Lewicki, M. S. (2007). A theory of retinal population coding. In
876 *Advances in neural information processing systems* (Vol. 19, pp. 353–
877 360). Cambridge, MA, USA: MIT Press.
- 878 dos Santos, N. A., & Alencar, C. C. G. (2010, August). Early malnutrition
879 diffusely affects children contrast sensitivity to sine-wave gratings of
880 different spatial frequencies. *Nutritional Neuroscience*, 13(4), 189–194.
- 881 Emmons, L. H., Whitney, B. M., & Ross, D. L. (n.d.). *Sounds of neotrop-*
882 *ical rainforest mammals [audio cd] (library of natural sounds, cornell*
883 *laboratory of ornithology, ithaca, new york, 1997).*
- 884 Enroth-Cugell, C., & Robson, J. G. (1966). The contrast sensitivity of retinal
885 ganglion cells of the cat. *Journal of Physiology*, 187(3), 517–552.
- 886 Fairchild, M. D. (2001). A revision of CIECAM97s for practical applications.
887 *Color Research & Application*, 26(6), 418–427.
- 888 Field, D. J. (1987). Relations between the statistics of natural images and
889 the response properties of cortical cells. *Journal of the Optical Society*
890 *of American, A*, 4(12), 2379–2394.
- 891 Field, D. J., & Brady, N. (1997). Visual sensitivity, blur and the sources of
892 variability in the amplitude spectra of natural scenes. *Vision Research*,
893 37(23), 3367–3383.
- 894 Friedman, J., Hastie, T., & Tibshirani, R. (2010). Regularization paths for
895 generalized linear models via coordinate descent. *Journal of Statistical*
896 *Software*, 33, 1–22.
- 897 Graham, D. J., Chandler, D. M., & Field, D. J. (2006). Can the theory of

898 “whitening” explain the center-surround properties of retinal ganglion
899 cell receptive fields? *Vision Research*, 46(18), 2901–2913.

900 Griffin, G., Holub, A., & Perona, P. (2007). *Caltech-256 object category*
901 *dataset* (Tech. Rep. No. 7694). California Institute of Technology. Re-
902 trieved from <http://authors.library.caltech.edu/7694>

903 Jolliffe, I. T. (2002). *principal component analysis*. Springer-Verlag.

904 Karklin, Y., & Lewicki, M. S. (2009). Emergence of complex cell properties
905 by learning to generalize in natural scenes. *Nature*, 457, 83–86.

906 Kiang, N. Y.-S., Watanabe, T., Thomas, E. C., & Clark, L. F. (1965).
907 *Discharge patterns of single fibers in the cat’s auditory nerve*. MIT
908 Press Cambridge.

909 Lamel, L. F., Kassel, R. H., & Seneff, S. (1986). Speech database de-
910 velopment: design and analysis of the acoustic-phonetic corpus. In
911 *Proceedings of the DARPA speech recognition workshop* (pp. 100–109).

912 Lee, T.-W., Wachtler, T., & Sejnowski, T. J. (2002). Color opponency con-
913 stitutes a sparse representation for the chromatic structure of natural
914 scenes. *Vision Research*, 42(17), 2095–2103.

915 Lewicki, M. S. (2002). Efficient coding of natural sounds. *Nature Neuro-*
916 *science*, 5(4), 356–363.

917 Lewicki, M. S., & Olshausen, B. A. (1999). A probabilistic framework for
918 the adaptation and comparison of image codes. *Journal of the Optical*
919 *Society of America A*, 16(7), 1587–1601.

920 Lewicki, M. S., & Sejnowski, T. J. (2000). Learning overcomplete represen-
921 tations. *Neural Computation*, 12(2), 337–365.

922 Mairal, J., Bach, F., Ponce, J., & Sapiro, G. (2010). Online learning for
923 matrix factorization and sparse coding. *Journal of Machine Learning*
924 *Research*, 11, 10–60.

925 Martin, D., Fowlkes, C., Tal, D., & Malik, J. (2001, July). A database
926 of human segmented natural images and its application to evaluating
927 segmentation algorithms and measuring ecological statistics. In *Pro-*
928 *ceedings of the 8th international conference on computer vision* (Vol. 2,
929 pp. 416–423).

930 Olmos, A. A., & Kingdom, F. A. A. (2004). *Mcgill calibrated colour image*
931 *database*. <http://tabby.vision.mcgill.ca>.

932 Olshausen, B. A., & Field, D. J. (1996). Emergence of simple-cell receptive
933 field properties by learning a sparse code for natural images. *Nature*,
934 381, 607–609.

- 935 Olshausen, B. A., & O'Connor, K. N. (2002). A new window on sound.
936 *Nature*, 5(4), 292–294.
- 937 Pillow, J. W., Shlens, J., Paninski, L., Sher, A., Litke, A. M., Chichilnisky,
938 E. J., & Simoncelli, E. P. (2008). Spatio-temporal correlations and
939 visual signalling in a complete neuronal population. *Nature*, 454(7207),
940 995–999.
- 941 Recio-Spinoso, A., Temchin, A. N., van Dijk, P., Fan, Y.-H., & Ruggero,
942 M. A. (2005). Wiener-kernel analysis of responses to noise of chinchilla
943 auditory-nerve fibers. *Journal of neurophysiology*, 93(6), 3615–3634.
- 944 Smith, E. C., & Lewicki, M. S. (2006). Efficient auditory coding. *Nature*,
945 439(7079), 800–805.
- 946 Smith, T., & Guild, J. (1931). The CIE colorimetric standards and their
947 use. *Transactions of the Optical Society*, 33, 73–134.
- 948 Srinivasan, M. V., Laughlin, S. B., & Dubs, A. (1982). Predictive coding: a
949 fresh view of inhibition in the retina. *Proceedings of the Royal Society
950 of London. Series B, Biological Sciences*, 216(1205), 427–459.
- 951 Stokes, M., Anderson, M., Chandrasekar, S., & Motta, R. (1996). A standard
952 default color space for the Internet-sRGB. *Microsoft and Hewlett-
953 Packard Joint Report*.
- 954 Tailor, D. R., Finkel, L. H., & Buchsbaum, G. (2000). Color-opponent
955 receptive fields derived from independent component analysis of natural
956 images. *Vision Research*, 40(19), 2671–2676.
- 957 van Hateren, J. H., & Ruderman, D. L. (1998). Independent component
958 analysis of natural image sequences yields spatio-temporal filters similar
959 to simple cells in primary visual cortex. *Proceedings of the Royal Society
960 B: Biological Sciences*, 265(1412), 2315–2315.
- 961 van Hateren, J. H., & van der Schaaf, A. (1998). Independent compo-
962 nent filters of natural images compared with simple cells in primary
963 visual cortex. *Proceedings of the Royal Society B: Biological Sciences*,
964 265(1394), 359–366.
- 965 Vincent, B. T., & Baddeley, R. J. (2003). Synaptic energy efficiency in retinal
966 processing. *Vision Research*, 43(11), 1283–1290.
- 967 Vincent, B. T., Baddeley, R. J., Troschianko, T., & Gilchrist, I. D. (2005).
968 Is the early visual system optimised to be energy efficient? *Network:
969 Computation in Neural Systems, special issue on Sensory Coding and
970 the Natural Environment*, 16(2/3), 1283–1290.
- 971 Wachtler, T., Lee, T.-W., & Sejnowski, T. J. (2001). Chromatic structure
972 of natural scenes. *Journal of the Optical Society of America A*, 18(1),

973 65–77.
974 Zou, H., Hastie, T., & Tibshirani, R. (2006). Sparse principal component
975 analysis. *Journal of computational and graphical statistics*, 15(2), 265–
976 286.