

Individual variability of neural computations in the primate retina

Nishal Shah^{1,2*}, Nora Brackbill^{2,3}, Ryan Samarakoon^{2,4}, Colleen Rhoades^{2,5}, Alexandra Kling^{2,4}, Alexander Sher⁶, Alan Litke⁶, Yoram Singer⁷, Jonathon Shlens⁸, and E.J. Chichilnisky^{2,4}

1 Stanford University, Department of Electrical Engineering

2 Hansen Experimental Physics Lab, Stanford University,

3 Department of Physics, Stanford University

4 Department of Neurosurgery, Stanford University

5 Department of Bioengineering, Stanford University

6 University of California, Santa Cruz

7 WorldQuant, LLC

8 Google Brain

Abstract

Variation in the neural code between individuals contributes to making each person unique. Using ~100 neural population recordings from major ganglion cell types in the macaque retina, we develop an interpretable computational representation of individual variability using machine learning. This representation preserves invariances, such as asymmetries between ON and OFF cells, while capturing individual variation and covariation in properties such as nonlinearity, temporal dynamics, and spatial receptive field size. The similarity of these properties across cell types was dependent on the similarity of their synaptic connections. Surprisingly, male retinas exhibited higher firing rates and faster temporal integration than female retinas. By exploiting data from previously recorded macaque retinas, a new macaque retina (and crucially, a human retina) could be efficiently characterized. Simulations indicated that combining a vast dataset of healthy macaque recordings with behavioral feedback could be used to identify the neural code and improve retinal implants for treating blindness.

Main

An emerging frontier in biomedicine is understanding variability between individuals, with implications ranging from the mathematical modeling of living systems to ethics and personalized medicine. In neuroscience, differences in mental function between individuals are substantial, yet little is known about the underlying variation in the information processing performed by neural circuits. Neuroimaging, gross structural, and behavioral measurements in humans cannot reveal the neural code at circuit resolution, and physiological measurements in invertebrates or rodents have uncertain applicability to humans. These limitations have led to a large gap in our understanding of variability in the neural code and its implications for translational medicine and neuroengineering. Two technical challenges have limited our understanding of this variability in higher animal models: high-resolution, large-scale physiological recordings from many animals, and methods for deciphering variability in complex circuit level computations.

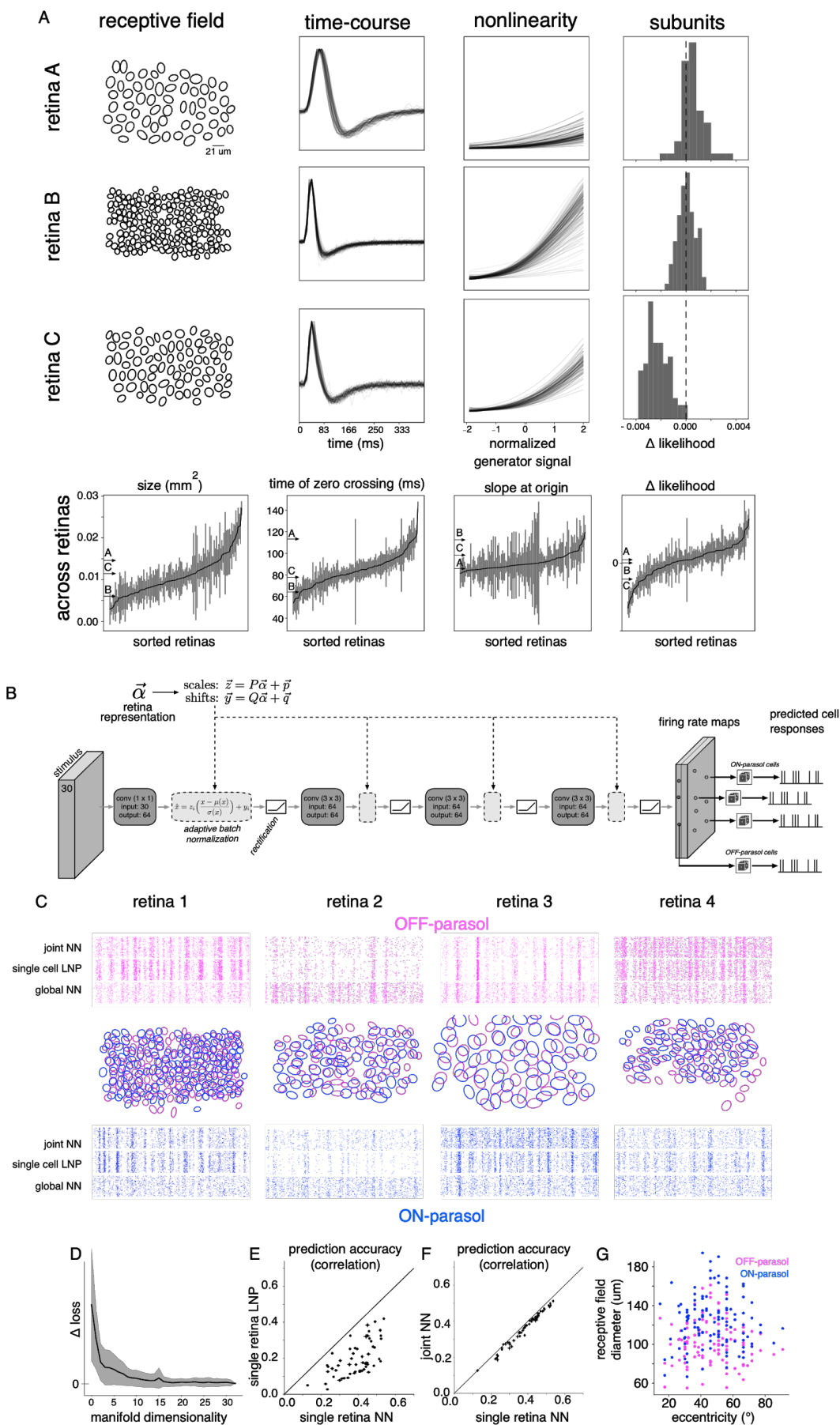
1 **Modeling the shared and individual components of neural coding variability**

2
3 To overcome these limitations, we exploited large-scale multi-electrode recordings from ~21626
4 retinal ganglion cells (RGCs) aggregated over 112 recordings from 66 isolated primate retinas,
5 in which the functional properties of diverse cell types have been extensively studied¹⁻⁵. These
6 data, gathered over a decade of experimentation, exhibited significant neural coding variability
7 across recordings. As a baseline, response properties in each recording were summarized by
8 the parameters of a linear-nonlinear-Poisson (LNP) encoding model. This widely used model⁶
9 captures light-evoked responses in RGCs using a spatiotemporal linear filter applied to the
10 stimulus, followed by an output nonlinearity and stochastic spike generation. Receptive field
11 sizes of two known RGC types -- ON and OFF parasol -- exhibited substantial variation across
12 recordings. This variability was evident across retinal eccentricities, but was also present at a
13 given eccentricity (Figure 1A, first column, Figure 1G). Diversity was also seen in the kinetics of
14 light response (Figure 1A, second column) and in the form of the output nonlinearity (Figure 1A,
15 third column). When the LNP model was extended to capture an additional computational motif,
16 nonlinear integration of spatial subunits within the receptive field⁷, the change in prediction
17 accuracy was also variable across recordings (Figure 1A, last column).

18
19 The structure and covariation of these response properties was therefore explored using a
20 flexible model that captured nonlinear spatial integration, and combined shared and recording-
21 specific parameters. The shared component was a multilayered convolutional neural network
22 (CNN), an extension of the LNP model, consisting of multiple alternating stages of spatio-
23 temporal filtering, normalization and rectification. The rectification captured nonlinear spatial
24 integration, and the convolutional structure captured the known translational invariance of visual
25 signals in each cell type (cells of the same type at different locations have very similar response
26 properties⁸). The model output consisted of multiple firing rate maps in each retina, one for each
27 cell type. To predict a given cell's responses, the model-predicted firing rate was read off from
28 the map at the cell's location (Figure 1B). Due to the translational invariance constraint, the
29 proposed model cannot capture differences between cells belonging to the same cell type,
30 resulting in poorer prediction accuracy compared to models that allow for cell-specific
31 parameters such as single-cell LNP⁶ or other state-of-the-art models^{9,10}. However, this
32 constraint enabled the proposed neural network architecture to predict responses across
33 different recordings with different numbers of cells, while focusing on variation across
34 recordings. When trained using ON and OFF parasol cell responses in each retina separately,
35 the CNN model exhibited performance superior to the single-retina LNP model (LNP model with
36 common parameters for all cells of a given type), as expected given its more flexible structure
37 (Figure 1E). However, when the CNN model was trained on multiple retinas together, it failed to
38 capture the responses of retinas with low firing rates, highly modulated responses, or other
39 features that varied between recordings (Figure 1C, rasters).

40
41 To capture the variation of light response properties across recordings in a compact and
42 tractable way, a small number of recording-specific parameters were used to reweight the
43 activations of different filters at each layer of the shared CNN. Data from all recordings were
44 used to learn the ~100K parameters of the shared CNN, while data from each recording were

1 used to obtain a small number recording-specific parameters. The collection of these recording-
2 specific parameters were interpreted as a *manifold of neural coding variability*. When learned
3 using 71 recordings, the low-dimensional manifold captured variations in background firing rate,
4 sustained vs. transient dynamics, and response nonlinearities (Figure 1C, rasters). The ability to
5 simultaneously predict responses across multiple retinas saturated at ~15 dimensions of the
6 learned manifold (Figure 1D), much lower than the total number of CNN parameters, and the
7 performance of this joint model based on the manifold was only slightly lower than a CNN model
8 trained for each retina separately (Figure 1F). Thus, a simple, low-dimensional representation
9 can efficiently and accurately capture the diversity of retinal computations.
10



1 **Figure 1. Modeling variability in the neural code** (A) Variability of response properties across preparations. Spatial
2 receptive fields (first column), time of zero-crossing of temporal filter (second column), nonlinearity (third column) and
3 change in response likelihood with five nonlinear subunits (fourth column) for OFF parasol cells in three
4 representative recordings (rows, same y-axis for each response property across recordings). The last row represents
5 the range of response parameters across 122 preparations, sorted according to their population means, and error
6 bars corresponding to the robust standard deviation. Arrows indicate the values of the chosen retinas in the first three
7 rows. (B) Architecture of the neural network for capturing response variation. The visual stimulus is passed through
8 multiple layers of convolution with spatial filters, with adaptive batch normalization and rectification at each layer,
9 producing two firing rate maps (one each for ON and OFF parasol cell types). The Poisson firing rate for each cell is
10 read off from the value at the cell's location in the firing rate map of its cell type. Retina-specific tuning of responses is
11 performed by adjusting the mean and standard deviation of the activation values at each layer, determined by a linear
12 transformation of the retina's location in the low-dimensional manifold. (C) Response prediction across 4
13 representative training retinas (columns). The receptive field mosaics are shown for each retina (middle row), along
14 with response predictions for a randomly selected OFF parasol (top row) and ON parasol cell (bottom row). Rasters
15 (60 trials) for predicted responses to a 3 sec long white noise stimulus using the LNP model; neural network model,
16 trained jointly on multiple retinas, with retina-specific parameters (15 dimensional manifold, joint NN) or without them
17 (no manifold, global NN). (D) Model error (log likelihood) on test stimuli with varying manifold dimensionality;
18 dimensions=0 indicates no retina-specific adaptation (global NN). Error bars indicate the standard deviation across
19 retinas. (E) Prediction accuracy, averaged across cells, for different retinas (points), using a neural network trained on
20 data from each individual retina (x-axis) and an LNP model with shared parameters across cells of a given type in
21 each retina (y-axis). Prediction accuracy is measured as correlation between predicted firing rate and recorded
22 responses smoothed with a Gaussian filter (σ : 11ms). (F) Similar to (E), comparing predictions from the neural
23 network with a 15-dimensional manifold, trained jointly (y-axis) vs. trained on each retina separately (x-axis). (G)
24 Range of eccentricities (x-axis) and the average receptive field sizes (y-axis) for ON parasol (blue) and OFF parasol
25 (magenta) cells across 122 preparations used in this study.

26 **Neural coding manifold smoothly captures systematic variation across recordings**

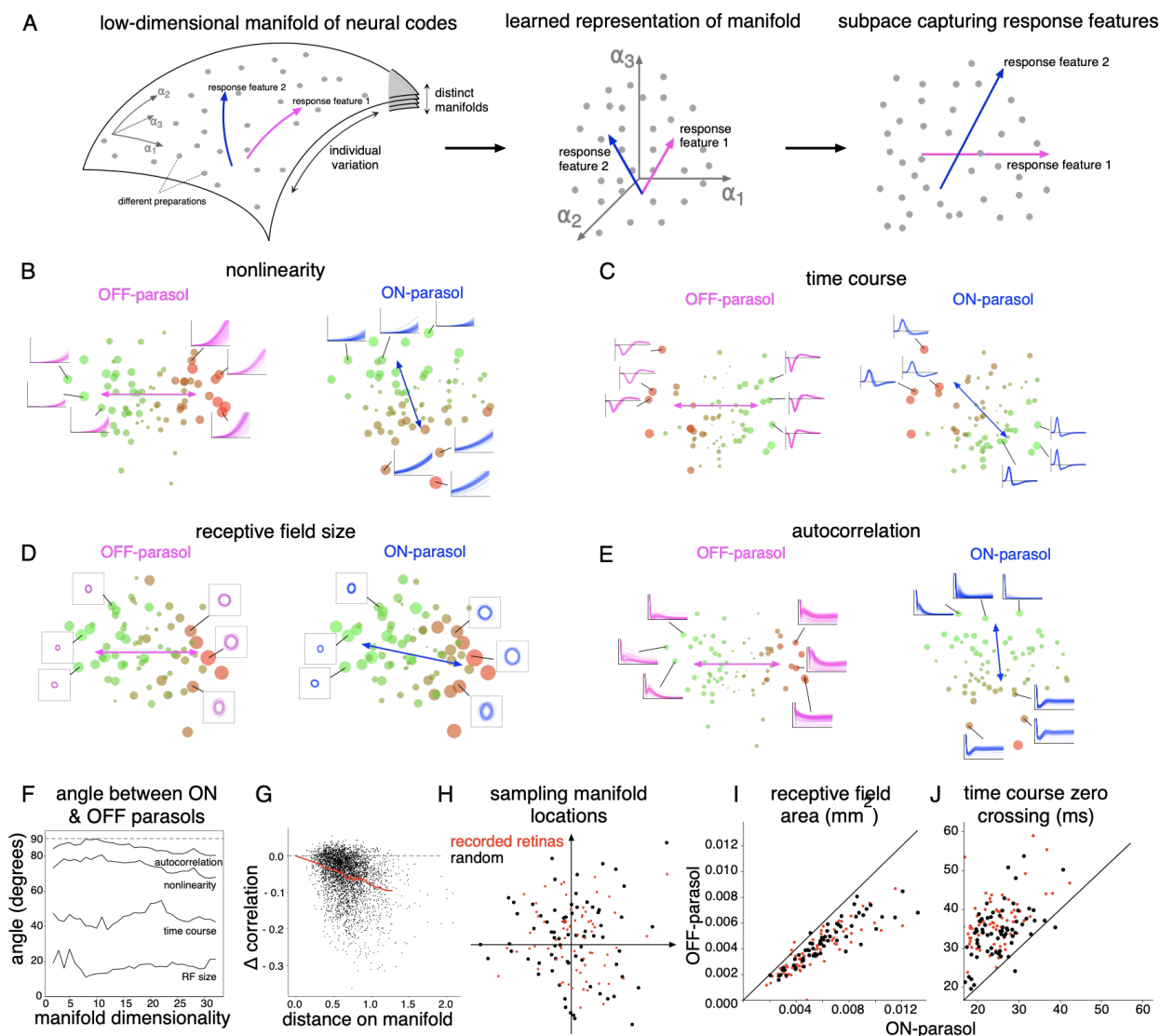
27 The learned manifold smoothly captured changes in neural code with a greater perturbation in
28 manifold location leading to a greater decrease in response prediction accuracy (Figure 2G).
29 The manifold geometrically represented variation in several light response properties, including
30 receptive field size, time course, output nonlinearity, and spike train autocorrelation. This was
31 observed by projecting the average response property for each retina onto its principal
32 component across all recordings, and then identifying the manifold direction with maximum
33 correlation to this projection. The Spearman rank correlation between these projections and the
34 projections along the identified manifold direction was significantly higher than the value
35 observed in random permutations of the data ($p < 0.001$ for all response properties) (Figure 2B,
36 C, D, E, F), indicating that the geometry of the manifold is well-suited to representing the
37 response properties linearly.

38 The geometry of the manifold also captured co-variation in response properties of different cell
39 types across recordings. For both spiking autocorrelation and response nonlinearity, the large
40 angles between the manifold directions for ON and OFF parasol cells (86° and 76° respectively)
41 were consistent with low Spearman rank correlation in these response properties (0.28 and 0.12
42 respectively). Conversely, for response time course and receptive field size, the small angles
43 between the directions for ON and OFF cells (45° and 16° respectively) were consistent with
44 larger Spearman rank correlation (0.56 and 0.95 respectively). Although high correlation is
45 expected for receptive field size based on variation in the eccentricity of different recordings, the
46 moderate correlation observed for response time course and the low correlation observed for

1 nonlinearity and autocorrelation have not been previously reported. While this analysis could be
 2 done directly in the space of LNP model parameters, the manifold presents a more flexible way
 3 to study these response properties for a wide range of encoding models.

4 The manifold also captured known *invariances* in retinal encoding, i.e. properties that were
 5 consistent across retinas. To assess whether these invariances were present at many
 6 intermediate manifold locations not directly sampled in the experiments, the manifold locations
 7 of recorded retinas were perturbed using Gaussian noise with standard deviation equal to the
 8 median nearest neighbor separation (Figure 2H). Light responses were then generated using
 9 these randomly sampled manifold locations, and the encoding properties were summarized by
 10 fitting a LNP model. For the simulated responses, OFF parasol cells had consistently smaller
 11 RFs (Figure 2I) and slower time courses (Figure 2J) than ON parasol cells, as expected based
 12 on previously reported asymmetries⁸.

13



14

1 Figure 2. **Geometry of manifold.** (A) Summary of the steps in subsequent analysis. Left: Schematic of the manifold
2 representation of variability. Each recording is summarized by its neural encoding function, indicated by a point (gray)
3 in space of all possible neural encoding functions. The observed neural encoding functions lie in a low-dimensional
4 manifold (depicted as curved surface, but could be more than two dimensional) within this space. Different manifolds
5 (other surfaces) would potentially correspond to different properties that are invariant across recordings. The training
6 procedure learns a coordinate system (α) within the manifold. Middle: Directions corresponding to response features
7 can be identified in the learned coordinate system for representing the manifold. Right: Geometry of the subspace
8 corresponding to the identified directions lead to interpretation of the variations. (B) Manifold locations for 95
9 preparations (points), projected onto the 2D subspace given by the first principal components of variation in the
10 output nonlinearity for OFF and ON parasol cells. Size and color of dots indicate deviation from the mean. Colored
11 lines indicate the direction of maximum nonlinearity variation for ON parasol (blue) and OFF parasol (magenta) cells.
12 Insets: Lines show output nonlinearities for all cells in representative retinas. (C, D, E) Similar to (B), for time course,
13 receptive field size and autocorrelation, respectively. (F) Angle between ON and OFF parasol directions for particular
14 response properties, as a function of manifold dimension. (G) Change in response prediction accuracy (y-axis) as the
15 manifold location is perturbed from the learned location. Each black dot represents a different perturbation, the red
16 line is the average. (H) Random manifold locations (red) were sampled by adding noise to the learned retina-specific
17 locations (black). Responses to a white noise stimulus of 100 cells of each type (ON and OFF parasol) were sampled
18 from random locations in the visual field using the firing rate maps for the two cell types associated with these
19 sampled manifold locations. These simulated responses were then used to fit a LNP model. (I, J) Relationship
20 between ON and OFF parasol cells for receptive field area and zero crossing of response time course, respectively,
21 for recorded (red) and randomly sampled (black) retinas

22 **Manifold reveals covariation associated with connectivity, and male-female differences**

23 The manifold of variability revealed two novel properties of retinal encoding. First, the RGC
24 types that receive synaptic input from bipolar cells at similar depths in the inner plexiform layer
25 (IPL) showed greater covariation in their response properties across recordings. To examine
26 covariation between cell types, the ON and OFF midget cell types were included with the ON
27 and OFF parasol cell types considered thus far. In 85 recordings (53 macaques), the similarity
28 of three response properties – nonlinearity, time course and autocorrelation – across different
29 pairs of cell types was measured either directly, or in the manifold (Figure 3A-C). Using both
30 methods, the highest correlation in these physiological properties was observed between cell-
31 type pairs with the same polarity (ON or OFF), consistent with the lamination of ON and OFF
32 cells in the inner and outer IPL, respectively (Figure 3D). Moreover, for cell type pairs with
33 opposite polarities, a higher correlation in physiological properties was observed for the ON-
34 parasol/OFF-parasol pair than for the ON-midget/OFF-midget pair, consistent with the
35 lamination of parasol cells closer to the middle of IPL (Figure 3A-D). These observations
36 support the approach of studying the response properties of newer cell types (such as ON and
37 OFF smooth monostratified cells⁴) after normalizing to the properties of more commonly studied
38 cells with similar synaptic inputs, thereby minimizing the effects of inter-retina variability.

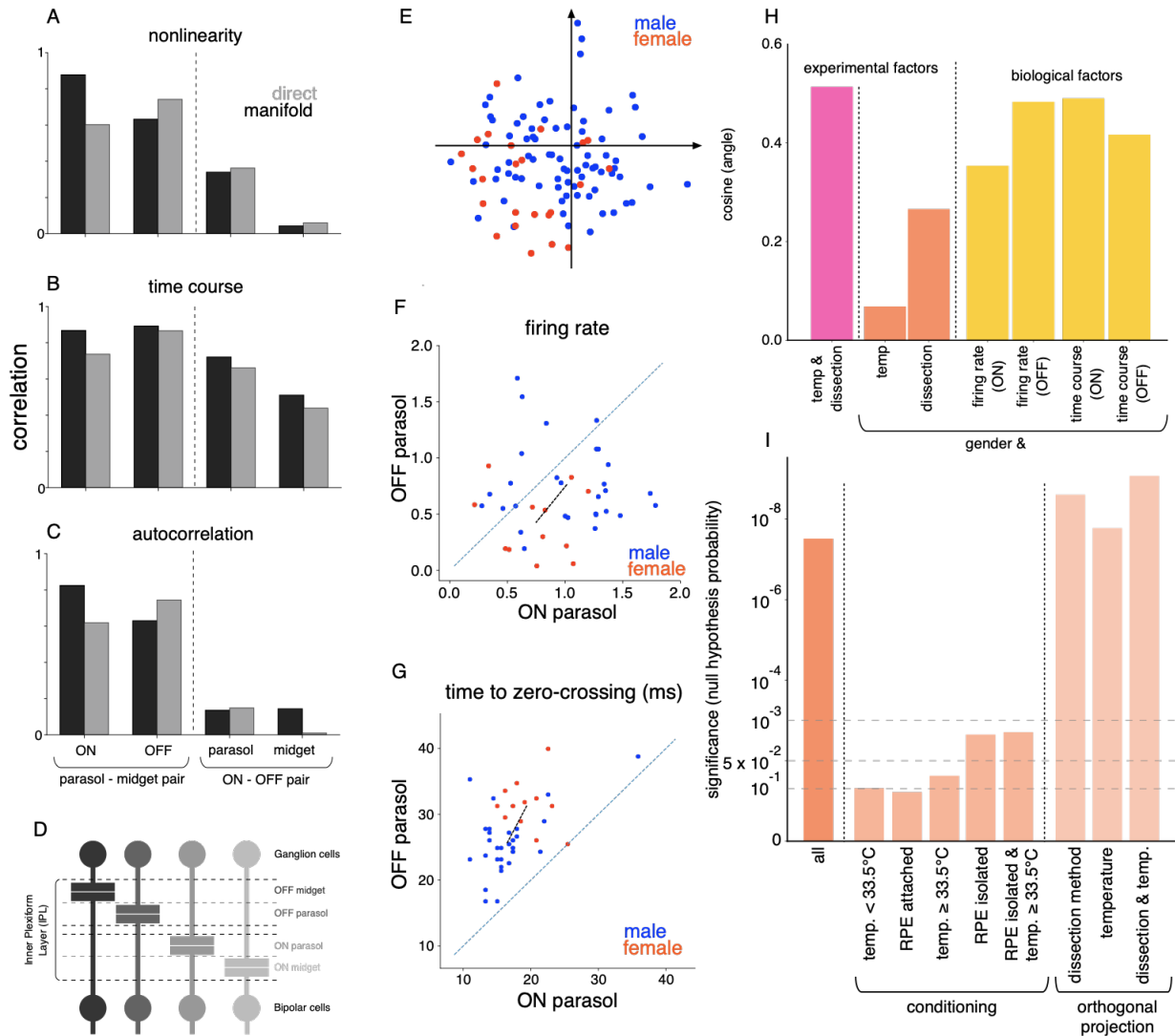
39 Second, the recordings from male and female retinas were separated in the manifold ($d'=1.8$ for
40 15 dimensional manifold, Figure 3E), in a way that could not be explained by variations in
41 experimental factors (Figure 3I). For both ON and OFF parasol cells, the differences between
42 male and female retinas were partially due to differences in firing rate and speed of temporal
43 filtering. This was determined by computing the direction separating the means of male and
44 female retinas in the manifold, and determining the angle between this direction and the
45 directions most aligned with variation in firing rate and response time course (cosine(angle) \sim
46 0.5 for both) (Figure 3H). Note that unlike gender, animal age did not separate different

1 recordings on the manifold. In principle, the observed gender differences could be confounded
2 by variation in experimental methods such as dissection procedure (isolation from the retinal
3 pigment epithelium, or RPE) or temperature (31°-36°C across recordings). Because higher
4 recording temperatures were associated with isolation from the RPE for technical reasons, the
5 directions in the manifold representing dissection method and temperature variation were
6 aligned (cosine(angle) ~ 0.57, Figure 3H). Compared to firing rate and response time course,
7 these experimental factors were less aligned to the direction of gender separation
8 (cosine(angle) ~ 0.26) (Figure 3H), suggesting that gender differences were probably not due to
9 differences in these experimental methods. To eliminate experimental factors more rigorously,
10 the separation of males and females in the manifold was measured after conditioning the data in
11 several ways. For each condition, a bootstrap rank test was performed to verify if the mean
12 locations of male and female recordings differed (see Methods). Significant separation ($p < 0.05$)
13 was observed for the recordings with RPE-isolated dissections and high temperature ($\geq 33.5^\circ\text{C}$),
14 whereas the separation was not significant ($p > 0.1$) for RPE-attached dissections and lower
15 temperatures ($< 33.5^\circ\text{C}$) (Figure 3I). For the RPE-isolated retinas, the male recordings exhibited
16 higher firing rates (Figure 3F) and faster temporal integration (Figure 3G) ($p < 0.01$, see Methods
17 for details).

18 Although this conditioning on specific experimental conditions revealed statistically significant
19 differences between males and females, the level of significance was lower when compared to
20 all the recordings (Figure 3I), potentially due to a reduction in the number of samples when
21 analysis is restricted to a particular set of experimental conditions. The manifold made it
22 possible to separate experimental variations more efficiently, without reducing the number of
23 data points. To accomplish this, the data were projected onto axes in the manifold orthogonal to
24 the two identified directions of experimental variability. This projection increased the statistical
25 separation between the male and female retinas ($p < 10^{-6}$) (Figure 3I). Thus, the geometry of the
26 manifold, which smoothly organizes the neural computation across retinas, makes it possible to
27 examine statistical trends in the data efficiently in spite of potential experimental confounds.

28

29



1
 2 **Figure 3 Biological factors underlying variability.** Relation between the first principal component of (A)
 3 nonlinearity, (B) response time course and (C) spike train autocorrelation variation for different pairs of cell types. The
 4 relationship is either measured directly using Spearman rank correlation, or the cosine of the angle between
 5 corresponding directions in the 15-dimensional manifold. (D) Distinct lamination depths for the bipolar-ganglion cell
 6 synapse for different ganglion cell types¹¹. (E) Two dimensional PCA projection of manifold locations for recordings
 7 from male (blue) and female (red) retinas. (F) The average firing rate for ON parasol (y-axis) and OFF parasol (x-
 8 axis) cells for preparations with isolated RPE. The mean manifold location of male (blue) and female (red) recordings
 9 were different ($p < 0.01$ for bootstrap and $p < 0.05$ for hierarchical bootstrap). Black line joins the mean male and female
 10 locations. (G) Similar to (F) for the time course of STA, with separation of male and female retinas ($p < 0.01$ for
 11 bootstrap and $p < 0.05$ for hierarchical bootstrap¹²). (H) Cosine of the angle (y-axis) between the manifold directions
 12 corresponding to different pairs of factors, which are either biological (gender, firing rate, time of zero crossing of time
 13 course) and experimental (temperature, dissection - whether retinal pigment epithelium (RPE) was attached or
 14 isolated). (I) Degree of separation of male and female recordings measured using a resampling test, for all the
 15 recordings, conditioned on the subset of recordings with specific dissection procedure or temperature, or all
 16 recordings with locations projected orthogonal to directions for dissection procedure and temperature variation.

17

1 **Manifold generalizes to a novel recording**

2 The manifold permitted efficient response modeling of a new, previously unseen retina by
3 leveraging trends in the large data set of retinas used for training. Response modeling can be
4 performed efficiently by identifying the manifold location of a new retina in several ways, using
5 limited data. In the absence of any new data, the “typical” manifold location for a new retina can
6 be obtained by merely *averaging* the locations of all training retinas. In the case of a
7 degenerated retina with no light evoked response, partial information such as the
8 autocorrelation function of the spiking of recorded neurons can still be identified from
9 spontaneous activity. In this case, the manifold location can be *approximated* by averaging the
10 locations of training retinas that have similar autocorrelation. Finally, in the presence of light
11 evoked responses, gradient descent can be used to *optimize* the manifold location based on the
12 likelihood of the data, leveraging the training retinas as prior information (Figure 4A).

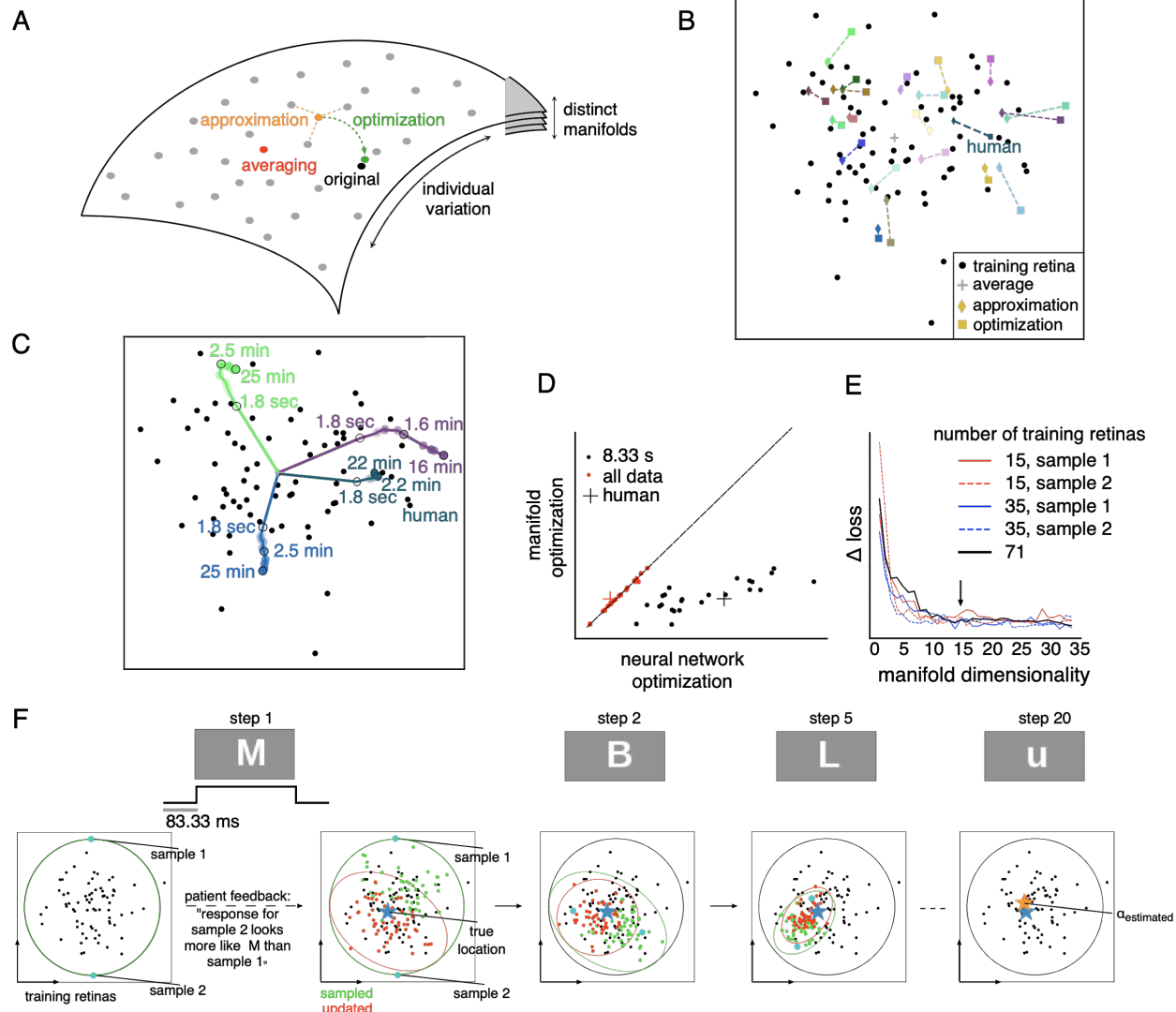
13 These three approaches were examined using a model trained with 71 preparations and tested
14 with 24 preparations. The proximity of manifold locations identified by approximation and
15 optimization suggest that these approaches accurately capture properties of the new retina
16 (Figure 4B). To examine the efficiency obtained by using the manifold, the optimization
17 approach with limited data was studied by measuring convergence of the manifold location and
18 of response prediction accuracy. First, the manifold location converged quickly as the recording
19 duration increased, with ~3 minutes of light response data producing a location similar to that
20 produced by ~30 minutes of data (Figure 4C). Second, whereas optimization of the manifold
21 location performed similarly to training the full model (along with the CNN) when performed with
22 a large amount of data, the manifold approach more accurately predicted light responses when
23 the data were limited (~8 sec) (Figure 4D).

24 The efficiency of the low dimensional manifold could come at a cost of its generalization
25 accuracy to new retinas. On varying the manifold dimensionality, the ability to predict responses
26 to previously unseen retinas saturated at ~15 dimensions (Figure 4E, black line), a value that
27 did not change with the number of retinas used for training (Figure 4E, colored lines). Hence, in
28 addition to the previous observation on generalization to new stimuli within the collection of
29 retinas used for learning (Figure 1D), the low dimensional manifold is also able to generalize to
30 new, previously unseen retinas.

31 These findings suggest that the manifold may aid in translating our understanding of the
32 macaque retina to the human retina, an important goal for biomedical research. Recent work^{13,14}
33 has shown that the receptive field properties of the four numerically dominant RGC types (ON
34 and OFF parasol and midget) are similar to those of their macaque counterparts. To test
35 whether light responses in the human retina fall within the range observed across many
36 macaque retinas, the manifold location of a single human retina was identified using the three
37 operations described above (averaging, approximation, optimization). For each method, the
38 estimated manifold location of the human retina was well within the span of manifold locations of
39 many macaque retinas (Figure 4B,C), and the responses of a new retina were predicted with
40 similar accuracy and efficiency in the two species (Figure 4D). Hence, the manifold reveals
41 similar light response properties of corresponding RGC types in macaques and humans.

1 Given that the macaque retinal code translates accurately to humans, it may provide a valuable
2 tool for the development of an advanced artificial retina for treating vision loss. However, a
3 challenging first step in restoring vision with such a device is to identify how the neurons in a
4 blind retina should encode visual stimuli with the implant. In this setting, the retina is no longer
5 light-sensitive, but the blind person with an implanted device could potentially report the
6 similarity of artificially induced images to a verbally described object. The correct neural
7 encoding could be identified by estimating the location of the retina on the low-dimensional
8 manifold, using a psychophysical discrimination task with the following iterative procedure: (i)
9 sample a few of the plausible manifold locations, (ii) use each of these manifold locations to
10 predict retinal responses for a particular visual stimulus, (iii) stimulate with the implanted device
11 to produce each of these responses, (iv) ask the subject which stimulus produced a sensation
12 that most closely matches a verbal description of it, and (v) update the set of plausible manifold
13 locations.

14 To illustrate the feasibility of this procedure, the above steps were simulated assuming (for
15 simplicity) an artificial retina that has perfect cellular selectivity of stimulation. For updating the
16 plausible locations, the perceptual accuracy was assumed to be governed by the Kullbeck-
17 Leibler divergence between the distribution of neural responses associated with the tested
18 manifold location and the distribution of responses associated with the true location. The set of
19 plausible locations identified by the procedure converged to the true location (given by
20 optimization; see above) in a small number of steps (see Methods, Figure 4F). Hence, the low-
21 dimensional manifold could provide valuable efficiency in translational applications.



1
2 **Figure 4. Generalization to new, previously unseen retinas.** (A) Efficient identification of manifold location.
3 Training retinas (gray points) are used to define the manifold. Then, the manifold location of a new retina is
4 determined in one of several ways - averaging (red): compute the mean location of all training retinas; approximation
5 (orange): compute the mean location of a subset of training retinas with specific features similar to the new retina;
6 optimization (green): gradient descent on the manifold location to maximize the prediction accuracy for measured
7 light responses of the new retina. (B) Identified manifold locations using averaging (+; same for all retinas),
8 approximation (\diamond) and optimization (\square) for testing retinas (colors, each pair joined with a line). Black points
9 correspond to the locations of training retinas. (C) Optimized manifold location for three retinas (colored lines), with
10 varying duration of recorded responses. Training retina locations (black points) and locations identified by averaging
11 (+) are shown. For (B, C), the 15 dimensional manifold is projected into two dimensions that capture autocorrelation
12 variation (same as Figure 2E) (D) Response prediction loss with optimization of manifold location (y-axis) vs. loss
13 with learning all neural network parameters (x-axis), using either 8.33ms (red) or 15-30 min (black) of data. (E)
14 Convergence of prediction loss with optimized manifold location, averaged across 24 testing retinas (y-axis) as a
15 function of the number of manifold dimensions (x-axis). Loss measured as negative log-likelihood, averaged across
16 cells. Colored lines indicate loss obtained using fewer training retinas. (F) Simulation of the discrimination task
17 used to identify manifold location in the retina of a blind person. At each step, the visual stimulus is a letter from
18 English alphabet (top row). Two dimensional projection of the manifold with locations of training retinas (black dots) is
19 shown (bottom row). Posterior over the set of feasible manifold locations approximated using a Gaussian (red circle). At
20 each step, random manifold locations are sampled from the posterior, and corresponding responses are reproduced

1 using the artificial retina, and feedback from the subject is used to update the posterior using Monte Carlo methods
2 (all samples in green, accepted samples in red). In ~20 steps, the estimated manifold location (orange) converges to
3 the true underlying location (blue).

4 **Conclusions and Discussion**

5 Recent studies have shown that neural responses in the retina, which have often been
6 described using pseudo-linear models^{6,15}, can be modeled more accurately using artificial neural
7 networks^{9,10}. However, the function of these complex models can be difficult to interpret. While
8 using machine learning models to more accurately capture shared (and complex) aspects of
9 neural response, the present work crucially maintained the information about variability between
10 recordings in an interpretable low-dimensional manifold.¹⁶ As a result, the manifold revealed
11 novel structure in inter-retina variation, such as gender-based differences and covariation of
12 several response properties across cell types. In principle, such an approach could be used for
13 other applications in which the varying neural computation of interest is captured using a simple
14 low-dimensional manifold, while machine learning models improve accuracy by capturing
15 complex but unchanging aspects of the computation that are of less immediate interest.

16 At least three sources of biological variability in the neural code may play a role in the present
17 data: variation between animals, differences between the two eyes, and variation across retinal
18 locations in a given eye. However, in the present data, these sources of variability are
19 confounded with variation in experimental procedures for euthanasia, for eye removal, and for
20 *ex vivo* recording, respectively. Thus, it is not possible to definitively isolate biological and
21 experimental variation in the present data, though the observed male-female differences were
22 separated from the clear influences of experimental procedures and therefore most likely reflect
23 true biological differences. The tools developed here make it possible to capture and analyze
24 both biological and experimental variability in a single framework, an asset for understanding
25 the neural code more completely using imperfect experimental data. However, future work will
26 be needed to parse the mechanisms of both biological and experimental variability, potentially
27 exploiting the geometry of the manifold. For example, the low-dimensional manifold made it
28 possible to project the data into a subspace orthogonal to known confounds, and thus to control
29 for them without losing data by conditioning on specific experimental variables.

30 The present findings on neural coding differences between male and female retinas add to a
31 large literature on gender-based differences in brain structure and function¹⁷⁻¹⁹. In the retina,
32 genetic differences between male and female primates (including humans) produce different
33 variation in cone photopigment spectral sensitivities, and thus different color vision, across the
34 population²⁰. However, to our knowledge, differences in neural coding between males and
35 females have not been reported, perhaps due to the lack of appropriate physiological recordings
36 and/or analysis tools. In the present work, the large data set and the geometry of the manifold
37 representation made it possible to establish the observed male-female differences in spite of
38 other potential confounds.

39 The manifold of neural coding variability may also be useful in other neuroengineering
40 applications, such as motor prostheses. The goal in motor prosthesis is usually to read out the
41 neural activity in a paralyzed person to control a computer cursor or a robotic limb²¹. Similar to

1 the problem of identifying the neural encoding in blind person, identifying the neural mapping in
2 a paralyzed person is limited by the absence of simultaneous neural recordings and limb
3 trajectory measurements, and thus may benefit from leveraging existing data that captures the
4 diversity of neural coding across individuals. Specifically, a manifold of inter-individual variation
5 may be useful for identifying neural decoding in a person with an implant, perhaps using a task
6 involving imagined movements to identify the manifold location²². The manifold may also be
7 useful for adjusting to the challenge of variability over time in chronic recordings²³.

8 The neural coding manifold may also be useful for harnessing brain plasticity, which could
9 improve vision with an artificial retina. Indeed, present-day retinal implants make little attempt to
10 reproduce the neural code of the retina, and thus implicitly rely heavily on plasticity to
11 compensate for device limitations²⁴. In motor prostheses, it has been shown that the brain can
12 more easily adjust its activity to accommodate perturbations in the artificial neural decoder if
13 these perturbations lie in a specific low-dimensional manifold²⁵. In the case of an artificial retina,
14 we hypothesize that the brain may more readily learn to interpret the neural activity produced by
15 the implanted device if the visual encoding that it uses lies within the manifold of retinal coding
16 variability.

17 **Acknowledgements**

18 We thank J. Carmena, K. Bankiewicz, T. Moore, W. Newsome, M. Taffe, T. Albright, E. Callaway, H. Fox,
19 R. Krauzlis, S. Morairty, and the California National Primate Research Center for access to macaque
20 retinas. We thank Kristy Berry, K. Williams, B. Morsey, J. Frohlich, and M. Kitano for accumulating and
21 providing macaque retina metadata. The human eye was provided by Donor Network West (San Ramon,
22 CA). We are thankful for the cooperation of Donor Network West and all of the organ and tissue donors
23 and their families, for giving the gift of life and the gift of knowledge, by their generous donations. We
24 thank E. Simoncelli, S. Mitra, V. Gupta, K. Talwar, V. Feldman, A. Gogliettino, E. Wu, S. Madugula and
25 the entire Stanford Artificial Retina team for helpful discussions. We thank Google internship and Student
26 Researcher programs (NPS), Research to Prevent Blindness Stein Innovation Award, Wu Tsai
27 Neurosciences Institute Big Ideas, NIH NEI R01-EY021271, NIH NEI R01-EY029247, and NIH NEI P30-
28 EY019005(EJC), NSF IGERT Grant 0801700 (NB) and NIH NEI F31EY027166 (CR), NSF GRFP DGE-
29 114747 (NB and CR) for funding this work.

30 **Methods**

31 **Recordings**

32 Preparation and recording methods are described elsewhere^{1,2,8}. Briefly, eyes were enucleated from
33 terminally anesthetized macaque monkeys (*M. Mulatta* or *M. Fascicularis*) used by other experimenters in
34 accordance with institutional guidelines for the care and use of animals. Immediately after enucleation,
35 the anterior portion of the eye and the vitreous were removed in room light. The eye was stored in
36 darkness in oxygenated Ames' solution (Sigma, St. Louis, MO) at 33°C pH 7.4. Segments of isolated or
37 RPE-attached peripheral retina (approximately 3mm x 3mm, taken from 6-15mm temporal equivalent
38 eccentricity⁸) were placed flat, RGC side down, on a planar array of 512 extracellular microelectrodes
39 arranged in an isosceles triangular lattice. The electrode spacing was 60µm in each row, and the array
40 covered a rectangular region measuring 1800 µm x 900 µm. While recording, the retina was perfused with
41 Ames' solution (31-36°C; typically 32 °C for RPE attached and 34°C for RPE isolated dissections),

1 bubbled with 95% O₂ and 5% CO₂, pH 7.4. Voltage signals on each electrode were bandpass filtered
2 (80Hz - 2kHz), amplified, and digitized at 20 kHz¹.

3 A custom spike-sorting algorithm was used to identify and segregate spikes from distinct cells¹. Briefly,
4 candidate spike events were detected using a threshold on each electrode, and voltage waveforms on the
5 electrode and nearby electrodes in the 4ms period surrounding the time of the spike were extracted.
6 Candidate neurons were identified by clustering the waveforms using a Gaussian mixture model.
7 Candidate neurons were retained only if the assigned spikes exhibited a 1 ms refractory period and had a
8 stable firing rate for the entire duration of recording. Duplicate spike trains were identified by temporal
9 cross-correlation and removed. For each cell, the autocorrelation of the recorded spike train was
10 computed and normalized by its value at zero time lag.

11 **Visual stimuli and cell type identification**

12 Visual stimuli were delivered using the optically reduced image of a CRT monitor refreshing at 120 Hz
13 and focused on the photoreceptor outer segments. The optical path passed through the transparent plug
14 and Ames' solution or through the mostly transparent electrode array and the retina. The relative
15 emission spectrum of each display primary was measured with a spectroradiometer (PR-701,
16 PhotoResearch) after passing through the optical elements between the display and the retina. The total
17 power of each display primary was measured with a calibrated photodiode (UDT Instruments). The mean
18 photoisomerization rates for the cone photoreceptors were estimated by computing the inner product of
19 the primary power spectra with the spectral sensitivity of each cone type, and multiplying by the effective
20 collecting area of primate cones ($\sim 0.6 \mu\text{m}^2$)^{26,27}, resulting in photoisomerization rates of approximately
21 800–2200, 800–2200, 400–900 for the long-, middle- and short-wavelength sensitive cones, respectively.
22 The stimulus pixel size on the retina was either 41.6 microns (8 monitor pixels), 52 microns (10 monitor
23 pixels) or 83.2 microns (16 monitor pixels). A new white noise frame was drawn at refresh rates of 60 Hz
24 or 30 Hz. The pixel contrast (difference between the maximum and minimum intensities divided by the
25 sum) was 96% for each display primary, with mean intensity of 50%. The white noise stimulus either
26 modulated the three display primaries independently, or coherently, at each spatial location.

27 In each recording, RGCs were classified into distinct types using properties of the spatial receptive field
28 and response time course obtained from the spike-triggered average (STA) stimulus^{68,28}. A two-
29 dimensional Gaussian fit to the spatial receptive field was used for determining the center location⁸. All
30 analyses used recordings with stable firing rates and nearly complete tiling of ON and OFF parasol cell
31 receptive field mosaics. For Figure 3, only recordings that also had nearly complete ON and OFF midget
32 cell mosaics were used. For model fitting, both the visual stimulus and spike times were binned at 8.33ms
33 (120Hz), and the visual stimulus was upsampled to 8 monitor pixels, resulting in a common 80x40 pixel
34 grid across recordings.

35 **Linear Nonlinear Poisson model**

36 The Linear Nonlinear Poisson (LNP) model consists of a linear spatio-temporal filter followed by a point
37 nonlinearity⁶. A filter that is separable in space and time was used⁸, which is equivalent to a cascade of a
38 spatial filter and a temporal filter. These filters were estimated for each cell in each recording as follows.
39 First, the STA was computed by averaging the stimulus preceding spikes, over all pixel locations and
40 250ms (30 frames at 120Hz) prior to the spike. Next, the spatial filter was computed by choosing the STA
41 frame with the single largest pixel magnitude. The spatial filter was restricted to a rectangular window
42 around the receptive field. The receptive field was defined as the set of pixels with absolute magnitude
43 greater than 2.5σ , contiguous with the strongest pixel, where σ is the robust standard deviation²⁹ of pixels

1 in the STA, an estimate of the measurement noise. Next, the temporal filter was identified by averaging
2 the time course of all pixels in the receptive field. Finally, the output nonlinearity was estimated by fitting a
3 5th order polynomial to the relationship between the observed responses and the generator signal, which
4 was computed by filtering the stimulus with the estimated spatial and temporal filters⁶.

5 For describing variation in spatial nonlinearities, subunit models were fitted using spike triggered
6 clustering d⁷. Briefly, neural responses were modeled by passing the stimulus through multiple linear
7 filters (subunits), followed by an exponential nonlinearity, and summation over the filter outputs. The
8 temporal filtering of each subunit was assumed to be identical to the time course of the STA. The spatial
9 filters were inferred by soft-clustering the collection of stimuli preceding a spike. To assess the degree of
10 subunit nonlinearity, the log-likelihood of a model with five subunits was compared to that of a model with
11 one subunit (which reduces to an LNP model).

12 **Neural network model**

13 A convolutional neural network was used to predict RGC responses across multiple recordings
14 simultaneously. Below, the model architecture and the fitting procedure are described in detail.

15 *Notation*

16 The model $f(S, \alpha_i, C_i)$, takes as its input the visual stimulus S , recording-specific information about the
17 collection of recorded cells C_i , and the recording-specific manifold location α_i , and yields as its output the
18 predictions for recording-specific response R_i .

19 The recent history of the visual stimulus is given by $S \in R^{d_x \times d_y \times d_z}$, where R is the set of real numbers,
20 $d_x \times d_y$ are the spatial dimensions (80 x 40), and d_z is the number of time bins (30). Stimuli presented at
21 different spatial or temporal resolution were upsampled or downsampled to these dimensions.

22 The recording specific manifold location is given by $\alpha_i \in R^n$, where n is the manifold dimensionality.

23 The recording specific information about recorded cells is given by $C_i = \{x(c), y(c), t(c)\}_{c=1}^{|C_i|}$, where
24 each cell c is described by its receptive location $(x(c), y(c))$ in the $d_x \times d_y$ visual space and its cell type
25 $t(c)$. For models with only two cell types, $t(c) \in \{0, 1\}$, for ON and OFF parasols respectively. For models
26 with four cell types, $t(c) \in \{0, 1, 2, 3\}$, corresponding to ON parasol, OFF parasol, ON midget, and OFF
27 midget cell types.

28 The responses are given by $R_i \in Z_+^{|C_i|}$, where Z_+ denotes non-negative integers and $|C_i|$ is the
29 collection of cells in recording i . Responses were binned at the same resolution as the stimulus (8.33ms).

30 *Model architecture*

31 The model $f(S, \alpha_i, C_i)$ passes the visual stimulus S through a multilayered convolutional neural network,
32 with each layer consisting of a convolution (stride 1), retina-specific normalization and a point-wise
33 (softplus) nonlinearity (see Figure 1B). The model output is Poisson firing rate. This firing rate is used to
34 predict the responses R_i . The number of channels and filter sizes are chosen by cross-validation, as
35 described below. Recording-specific normalization and challenges associated with predicting responses
36 for varying numbers of cells across recordings are also given below.

1 Recording-specific normalization is inspired by previous work³⁰, in which a translationally-invariant affine
2 transformation of the layer activations adapts the model to each recording. The scale and shift
3 coefficients for this affine transform are determined linearly from the manifold location α_i . Let $\hat{a}(x, y, l, t)$ be
4 the activation after convolution at location x, y in the channel l of layer t . First, the mean μ and standard
5 deviation σ across samples in a batch are computed, and used to calculate normalized activations:
6 $\tilde{a}(x, y, l, t) = \frac{\hat{a}(x, y, l, t) - \mu}{\sigma}$. Next, using the manifold location α_i , a learned affine transform determines the
7 desired mean ($\tilde{\mu} = P\alpha_i + p$) and standard deviation ($\tilde{\sigma} = Q\alpha_i + q$) for each layer. Finally, the
8 normalized activations are transformed to give recording-specific activations $a(x, y, l, t) \leftarrow \tilde{a}(x, y, l, t)\tilde{\sigma}_{l,t} +$
9 $\tilde{\mu}_{l,t}$. Note that the retina-specific scales and shifts are the same for each location in visual space,
10 preserving the translational invariance of convolutional networks and reflecting the homogenous response
11 properties of the RGCs belonging to a single type.

12 A potential challenge is that the number of recorded neurons, and hence the number of outputs of $f(\cdot)$, is
13 variable across recordings. To address this issue, the model predicts multiple response maps, one for
14 each cell-type, with the same spatial dimensions as the visual stimulus. The response for each cell is
15 read off from its cell location in the response map of the corresponding cell type. Specifically, $f(\cdot)$ outputs
16 $m_i(x, y)$, which corresponds to firing rate map of cell-type i , and for a cell with type $t(c)$ and centered at
17 $x(c), y(c)$, the Poisson firing rate is given by $m_{t(c)}(x(c), y(c))$.

18 *Model fitting*

19 Estimation of recording-specific parameters (α_i) and the shared parameters are performed by maximizing
20 the log-likelihood of observed responses, summed across all the cells, recordings and stimuli. This is
21 performed by stochastic gradient descent, where at each step, a randomly sampled batch of stimuli and
22 corresponding responses from a particular recording are used to update the shared and the
23 corresponding recording-specific parameters. The batch size was 250 and the updates were performed
24 using the Adam³¹ update algorithm with learning rate of 0.1. For each training retina, the first 4 min of
25 white noise data were used for testing and the remainder was used for training. The duration of the
26 stimulus varied from 15-90 min (median 30 min) across experiments. A model with 4 layers, 3 x 3 or 1 x 1
27 filter size, 64 channels per layer and a 15 dimensional manifold was chosen based on cross validation
28 and used for subsequent analysis (see Figure 1B for architecture).

29 **Variation of neural coding on the manifold**

30 The following steps were used to test if the manifold captured variations in neural response properties
31 across recordings. First, the manifold direction that was maximally correlated to the variations of a
32 particular response feature was identified by linear regression. Second, recordings were projected along
33 this direction, and the Spearman rank correlation with the response property was measured. Statistical
34 significance was measured with a permutation test, where the null distribution was generated by
35 permuting the recordings with the manifold locations fixed. In Figure 2, the Spearman rank correlation and
36 its statistical significance was measured for the first principal component projection of various response
37 properties such as receptive field size (ON: 0.78, $p < 0.0001$; OFF: 0.76, $p < 0.0001$), time course (ON:
38 0.83, OFF: 0.85; $p < 0.0001$), output non-linearity of the LNP model (Spearman rank correlation for ON:
39 0.92; OFF: 0.92; $p < 0.0001$) and normalized auto-correlation (ON: 0.97, OFF: 0.94, $p < 0.0001$). The
40 interdependence between response properties was either measured directly in raw data using Spearman
41 rank correlation, or with the angle between the corresponding manifold directions.

42 **Differences between male and female recordings**

1 For analysis of gender differences, only the subset of recordings (102) from one species (M. Mulatta)
2 were used. First, the separation between male and female recordings was measured by computing the d'
3 value of the projection of the two distributions onto the difference in the means. The d' value observed
4 (~ 1.8) indicated that the gender based differences were not large on an individual basis. Second, a
5 bootstrap test was performed to test whether the mean locations of the male and female recordings were
6 statistically distinguishable. The distance between the mean manifold locations of male and female
7 retinas was measured and compared to a null distribution of distances generated by resampling (with
8 replacement) of the manifold locations. The null distribution was fitted with a normal distribution and the
9 significance level was measured as the probability mass greater than the observed distance in data.
10 Because multiple preparations were frequently recorded from the same animal, a hierarchical variant of
11 this bootstrap test was performed, in which the resampling was performed according to the hierarchical
12 structure of the data¹², by first sampling an animal and then sampling the manifold location of one of the
13 recordings from that animal, both with replacement. Hierarchical bootstrap is more conservative and
14 biased towards accepting the null hypothesis¹². Mean manifold locations for male and female retinas were
15 significantly different ($p < 10^{-6}$ for bootstrap and $p < 0.05$ for hierarchical bootstrap). Identical tests were
16 applied for assessing male-female differences in firing rate and the speed of temporal filtering ($p < 0.01$ for
17 bootstrap and $p < 0.05$ for hierarchical bootstrap for both quantities).

18 **Invariance of neural coding on the manifold**

19 The ability of the manifold to preserve previously reported invariances of the neural code was tested as
20 follows. First, random manifold locations were sampled by perturbing the learned locations of training
21 retinas with a Gaussian noise of standard deviation equal to their median nearest-neighbor distance.
22 Second, a ~ 800 sec long white noise stimulus was sampled, and ON and OFF parasol firing rate maps
23 were computed using the neural network, which was adapted using the manifold locations. Third, the
24 Poisson firing rates for 200 cells (100 of each type) with random receptive field locations were read off
25 from the firing rate maps. Finally, the cell responses were sampled, and used to estimate a Linear-
26 Nonlinear Poisson model, which served as an interpretable summary of neural encoding captured by the
27 manifold location. Comparison of average receptive field size and the zero crossing time of the temporal
28 filter revealed known invariances between ON and OFF parasols (Figure 2J, K).

29 **Estimation of the manifold location of a previously unseen retina**

30 By fixing the shared parameters after learning, and estimating the recording-specific representation on
31 the manifold, the trained model was adapted to predict responses in a new, previously unseen recording.
32 Based on the amount of data available, several methods can be employed to identify the manifold
33 location (Figure 4). These methods are described below in detail.

34 *Averaging:* When no data about the new retina are available, the simplest approach is to average the
35 locations of all the retinas used for training.

36 *Approximation:* This is similar to averaging, but only using the subset of training retinas with similar
37 response properties as the new retina. For Figure 4, locations of five training retinas with the most similar
38 autocorrelation function were used for approximation.

39 *Optimization:* When light response data are available, the manifold location α_i was determined by
40 Bayesian inference. Bayesian inference combines a Gaussian prior ($P_{prior}(\alpha) \sim N(\mu_{prior}, \sigma_{prior})$) over
41 manifold locations determined from the training retinas and the likelihood ($P(R_i^t | S^t, \alpha_i)$) of stimulus-

1 response data for the new retina. The posterior $P(\alpha_i | \{S^t, R_i^t\})$ was maximized using gradient ascent
 2 (learning rate 0.1) :

$$3 \quad \alpha_i^* = \operatorname{argmax}_{\alpha} P(\alpha_i | \{S^t, R_i^t\}) = \operatorname{argmax}_{\alpha} \sum_t \log P(R_i^t | S^t, \alpha_i) + \log P_{\text{prior}}(\alpha_i)$$

4 *Discrimination task*: Simulation of manifold location estimation in a blind person implanted with a retinal
 5 prosthesis was performed using a discrimination task. For a given visual stimulus, the task involves using
 6 the implanted retinal prosthesis to reproduce responses corresponding to two manifold locations and the
 7 subject selects the response that yields perception most closely matching a verbally described stimulus.
 8 Multiple rounds of this task are used to update the posterior on manifold locations.

9 The discrimination task was simulated under the assumption that the perceptual difference of the
 10 responses generated by hypothetical retinas at two manifold locations α_1 and α_2 for a stimulus S is equal
 11 to the KL-divergence between the corresponding response distributions $P(R|S, \alpha_1)$ and $P(R|S, \alpha_2)$. Given
 12 α_{true} as the true underlying manifold location, the blind person's feedback $Y(\alpha_{\text{true}}, \alpha_1, \alpha_2) = 0$ if

$$13 \quad D_{KL}(P(R|S, \alpha_1) || P(R|S, \alpha_{\text{true}})) \leq D_{KL}(P(R|S, \alpha_2) || P(R|S, \alpha_{\text{true}}))$$

14 and $Y = 1$ otherwise. For simplicity, sampled responses were used to compute an unbiased estimate of
 15 the KL-divergence :

$$16 \quad D_{KL}(P(R|S, \alpha_1) || P(R|S, \alpha_{\text{true}})) \approx \sum_{R_i \sim P(R|S, \alpha_1)} \log \left(\frac{P(R_i|S, \alpha_1)}{P(R_i|S, \alpha_{\text{true}})} \right).$$

17 Hence, the posterior over manifold location after t steps of the task is given by:

$$18 \quad P_{\text{posterior}}(\alpha | \{S_t, \alpha_{1,t}, \alpha_{2,t}, Y_t\}_{t=1}^{t=T}) \propto \prod_{t=1}^{t=T} P(Y_t | S_t, \alpha_{1,t}, \alpha_{2,t}) P_{\text{prior}}(\alpha)$$

19 where the prior is estimated from training retinas as a Gaussian distribution: $P_{\text{prior}}(\alpha) \sim N(\mu_{\text{prior}}, \sigma_{\text{prior}})$.

20 In the simulations, the visual stimulus S consisted of letters of English alphabet, flashed for 100 ms and
 21 preceded and succeeded by 50 ms of gray screen. At each step, a Gaussian approximation of the
 22 posterior $P_{\text{posterior}}(\alpha) \sim N(\mu_{\text{posterior}}, \sigma_{\text{posterior}})$ was maintained, and updated using Monte-Carlo
 23 sampling. In summary, the steps for the t^{th} iteration of the algorithm are as follows :

- 24 1. Sample symmetric $\alpha_{1,t}, \alpha_{2,t}$ around posterior mean: $\alpha_{1,t} \sim P_{\text{posterior}}(\alpha)$; $\alpha_{2,t} = 2\mu_{\text{posterior}} - \alpha_{1,t}$.
- 25 2. Sample an English letter and a target stimulus S_t .
- 26 3. Sample responses $R_{1,t} \sim P(R|S_t, \alpha_{1,t})$; $R_{2,t} \sim P(R|S_t, \alpha_{2,t})$.
- 27 4. Get patient feedback $Y(\alpha_{\text{true}}, \alpha_{1,t}, \alpha_{2,t})$, based on an estimate of the KL divergence using sampled
 28 responses $R_{1,t}, R_{2,t}$.
- 29 5. Update the posterior of plausible manifold locations.
 - 30 a. Sample N retina locations $\alpha_i \sim P_{\text{posterior}}(\alpha)$ for $i \in [1, \dots, N]$.
 - 31 b. For the set of sampled manifold locations, find the subset that matches user feedback,
 32 i.e., with $Y(\alpha_i, \alpha_{1,l}, \alpha_{2,l}) = Y(\alpha_{\text{true}}, \alpha_{1,l}, \alpha_{2,l})$ for all $l (= 1, \dots, t)$ previous steps. Let this
 33 subset of be $\{\tilde{\alpha}_j\}$.

- 1 c. Update posterior distribution with $\mu_{posterior} = \langle \tilde{\alpha}_j \rangle$ and $\sigma_{posterior}^2 = \langle \alpha_j \alpha_j^T \rangle - \mu_{posterior} \mu_{posterior}^T$.

3 For results shown in Figure 4F, α_{true} was set as the result of optimizing the manifold location using light-
4 evoked responses. In the simulations, the posterior distribution converged in ~20 steps, suggesting that
5 the low dimensional manifold can be used for efficiently identifying the expected neural code in a blind
6 person. However, the amount of noise in the simulation is probably lower compared to what would be
7 encountered in practice, leading to a larger number of steps to identify the true manifold location and may
8 perhaps require changes to the estimator of KL divergence and the method to update the posterior of α .

9 **Data/Code Availability**

10 The data/code that support the findings of this study are available from the corresponding author upon
11 reasonable request.

12 **References**

- 13 1. Litke, A. M. *et al.* What does the eye tell the brain?: Development of a system for the large-
14 scale recording of retinal output activity. *IEEE Trans. Nucl. Sci.* **51**, 1434–1440 (2004).
- 15 2. Frechette, E. S. *et al.* Fidelity of the ensemble code for visual motion in primate retina. *J.*
16 *Neurophysiol.* **94**, 119–135 (2005).
- 17 3. Field, G. D. *et al.* Spatial properties and functional organization of small bistratified ganglion
18 cells in primate retina. *J. Neurosci.* **27**, 13261–13272 (2007).
- 19 4. Rhoades, C. E. *et al.* Unusual Physiological Properties of Smooth Monostratified Ganglion
20 Cell Types in Primate Retina. *Neuron* **103**, 658–672.e6 (2019).
- 21 5. Greschner, M. *et al.* A polyaxonal amacrine cell population in the primate retina. *J.*
22 *Neurosci.* **34**, 3597–3606 (2014).
- 23 6. Chichilnisky, E. J. A simple white noise analysis of neuronal light responses. *Network* **12**,
24 199–213 (2001).
- 25 7. Shah, N. P. *et al.* Inference of nonlinear receptive field subunits with spike-triggered
26 clustering. (2020) doi:10.7554/eLife.45743.
- 27 8. Chichilnisky, E. J. & Kalmar, R. S. Functional asymmetries in ON and OFF ganglion cells of
28 primate retina. *J. Neurosci.* **22**, 2737–2747 (2002).

- 1 9. McIntosh, L. T., Maheswaranathan, N., Nayebi, A., Ganguli, S. & Baccus, S. A. Deep
2 Learning Models of the Retinal Response to Natural Scenes. *Adv. Neural Inf. Process.*
3 *Syst.* **29**, 1369–1377 (2016).
- 4 10. Batty, E. *et al.* Multilayer Recurrent Network Models of Primate Retinal Ganglion Cell
5 Responses. in *International Conference on Learning Representations* (2017).
- 6 11. Wassle, H. & Boycott, B. B. Functional architecture of the mammalian retina. *Physiological*
7 *Reviews* vol. 71 447–480 (1991).
- 8 12. Saravanan, V., Berman, G. J. & Sober, S. J. Application of the hierarchical bootstrap to
9 multi-level data in neuroscience. *bioRxiv* 819334 (2019) doi:10.1101/819334.
- 10 13. Soto, F. *et al.* Efficient Coding by Midget and Parasol Ganglion Cells in the Human Retina.
11 *Neuron* **107**, 656–666.e5 (2020).
- 12 14. Kling, A. *et al.* Functional Organization of Midget and Parasol Ganglion Cells in the Human
13 Retina. 2020.08.07.240762 (2020) doi:10.1101/2020.08.07.240762.
- 14 15. Pillow, J. W. *et al.* Spatio-temporal correlations and visual signalling in a complete neuronal
15 population. *Nature* **454**, 995–999 (2008).
- 16 16. Schneidman, E., Brenner, N., Tishby, N., van Steveninck, R. R. de R. & Bialek, W.
17 Universality and Individuality in a Neural Code. in *Advances in Neural Information*
18 *Processing Systems 13* (eds. Leen, T. K., Dietterich, T. G. & Tresp, V.) 159–165 (MIT
19 Press, 2001).
- 20 17. Cahill, L. Why sex matters for neuroscience. *Nat. Rev. Neurosci.* **7**, 477–484 (2006).
- 21 18. Poplin, R. *et al.* Prediction of cardiovascular risk factors from retinal fundus photographs via
22 deep learning. *Nat Biomed Eng* **2**, 158–164 (2018).
- 23 19. Choleris, E., Galea, L. A. M., Sohrabji, F. & Frick, K. M. Sex differences in the brain:
24 Implications for behavioral and biomedical research. *Neurosci. Biobehav. Rev.* **85**, 126–145
25 (2018).
- 26 20. Color Blindness | National Eye Institute. <https://www.nei.nih.gov/learn-about-eye->

- 1 health/eye-conditions-and-diseases/color-blindness.
- 2 21. Bensmaia, S. J. & Miller, L. E. Restoring sensorimotor function through intracortical
3 interfaces: progress and looming challenges. *Nat. Rev. Neurosci.* **15**, 313–325 (2014).
- 4 22. Shenoy, K. V. & Carmena, J. M. Combining decoder design and neural adaptation in brain-
5 machine interfaces. *Neuron* **84**, 665–680 (2014).
- 6 23. Chestek, C. A. *et al.* Long-term stability of neural prosthetic control signals from silicon
7 cortical arrays in rhesus macaque motor cortex. *J. Neural Eng.* **8**, 045005 (2011).
- 8 24. Beyeler, M., Rokem, A., Boynton, G. M. & Fine, I. Learning to see again: Biological
9 constraints on cortical plasticity and the implications for sight restoration technologies. *J.*
10 *Neural Eng.* **14**, 051003 (2017).
- 11 25. Golub, M. D. *et al.* Learning by neural reassociation. *Nat. Neurosci.* **21**, 607 (2018).
- 12 26. Angueyra, J. M. & Rieke, F. Origin and effect of phototransduction noise in primate cone
13 photoreceptors. *Nat. Neurosci.* **16**, 1692–1700 (2013).
- 14 27. Schnapf, J. L., Nunn, B. J., Meister, M. & Baylor, D. A. Visual transduction in cones of the
15 monkey *Macaca fascicularis*. *J. Physiol.* **427**, 681–713 (1990).
- 16 28. Field, G. D. & Chichilnisky, E. J. Information processing in the primate retina: circuitry and
17 coding. *Annu. Rev. Neurosci.* **30**, 1–30 (2007).
- 18 29. Rousseeuw, P. J. & Croux, C. Alternatives to the Median Absolute Deviation. *J. Am. Stat.*
19 *Assoc.* **88**, 1273–1283 (1993).
- 20 30. Dumoulin, V., Shlens, J. & Kudlur, M. A Learned Representation For Artistic Style.
21 *International Conference on Learning Representations* (2017).
- 22 31. Kingma, D. P. & Ba, J. Adam: A Method for Stochastic Optimization. (2014).