

RESEARCH ARTICLE

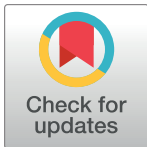
Normalization by orientation-tuned surround in human V1-V3

Zeming Fang^{1,2*}, Ilona M. Bloem¹, Catherine Olsson¹, Wei Ji Ma¹, Jonathan Winawer¹

1 Department of Psychology and Center for Neural Science, New York University, New York City, New York, United States of America, **2** Department of Cognitive Science, Rensselaer Polytechnic Institute, Troy, New York, United States of America

✉ Current address: Anthropic, San Francisco, California, United States of America

* zemingfang11@gmail.com



OPEN ACCESS

Citation: Fang Z, Bloem IM, Olsson C, Ma WJ, Winawer J (2023) Normalization by orientation-tuned surround in human V1-V3. PLoS Comput Biol 19(12): e1011704. <https://doi.org/10.1371/journal.pcbi.1011704>

Editor: Robbe L. T. Goris, University of Texas at Austin, UNITED STATES

Received: December 15, 2021

Accepted: November 20, 2023

Published: December 27, 2023

Copyright: © 2023 Fang et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Data Availability Statement: Analysis, data, and stimulus files are publicly available on the Open Science Framework site, accessible at <https://osf.io/vxjya/> (DOI [10.17605/OSF.IO/VXJYA](https://doi.org/10.17605/OSF.IO/VXJYA)). Analysis code is publicly available on GitHub under: <https://github.com/WinawerLab/StdVisualModel>. Data set 3 and 4 from the previously published study are available under: <http://kendrickkay.net/socmodel/>.

Funding: This study was supported by the National Eye Institute of the National Institutes of Health grant R01 EY027401 and the National Institute of Mental Health of the National Institutes of Health

Abstract

An influential account of neuronal responses in primary visual cortex is the normalized energy model. This model is often implemented as a multi-stage computation. The first stage is linear filtering. The second stage is the extraction of contrast energy, whereby a complex cell computes the squared and summed outputs of a pair of the linear filters in quadrature phase. The third stage is normalization, in which a local population of complex cells mutually inhibit one another. Because the population includes cells tuned to a range of orientations and spatial frequencies, the result is that the responses are effectively normalized by the local stimulus contrast. Here, using evidence from human functional MRI, we show that the classical model fails to account for the relative responses to two classes of stimuli: straight, parallel, band-passed contours (*gratings*), and curved, band-passed contours (*snakes*). The snakes elicit fMRI responses that are about twice as large as the gratings, yet a traditional divisive normalization model predicts responses that are about the same. Motivated by these observations and others from the literature, we implement a divisive normalization model in which cells matched in orientation tuning (“tuned normalization”) preferentially inhibit each other. We first show that this model accounts for differential responses to these two classes of stimuli. We then show that the model successfully generalizes to other band-pass textures, both in V1 and in extrastriate cortex (V2 and V3). We conclude that even in primary visual cortex, complex features of images such as the degree of heterogeneity, can have large effects on neural responses.

Author summary

How does the nervous system transform images into patterns of neural responses? A test of our understanding of this process is whether we can implement a computational model that takes digital images as input and accurately predicts responses in some part of the visual pathway. A widely used model of the transformation from image to neural response in primary visual cortex is the normalized energy model, in which image contrast drives neural activity, but the activity is self-limited because neural responses inhibit one another (“normalization”). We used functional MRI to measure the responses of human visual

grant R01 MH111417 to JW. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Competing interests: The authors have declared that no competing interests exist.

cortex while participants viewed a variety of images, and asked how accurately a normalized energy model could predict these responses. We find that if the model assumes that all neurons mutually inhibit one another (“untuned normalization”) it is not accurate, but if it assumes that only neurons tuned to similar image features inhibit one another (“tuned normalization”), predictions are much more accurate. We speculate that feature-tuned normalization helps the brain extract meaningful information about higher-order image statistics that are prevalent in many natural images.

1. Introduction

Primary visual cortex (“V1”) has served as a testing ground for studying physiology, anatomy, brain development, and neuroimaging. There has been considerable success in developing general model forms that capture many of the encoding properties of V1 neurons reasonably well over a range of stimulus conditions, including the normalized energy model of V1 complex cells [1]. This type of model, like many others [reviewed in chapter 6 of 2,3], includes a linear filter as the first stage, i.e., a weighted sum of the stimulus intensity over space and time. In a second stage, the outputs of the filter are squared and summed across nearby spatial locations or across phase [4–6]. If the outputs are summed across a pair of linear filters tuned to the same frequency, orientation, and location, but differing in phase by 90 deg, it is called an energy model. In the third stage, the response of each neuron is normalized (divisively suppressed) by the second-stage outputs of the nearby neural population [1,7]. This effectively adjusts the gain based on the contrast energy in the image patch. There is substantial evidence that each of these three operations—linear filtering, energy, and normalization—contributes to the responses of V1 neurons [reviewed by 8].

The normalized contrast energy model, though initially developed to explain the outputs of single neurons, has also been successfully applied to functional MRI data in human visual cortex. First, a contrast energy model without normalization, applied to voxels in V1, V2, and V3, was used to predict BOLD responses (encoding) and to infer the viewed images from the BOLD responses (decoding) [9]. Subsequent work showed that incorporating a normalization-like non-linearity improved model accuracy when testing stimuli that varied substantially in size [10] or pattern [11], and that normalization could account for the BOLD contrast response function for gratings with and without masking stimuli at other orientations [12]. These models have also shown good prediction accuracy for similar stimulus sets used in human intracranial electrode recordings of visual cortex [13,14].

The normalized contrast energy model, although successful at accounting for responses to a range of stimuli, nonetheless fails to explain some phenomena. There is some evidence that the standard model fit to artificial stimuli generalizes poorly to natural images [15–17] [but see also 18]. Even testing with simple patterns, early studies of V1 and extrastriate electrophysiology showed that some cells had tuning properties differing from simple or complex cells, called “hypercomplex” cells, many of which were associated with “end-stopping” [19,20]. Recent V1 two-photon calcium recordings included a large stimulus set and found that many cells, with or without end-stopping, were surprisingly sparsely tuned, often sensitive to complex patterns such as crosses or composite features [21]. It is unlikely that the standard energy model would predict the kind of tuning they observed, although they did not fit this model to their data. Other studies with single-unit electrophysiology found that a normalization model could be successful but only if the normalization was flexible, such that its strength depended on statistical dependencies in the image [22]. In closed-loop experiments in which models are used in

real time to find the optimal stimulus for mouse V1 cells, the most effective stimuli were often quite complex, differing from simple Gabor patterns [23]. In human fMRI studies, the responses to natural/complex images in V1 also appear to be influenced by statistical dependencies and image context [24,25], factors unlikely to influence the predictions of the normalized energy model. Even in relatively simple artificial images with a fixed amount of total contrast energy, the BOLD response is lower when there is a single orientation compared to when there are two divergent orientations [26,27].

In visual areas beyond V1, the normalized energy model is expected to be incomplete, as circuits in these areas contribute new computations. There are no widely adopted encoding models for these areas analogous to the normalized energy model for V1, but there has been some success in modeling patterns in the extrastriate responses by incorporating higher-order statistical dependencies of the modeled V1 outputs [28–30], higher-order statistical dependencies learned from natural image statistics [31], or sensitivity to second-order contrast [11]. There has also been progress in predicting V4 and IT responses from deep convolutional networks [32–34].

Here, using evidence from human functional MRI, we show that the classical normalized energy model fails to account for the relative responses to two classes of stimuli: straight, parallel, band-passed contours (*gratings*), and curved, band-passed contours (*snakes*) (Fig 1). The snakes elicit fMRI responses that are about twice as large as the gratings, yet traditional energy models, including normalized energy models, predict responses that are about the same. This is a large model failure which, in conjunction with the other failures of the simple normalization model described above, motivated us to implement a model in which the normalization is tuned, meaning that the normalization pool for a given neural channel has the same orientation tuning as the channel being normalized. We also developed and implemented a computational model that achieves tuned normalization in a different way, in which responses are normalized not by the sum of the contrast energy, but by the anisotropy (standard deviation)

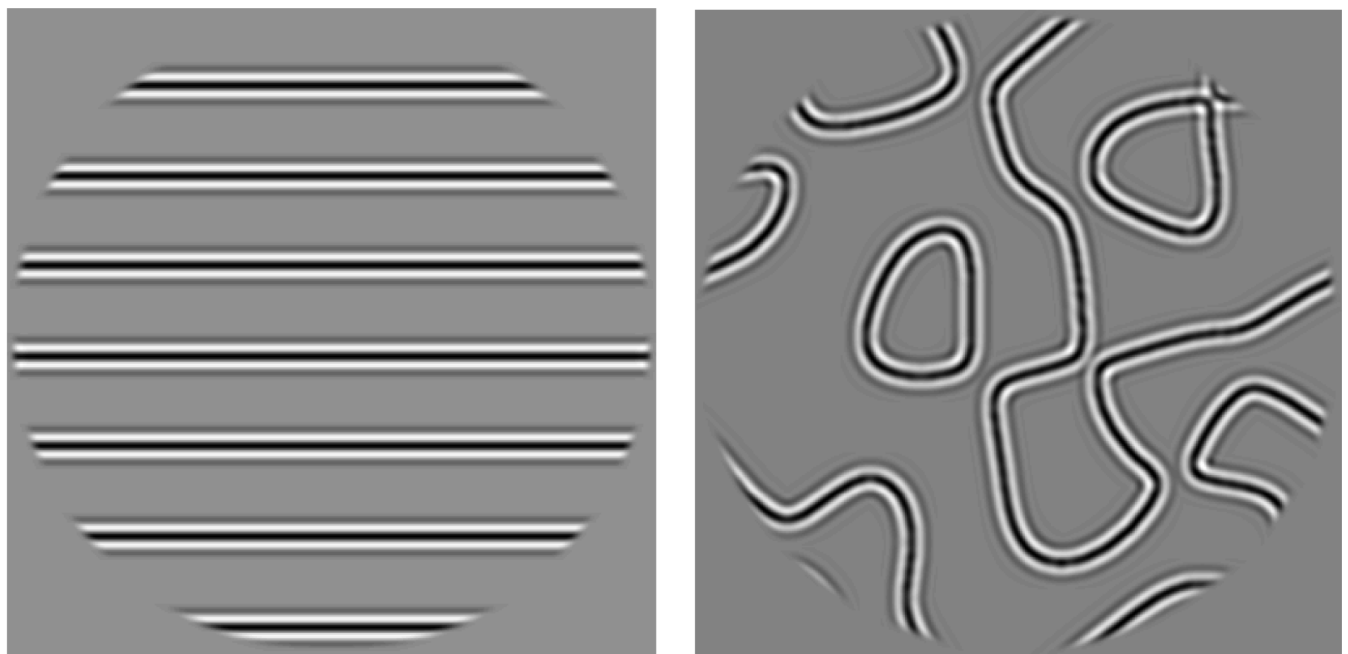


Fig 1. Example straight-line and curved-line stimuli from our experiments. We observe that in human V1, V2, and V3, stimuli with long straight lines (gratings) reliably evoke a smaller fMRI response than similar stimuli with curved lines.

<https://doi.org/10.1371/journal.pcbi.1011704.g001>

in contrast energy computed across orientation channels. Both models account for the differences in responses to snakes vs gratings, supporting the proposal that normalization depends on the spatial arrangement of image features, and not just the total amount of contrast energy.

2. Results

2.1. fMRI BOLD responses to snakes are larger than to gratings

We first consider an observation about fMRI responses to two classes of simple, grayscale, band-passed, static images. For one class, the stimuli contain several curved contours, which we refer to throughout as snakes. For the other class, the stimuli contain several straight, parallel contours, which we refer to as gratings. We refer to these classes together as the *target stimuli*. The surprising observation is that for V1, V2 and V3, the fMRI responses are substantially larger for the snakes than for gratings (Fig 2). The responses to the gratings, irrespective of density, are only about as high as the lowest response to the snakes. We confirm this pattern with three additional fMRI data sets, which also show larger responses to snakes than gratings in V1, V2 and V3 (Figs A1–A4 in S5 Appendix).

To check whether this pattern of results is specific to fMRI, we replot published intracranial data from Hermes et al., 2019 [13], in which human subjects with ECoG recordings viewed a similar set of stimuli (Fig 2, bottom panel). The ECoG data, plotted as the percent increase in broadband power over baseline, show the same general effect as the fMRI data: The responses to snakes are much larger than to gratings, irrespective of texture density, indicating that the effect is not limited to the fMRI BOLD response. (See Hermes et al., 2019 [13] for methodological details.)

In the next four sections, we describe the four models that are fit to the data. For tractability, we fit each model's parameters to the aggregate (average) data within a visual area, rather than to each voxel individually. All the stimuli are texture-like, meaning that they have similar properties across the whole image aperture, and for each stimulus class, nine different exemplars were shown per 3-s trial. The different exemplars have the same higher-order statistics but vary in their precise spatial distributions. Hence, model variables, such as contrast energy, would have similar values for spatially localized portions of the image (as one would compute for an individual voxel) as for the whole image (as we compute to model the aggregate response of a visual area). For this reason, we did not include model parameters for the spatial location or spatial extent of the receptive fields for individual voxels. We first describe models fit to the target stimuli. We then summarize fits to the larger set of stimuli viewed by the subjects.

All four models consist of three components, *filtering*, *spatial pooling*, and a *power-law nonlinearity*. In the filtering stage, we computed the contrast energy of the stimuli at 8 orientations and 4 spatial frequencies (see Methods for details). Although stimuli were designed to have power concentrated close to 3 cycles per degree, there is some spillover to lower and higher frequencies, which is why we use multiple spatial frequency bands in our models. The spatial pooling stage pools the contrast energy to yield a total contrast energy. In some models, the pooling stage also includes divisive normalization. The output of the pooling is a scalar, which is then passed through a power-law nonlinearity to predict BOLD amplitude in units of percent signal change. The power-law nonlinearity achieves compressive spatial summation [10]. All models have the same filtering (first stage) and output non-linearity (third stage). They differ in the spatial pooling stage.

2.2. The larger response to snakes is not captured by a simple contrast energy model

A *standard contrast energy* model pools the contrast energy by simply summing it across space, orientations, and spatial frequencies to give a total contrast energy. It predicts that

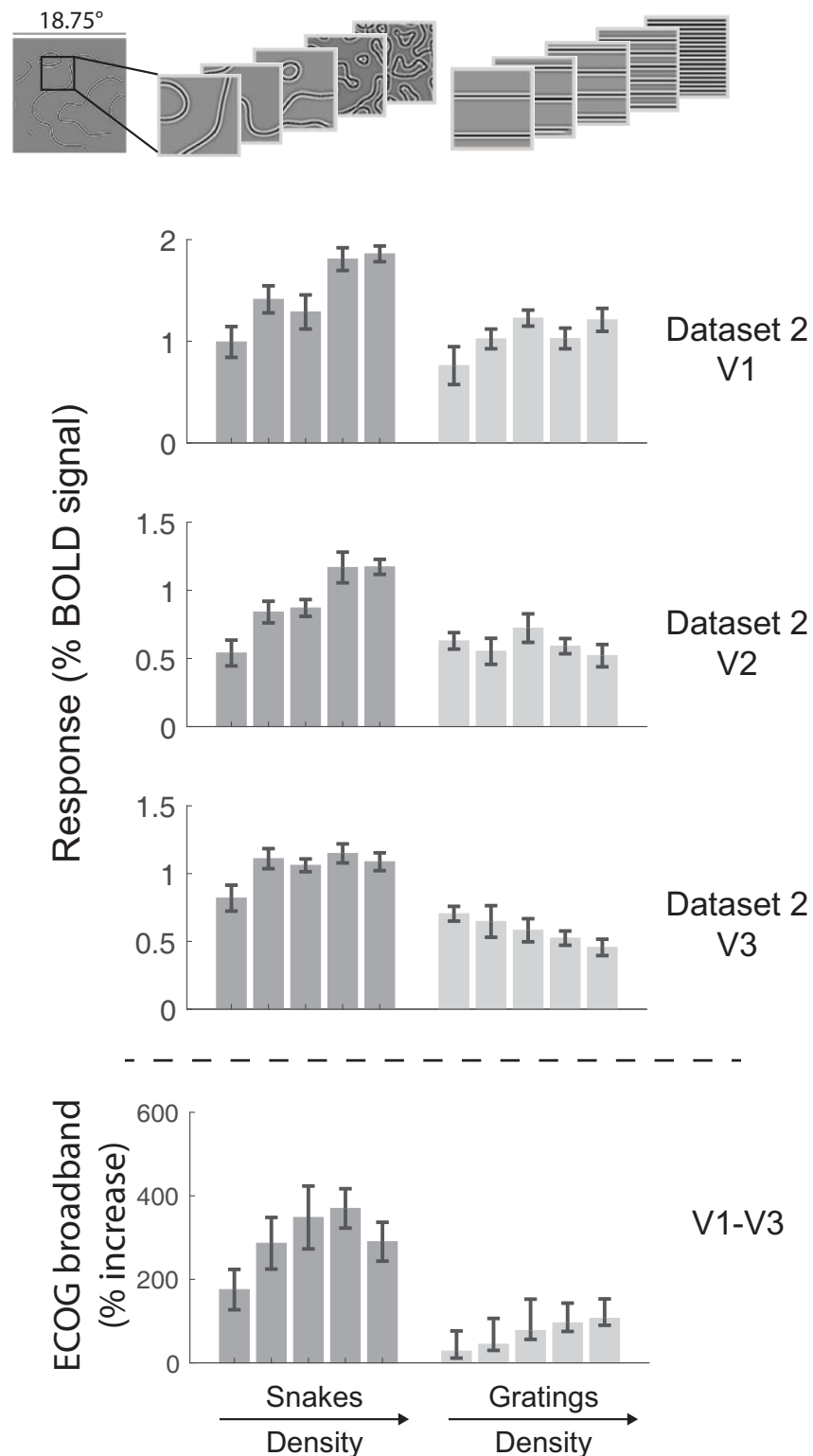


Fig 2. FMRI responses are larger for curved patterns than for straight patterns. The mean fMRI responses within a visual area are plotted for V1, V2, and V3 in Data set 2 of Table A in S1 Appendix. For both curved stimuli (snakes, dark bars) and straight stimuli (gratings, light bars), data are plotted in order of increasing texture density from left to right. Examples of the 5 densities are shown above. Error bars are the standard deviation of the mean, bootstrapped across fMRI runs. Responses are larger for the snakes than gratings. The same effect is also observed for Data sets 1, 3

and 4 of Table A in [S1 Appendix](#) (Fig A in [S5 Appendix](#)). Bottom panel: For comparison, we replot data from intracranial recordings (ECoG) from Hermes et al 2019 [13], which also show larger responses for snakes than for gratings. The full set of stimuli is shown in [S2 Appendix](#). Data plotted from the function `s4_visualize('figure 2')` in the [GitHub code repository](#).

<https://doi.org/10.1371/journal.pcbi.1011704.g002>

responses should increase with both stimulus contrast and with density of the pattern. This model does not predict a larger response to the snakes than to gratings, contrary to the data ([Fig 3](#)). In fact, the cross-validated variance explained is low (V1, V2) or even negative (V3) in the example data, meaning that the model prediction is less accurate than it would have been if it simply predicted the mean response across all stimuli. (The data are cross-validated, which is why the variance explained can be negative). In short, the contrast energy model provides a poor fit to the fMRI data in V1-V3 for these classes of stimuli. It is also a poor fit to the target stimuli in the other three data sets ([Table A in S3 Appendix](#) and [Fig A1 in S5 Appendix](#)). This does not mean that contrast energy models are always poor fits to fMRI responses in V1-V3. For example, when stimuli vary in how contrast energy is distributed across space, a contrast energy model can capture a lot of the variance in the responses across images, as shown by Kay et al., 2013 [9].

2.3. The larger response to snakes is not captured by an untuned normalization model

We then add divisive normalization to the model. After computing contrast energy, we normalize the outputs by dividing the contrast energy at each pixel by the contrast energy of a normalization pool ([Fig 4, upper panel](#)). The normalization pool includes nearby locations, all spatial frequencies, and all orientations, giving it the name *untuned normalization model* ([Fig 4, lower left panel](#)).

The untuned normalization model and the contrast energy model make similar predictions and explain a similar proportion of the variance in the data. The normalization model results in more saturation at high contrast, as expected from divisive normalization [7]. This is especially evident in V2 and V3 at the highest stimulus contrast. The reason that the normalization model and the contrast energy model have a similar overall pattern of predictions is that the power law output nonlinearity, included in all models, can partially mimic the effects of normalization [10, section “Relationship to Divisive Normalization”].

Like the contrast energy model, the untuned normalization model predicts a similar BOLD amplitude for snakes and gratings, thereby failing to account for the data ([Fig 5, upper panel](#)). The model is not sensitive to heterogeneity across orientations because the normalization pool equally weights all orientation channels. Therefore, the output does not depend on whether the contrast energy is concentrated in one orientation channel, as in the gratings, or spread across many channels, as in the snakes. The untuned normalization model’s failure to account for the greater response to snakes is reflected by low variance explained for the target stimuli in each of the four data sets ([Table A in S3 Appendix](#) and [Fig A2 in S5 Appendix](#)).

2.4. The larger response to snakes is captured by an orientation-tuned normalization model

The untuned normalization model implements a surround suppression that includes energy at all orientations. Findings from electrophysiology [35–39], psychophysics [38,40–42], neuroimaging [26,27,43–45], and theory [46], however, suggest that surround suppression is orientation-tuned. For example, Cavanaugh, Bair, and Movshon, 2002 [35] reported that the response of a neuron to a stimulus at its preferred orientation in its receptive field is suppressed more

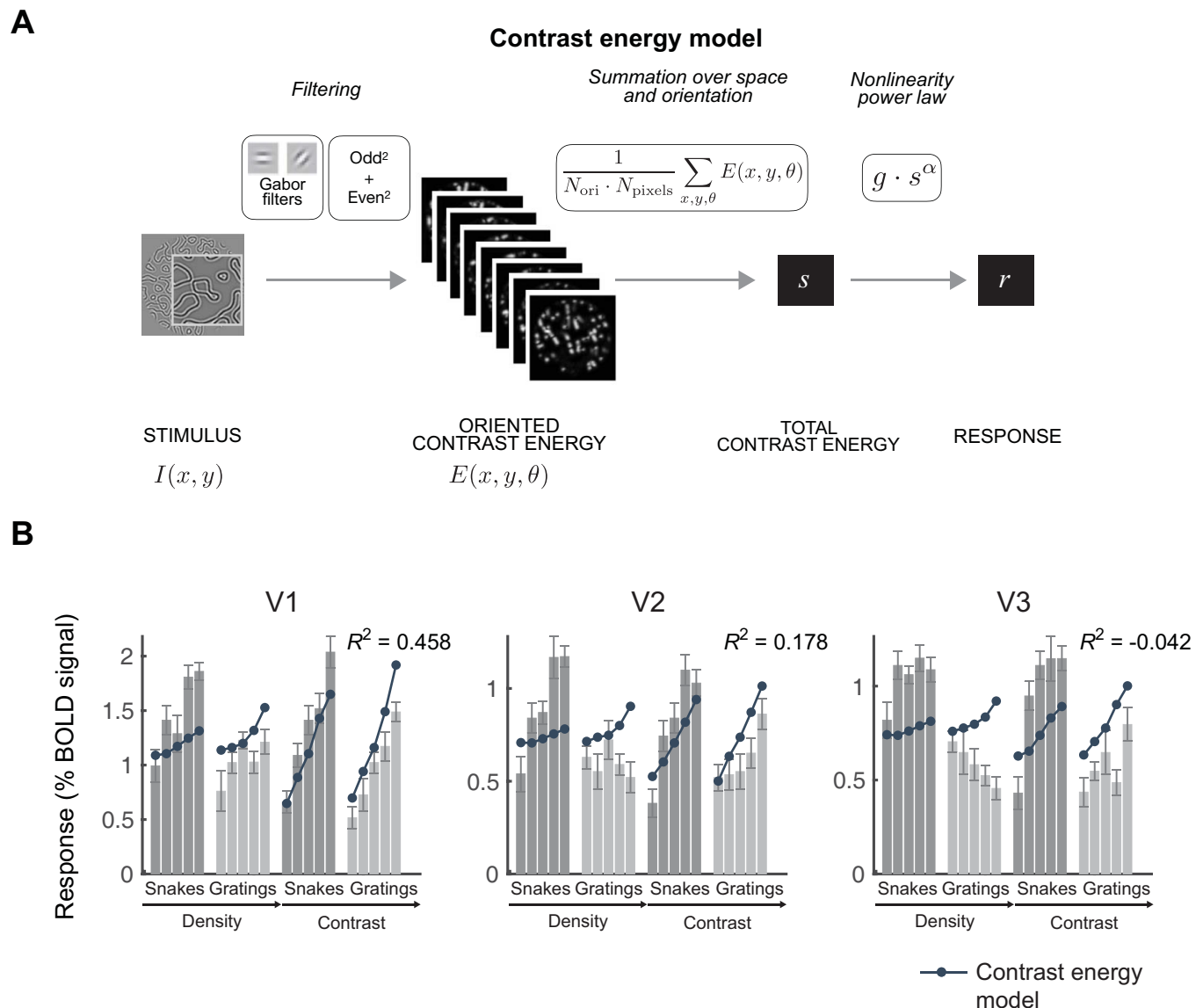


Fig 3. The contrast energy model does not account for V1-V3 responses to snakes and gratings. (A) Schematic representation of the contrast energy model. For simplicity, we show E summed across spatial frequency bands. (B) Mean fMRI responses in V1, V2, and V3 to snake and grating stimuli that vary in density and contrast (Data set 2). Bars: mean and standard error of the responses. Dark bars represent snake stimuli and light bars represent grating stimuli. Each group of stimuli is arranged in increasing order of either density or contrast. Dots: cross-validated predictions from the contrast energy model. See Fig A1 in S5 Appendix for fits to all 4 data sets. See S4 Appendix for model parameters. Data and model fits plotted using the function `s4_visualize('figure 3')` in the [code repository](https://doi.org/10.1371/journal.pcbi.1011704.g003).

<https://doi.org/10.1371/journal.pcbi.1011704.g003>

when the surrounding region contains contrast at the same orientation compared to different orientations. Because our grating stimuli have contrast energy concentrated at a single orientation, and the snake stimuli do not, one might surmise that an *orientation-tuned normalization* model would show greater suppression for the gratings, where the RF centers and surrounds will have matched orientations, than for the snakes, where the orientations are more likely to differ between center and surround. If so, this could then account for our observed effect.

The untuned normalization model is the same as the tuned model except for one difference: At each pixel in the oriented contrast energy images, the tuned model normalizes the contrast

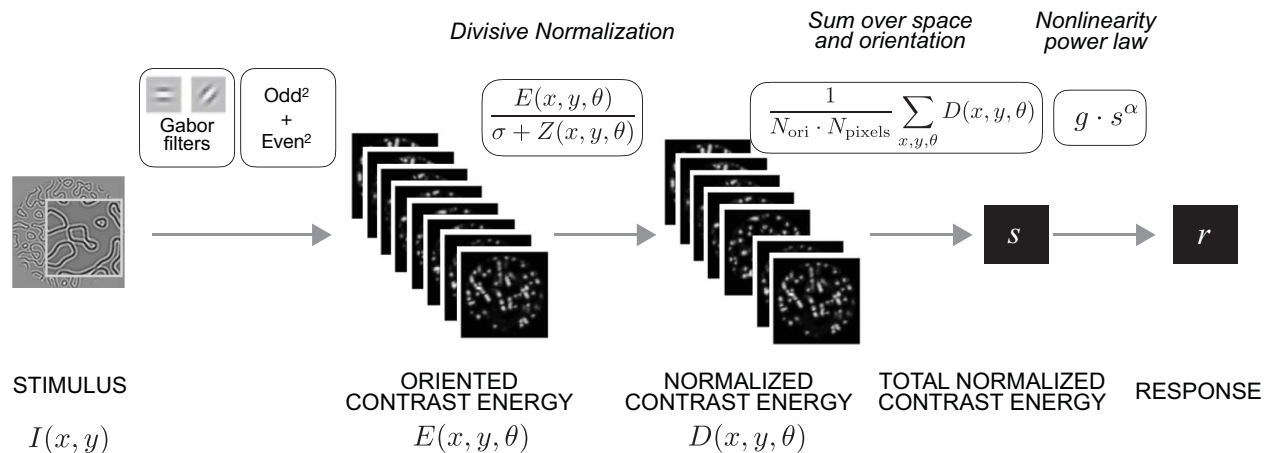
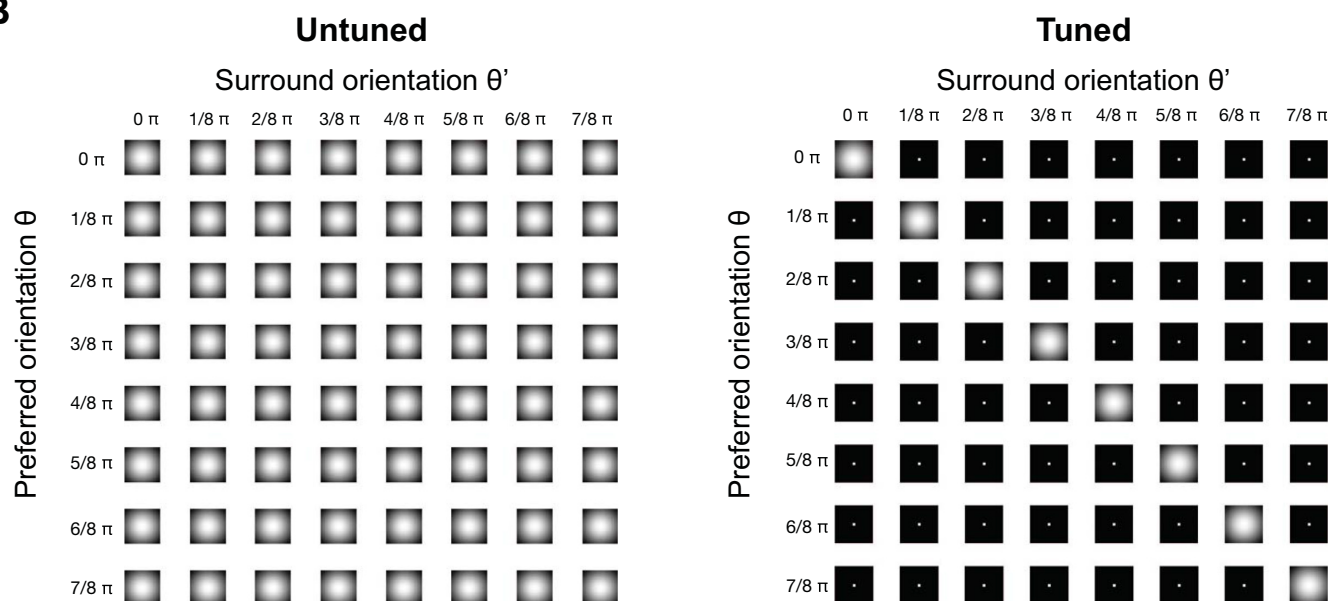
A**B**

Fig 4. Two divisive normalization models. (A) Schematic of divisive normalization models. (B) The weights used to calculate normalization for the untuned (left) and tuned (right) models, plotted using the function `s4_visualize('figure 4')` in the [code repository](https://doi.org/10.1371/journal.pcbi.1011704.g004).

<https://doi.org/10.1371/journal.pcbi.1011704.g004>

energy across nearby locations only at the preferred orientation (orientation-tuned surround) (Fig 4, lower right panel, diagonals). Within a single location (i.e., at each pixel), the normalization is untuned (off-diagonals in the same panel), also called cross-orientation suppression [47].

The orientation-tuned normalization model captures the large and systematic differences in response amplitude between gratings and snakes (Fig 5, lower panel). In the example data set, the BOLD amplitude and the orientation-tuned normalization model predictions for the gratings are about half of those to the snakes, for both the density and contrast manipulations, and for all three visual areas. In addition to capturing this difference in the means between the two stimulus classes, the model also captures the difference in slope. As the density or contrast increases, the model predicts steeper slopes for snakes than gratings. These patterns in the

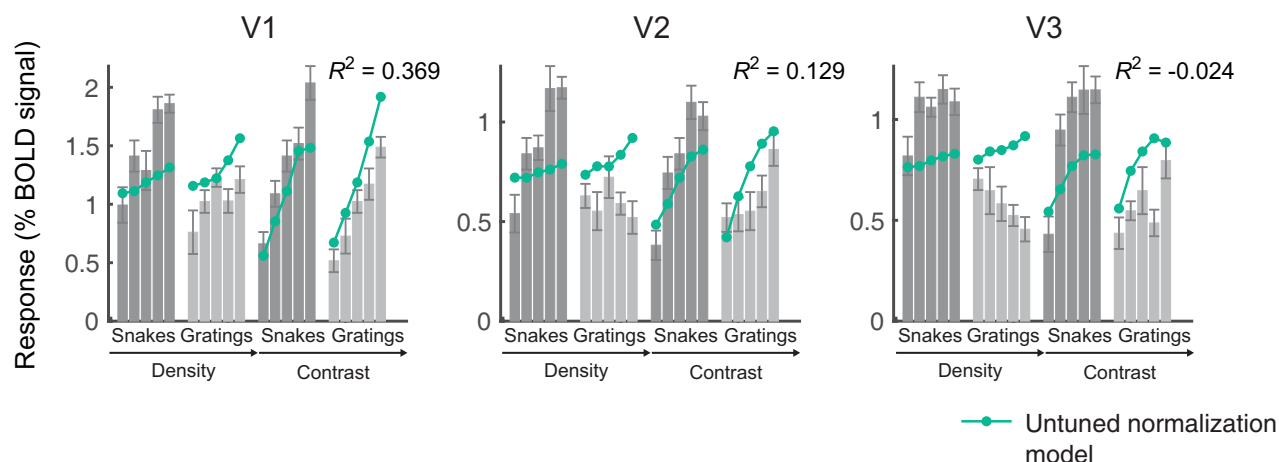
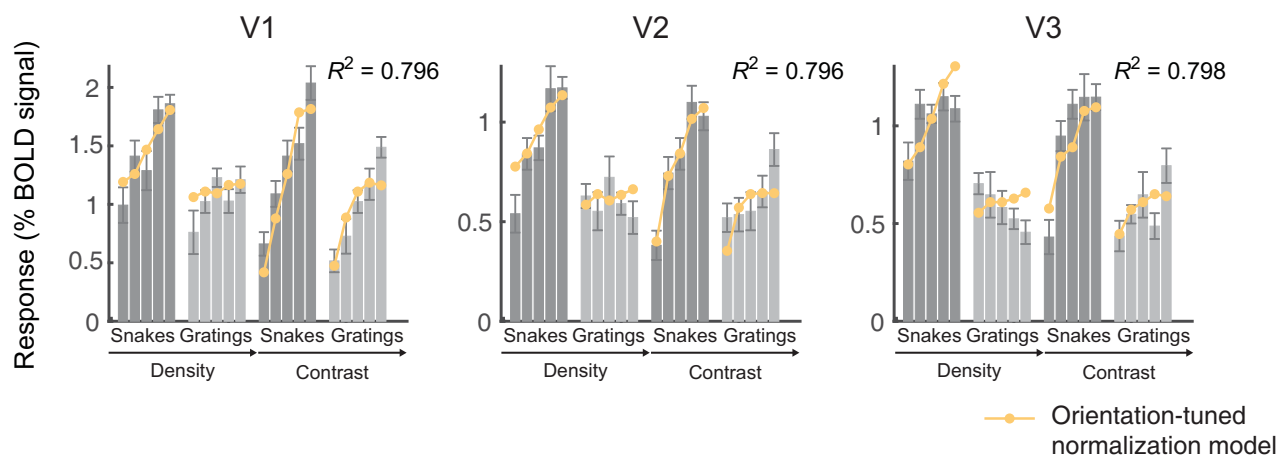
A**Untuned normalization model****B****Orientation-tuned normalization model**

Fig 5. The tuned normalization model is much more accurate than the untuned normalization model. The mean fMRI responses from V1, V2, and V3 are replotted from Fig 3. The dots are the cross-validated predictions from the untuned normalization model (above) or the tuned normalization model (below). Data and model fits plotted from the function `s4_visualize('figure 5A')` and `s4_visualize('figure 5B')` in the [code repository](#). See Figs A2-A3 in S5 Appendix for fits to all 4 data sets.

<https://doi.org/10.1371/journal.pcbi.1011704.g005>

model fits are found across the four data sets (Fig A3 in S5 Appendix). The orientation-tuned normalization model is more accurate than the previous two models in all cases (Table A in S3 Appendix, four data sets and three ROIs). This result holds up for different size surrounds—what mattered was whether the surround was tuned or not tuned, not its size. The indifference to the size of the surround almost certainly reflects the properties of the stimulus set, not neural tuning: the stimuli are all textures, with similar properties across the image. Had the stimuli varied systematically across location, the size of the surround would likely have had a large effect on model accuracy. Note that the data sets sometimes show a negative slope with increasing density (e.g., V3, gratings varying in density). The model is unable to capture this effect. The study in which these sparse stimuli were first used [11] showed that a second-order contrast model could account for the decreased response with increasing sparsity, as the

sparser stimuli have more second-order contrast. It is possible that the orientation-tuned normalization model would also be able to do so if it included spatial receptive fields per voxel that were small relative to the image.

2.5. Normalization by orientation anisotropy

The large advantage in prediction accuracy of the tuned over the untuned normalization model supports the idea that suppression is feature-specific. As with any model, we chose a specific instantiation of a more general idea, namely feature-specific suppression. The specific instantiation entailed a minimal change from the untuned normalization model, requiring only a change in normalization weights, and builds on the tradition of feedforward, filter-based models. Feature-specific tuning can also be implemented in other ways, for example based on more abstract ideas like predictability or redundancy in the image. There is some evidence for models like these [22,48]. We implemented a second method of achieving orientation-tuned normalization, in which normalization was proportional to *orientation anisotropy* (Fig 6, Normalization by orientation anisotropy, “NOA”). Normalization in this model is most pronounced when an image patch has a single orientation, without a specification in terms of a match between center and surround (See 3.3 *What is the tuning in orientation-tuned normalization?*). Specifically, in the normalization by orientation anisotropy model, the contrast energy is normalized by the standard deviation across the outputs of the orientation channels. This normalization by anisotropy model applies greater normalization when the contrast energy is concentrated in a single orientation channel, resulting in a lower response for gratings. There is no explicit representation of centers, surrounds, or feature matching in the normalization pool. This implementation is consistent with the idea that responses are reduced by the amount of redundancy in the image.

The normalization by anisotropy model exhibits similar predictions to the orientation-tuned normalization model, capturing the larger response to the snakes (Fig 6). Both models predict that the responses to snakes are about double the responses to gratings, similar to the data. It also predicts a higher slope for the snakes than the gratings, both as a function of density and contrast. The success of these models validates the idea that normalization depends on how contrast energy is distributed across orientations, not just on the overall contrast energy.

2.6. The two normalization models that are sensitive to orientation accurately predict responses to a wide range of stimuli

Both models without orientation sensitivity fail to account for the higher responses to snakes than gratings (Fig 7), and have low variance explained (Fig 8). Models that include orientation sensitivity capture the higher responses to snakes (Fig 7), and fit the data accurately (Fig 8). This suggests that sensitivity to orientation should be incorporated into normalization models of visual cortex. The accuracy differences across models are not due to the number of free parameters: the untuned normalization model has the same number of free parameters (three) as the tuned normalization model and the anisotropy model. Moreover, the prediction accuracy was computed using cross-validation, so that having more parameters does not necessarily lead to better predictions. The tuned normalization model has a numerically higher accuracy than the anisotropy model for nearly all data sets in all conditions (S3 Appendix), but the advantage of the two models with orientation sensitivity dwarfs the slight difference between these two models. Interestingly, for the two untuned models, prediction accuracy declines from V1 to V2 to V3, whereas for the two tuned models, accuracy increases (target stimuli) or stays flat (all stimuli) from V1 to V2/V3 (black lines in Fig 8). This pattern is consistent with the notion that along the visual hierarchy, neural responses become increasingly

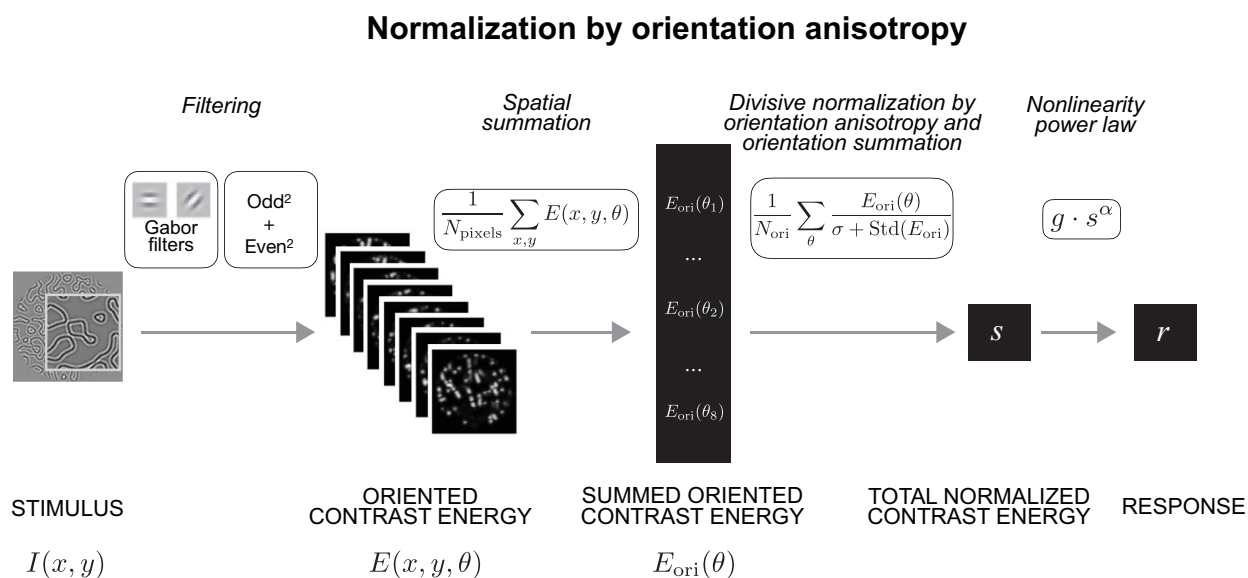
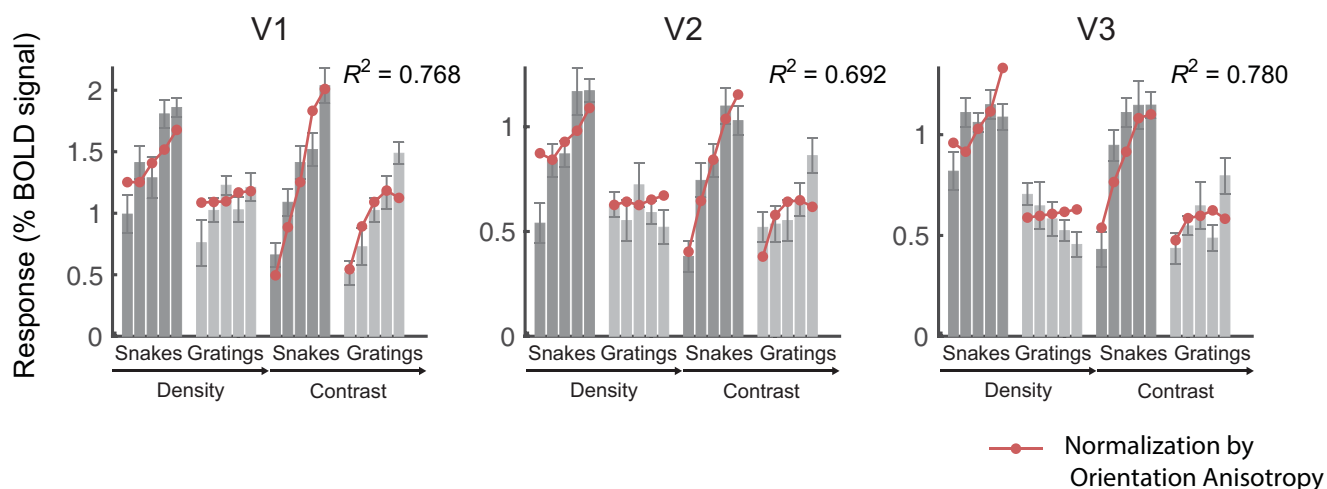
A**B**

Fig 6. The normalization by orientation anisotropy model also accounts for the responses in V1-V3. (A) Schematic of the normalization by orientation anisotropy model. (B) The mean fMRI responses from V1, V2, and V3 are replotted from Fig 3. The red dots are the cross-validated predictions from the normalization by orientation anisotropy model. Data and model fits plotted from the function `s4_visualize('figure 6')` in the [code repository](#). See Fig A1 in S5 Appendix for fits to all 4 data sets.

<https://doi.org/10.1371/journal.pcbi.1011704.g006>

sensitive to statistical regularities, such as similarity in features across the image (see [Discussion](#) section 3.4).

Because the two orientation-sensitive model were motivated by the need to explain the greater response to snakes than gratings, it is important to test the models on other stimuli as well. We refit all 4 models to the full data sets, which consisted of 50 (data set 1), 48 (data set 2) and 39 (data set 3, data set 4) stimuli, spanning a variety of texture types. In addition to the snakes and gratings, there are textures we refer to as *noise bars*, *waves*, *plaids*, and *circular* ([S2 Appendix](#) and [Table C in S1 Appendix](#)).

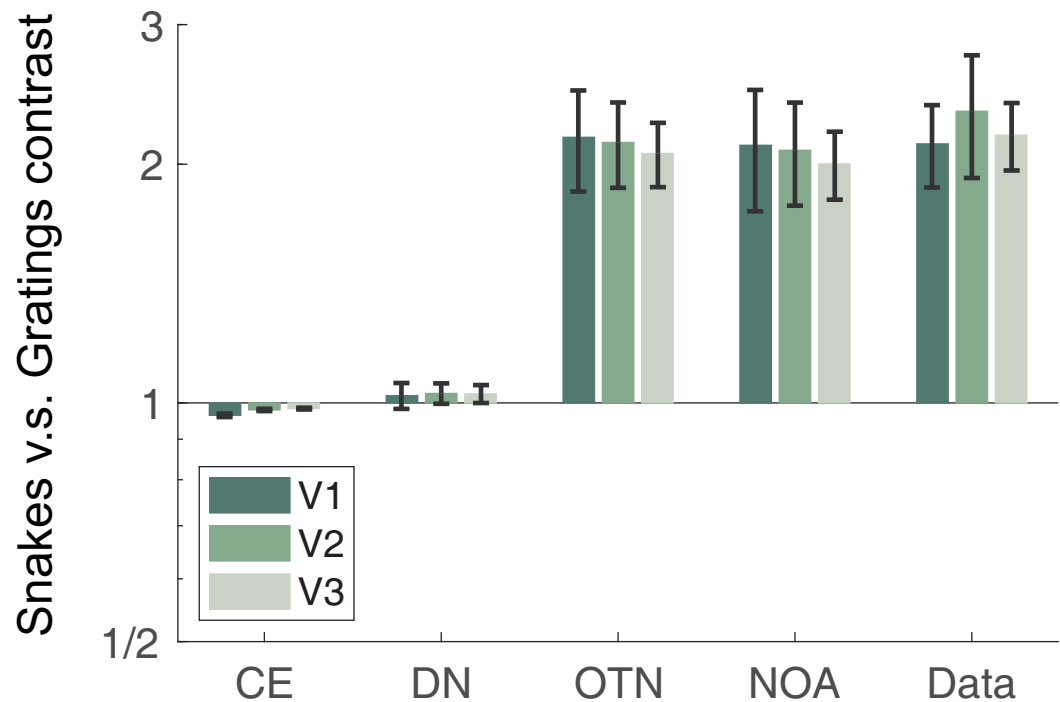


Fig 7. The orientation dependent models account for the higher responses to snakes than gratings. For each visual area, V1-V3, we averaged the response to snakes and to gratings across stimuli in the target set and across the 4 data sets to compute the ratio of snakes to gratings: $\text{mean}(\text{snakes}) / \text{mean}(\text{gratings})$. We computed this value for each data set and plotted the average and standard error across the four data sets. In the data, the response to snakes is about double to gratings. This is matched in the two normalization models that are sensitive to orientation, but not the other models. Data and model fits plotted from the function `s4_visualize('figure 7')` in the [code repository](#). CE = contrast energy; DN = untuned normalization; OTN = orientation-tuned normalization; NOA = normalization by anisotropy.

<https://doi.org/10.1371/journal.pcbi.1011704.g007>

Just as with the target stimuli, across the full sets of stimuli, the tuned normalization model and the anisotropy model made accurate predictions, explaining 63%-77% and 49%-66% of the cross-validated variance in V1-V3 for the example data set (Fig 9). These two models also provide good fits to the other three data sets, shown in S5 Appendix. The fits to the larger stimulus sets, like the fits to the target stimuli alone, capture the observation about the two stimulus classes, meaning a larger predicted response for snakes than gratings. The two models also accurately predict lower responses to waves (one dominant orientation) than noise bars (many orientations). The two models also predict increasing response amplitudes from gratings (one orientation) to plaids (two orientations) to circular (16 orientations), as evident in stimulus sets 3 and 4 (Figs D3-D4 and E3-E4 in S5 Appendix). This pattern of predictions matches the data. The untuned models do not differ in their predictions for these three stimulus categories.

Across ROIs and data sets, the orientation-tuned normalization model accounts for the highest variance, with R^2 ranging from 55% to 77% (Fig 8, bottom; S3 Appendix, the third row). The anisotropy model ranks second in all cases, substantially outperforming the two baseline models. Similar to the pattern with the target stimuli, when fitting to all stimuli the two untuned models show substantially decreasing accuracy from V1 to V2 to V3. The tuned normalization model and the anisotropy model decrease only slightly in accuracy from V1 to V2 to V3, meaning that the advantage for the tuned models is largest in extra-striate areas.

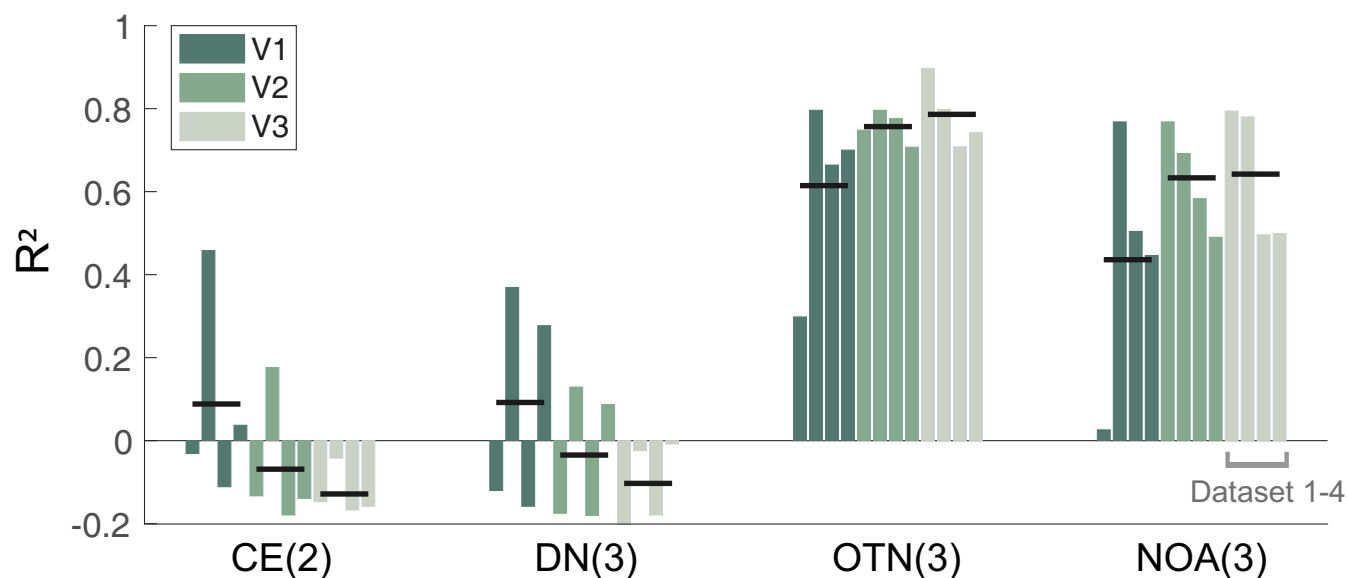
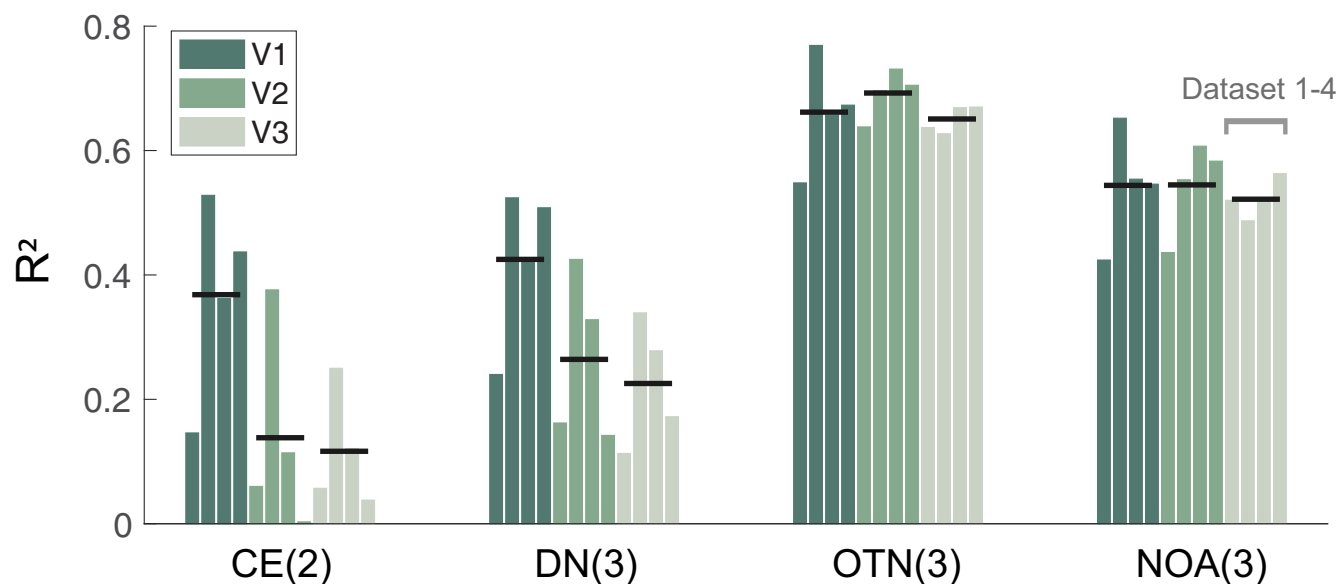
A**Cross-validated variance explained on *target* stimuli****B****Cross-validated variance explained on *all* stimuli**

Fig 8. Cross-validated variance explained for 4 models across all ROIs and data sets. (A) R^2 for the target data set, comprising 18 stimuli for data set 1 (DS1) and data set 2 (DS2) and 17 for data set 3 (DS3), data set (DS4). The number of fitted model parameters (degrees of freedom) is indicated in parentheses for each model type. The four bars in each group correspond to data sets 1–4. The black horizontal lines are the means across the 4 data sets. (B) Same as panel A, but for a larger set of stimuli (50 for data set 1; 48 for data set 2; 39 for data set 3 and 4). The R^2 values are also reported in S3 Appendix. Data plotted from the function `s4_visualize('figure 8')` in the [code repository](https://doi.org/10.1371/journal.pcbi.1011704.g008). Abbreviations as in Fig 7.

<https://doi.org/10.1371/journal.pcbi.1011704.g008>

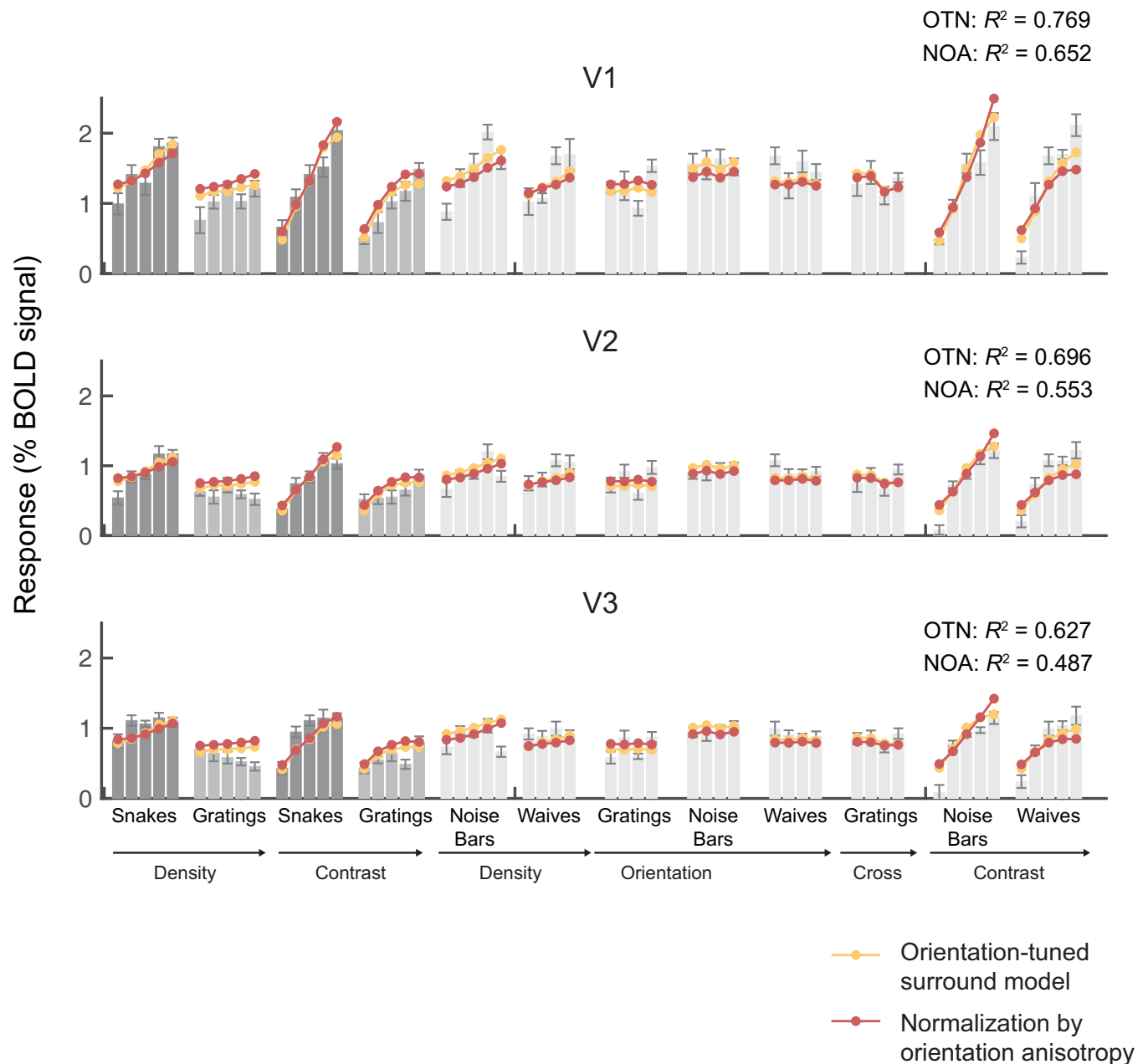


Fig 9. The orientation dependent models fit a wide variety of stimuli. The mean fMRI responses from V1, V2, and V3 are shown for the full set of stimuli from data set 2 (See [S2 Appendix](#)). The data for snakes and gratings are replotted from Figs 3, 5 and 6, and again shown in dark and light gray. An even lighter gray is used for all other stimulus classes. The red and yellow dots are the cross-validated predictions from the orientation-tuned normalization model and the normalization by orientation anisotropy model. See Figs B-E in [S5 Appendix](#) for similar plots for data sets 2–4. Data and model fits plotted from the function `s4_visualize('figure 9')` in the [code repository](#).

<https://doi.org/10.1371/journal.pcbi.1011704.g009>

3. Discussion

3.1 Why do some models fail to account for the much larger response to snakes?

We began with the observation that the BOLD response is about twice as large for patterns with curved contours (snakes) as for similar stimuli with straight, parallel contours (gratings).

The difference in the responses is not a peculiarity of the BOLD signal, as a similar pattern was also observed in human intracranial measures of the field potential (broadband power from ECoG electrodes). The contrast energy model without normalization and the contrast energy model with untuned normalization did not predict large differences between the responses to the two stimulus classes. Both models pool contrast energy over space and orientation without an orientation-specific normalization or other orientation specific non-linearity. If two images have the same total contrast energy, then the way the energy is distributed across orientation channels will not matter, either for the measure of total energy (contrast energy model) or for the amount of normalization (untuned normalization model).

The implementation of the orientation-tuned normalization model was motivated in part from electrophysiology results showing that in V1, surround suppressive fields tend to be tuned to orientations close to the RF center's preferred orientation [35]. Psychophysical experiments also show interactions between a target and surround that depend on matched orientation [e.g., 42]. While there is a lot of evidence from electrophysiology, fMRI, and psychophysics in support of feature tuning in normalization, here we explicitly compared image-computable models with and without feature-tuned normalization. Implementing the model and fitting it to data enabled us to assess whether (1) it quantitatively accounts for two-fold difference in response to snakes vs gratings (2) whether it also provides good fits to a wide range of other stimuli that it was not explicitly implemented for and (3) whether it can capture other more subtle effects, like the difference in the slope of the contrast response function between snakes and gratings. In all three cases, the answer was yes. We also implemented a second model with tuned normalization (the normalization by orientation anisotropy model), and this model also provided an affirmative answer to these questions, confirming the importance of including orientation dependence in the normalization computation. We discuss similarities (3.2) and differences (3.3) between the two models below.

3.2 Model behavior: Models with orientation dependent normalization capture the differences in both *mean* and in *slope* between snakes and gratings

The two models with orientation dependent normalization predict larger outputs, *on average*, for snakes than gratings, even when the contrast energy of the stimuli is approximately matched. This is due to how the normalization is computed. The grating stimuli elicit large outputs in the orientation channels that are matched to the stimulus, moderate outputs in adjacent orientation channels, and little to no response in other channels. As a result, there is high anisotropy (standard deviation across channel outputs), resulting in more suppression in the normalization by orientation anisotropy model. There is also high suppression from the surround in the orientation-tuned normalization model. The two models were implemented to capture the difference in mean between the two stimulus classes, so perhaps it is not so surprising that they do so.

The two models, unlike the contrast energy model and the untuned normalization model, also accurately predict *steeper slopes* for snakes than gratings (with respect to both contrast and density), a pattern that the models were not explicitly motivated to capture. They predict the difference in slope because at low total contrast energy for the image, there is little normalization, and hence the response to a snake and a grating stimulus will be comparable. This explanation applies to both low contrast stimuli and low-density stimuli, because in both cases the summed contrast energy is low, meaning the normalization term in the denominator is small. Hence, little normalization at low contrast is expected from the model, and is also confirmed by empirical measures of spatial summation and surround suppression at different contrast

levels [49,50]. The more nuanced prediction from these models is that for stimuli with high total contrast energy (high stimulus contrast and high density), there is a lot of normalization for gratings and much less for snakes, resulting in a more pronounced difference in predicted response. The difference in predicted responses at high contrast but not at low contrast causes a difference in slope. For the two untuned models, normalization increases with contrast energy, but it increases similarly for the snakes and gratings, hence predicting similar slopes.

Interestingly, the data and the two models with orientation dependent normalization also show a greater slope for “noise bar” stimuli than “waves” (Figs B3-B4 and C3-C4 in [S5 Appendix](#)). The noise bars, like the snakes, have many orientations in first order contrast (but unlike the snakes, have only one dominant orientation for second-order contrast). The waves are the complement, with one dominant orientation for first order contrast (like gratings) but many orientations for second-order contrast. The pattern in the data and model predictions is that the stimuli with many orientations (noise bars) increase more steeply as a function of contrast than the stimuli with a narrower range of orientations (waves), supporting the observations made with snakes and gratings, but adding the further nuance that what seems to matter most is the orientation distribution (wide vs narrow) for first order rather than second-order contrast.

The difference in slopes between snakes and gratings in the data is related to, but not identical to, results observed for single unit V1 cells. For a typical V1 cell, the contrast response function has a higher slope for preferred than for non-preferred stimuli [51–53], a pattern also predicted by normalization models [7]. The pattern is predicted because of the difference in the *numerator*, which is high for preferred stimuli and low for non-preferred. The denominator is about the same for the two stimuli, proportional to the total contrast energy in the stimulus. We attribute the difference in slope between snakes vs gratings to a difference in the *denominator* of the normalization equation, which is high for gratings, low for snakes. The numerator is about the same for the two classes, proportional to the total contrast energy. In both cases—single unit responses to non-preferred stimuli, and population responses to grating stimuli—the reduced slope is due to a large amount of normalization relative to the driven response.

3.3 What is the tuning in orientation-tuned normalization?

We implemented two models with orientation dependent normalization. Both capture the same tendency for more normalization for homogeneous stimuli like gratings. They differ in implementation, however. The tuned normalization model differs only slightly from more typical normalization models: the only difference is that the normalization weights happen to conform to a specific pattern, such that cells with similar feature tuning (here, orientation) with nearby receptive fields have high weights, and cells with different feature tuning have low weights. As a result, surround suppression is most effective according to this model when the orientation of a surrounding region is matched to the preferred orientation of a cell. This is justified by findings from single unit data in macaque V1 that

“the surround influence was always suppressive when the surround grating was at the neuron’s preferred orientation” [35].

The normalization by orientation anisotropy model is a larger departure from the standard normalization model, since the computation of anisotropy is not a simple weighting of the outputs of nearby cells. Unlike the tuned normalization model, it has the greatest suppressive effect when the features of a stimulus are anisotropic (like a grating) irrespective of whether the stimulus orientation matches the center tuning of a cell. Interestingly, there is also empirical support for this pattern from the same study by Cavanaugh et al., 2002 [35]: specifically, evidence for “the tuning of the surround being dependent to some degree on the stimulus

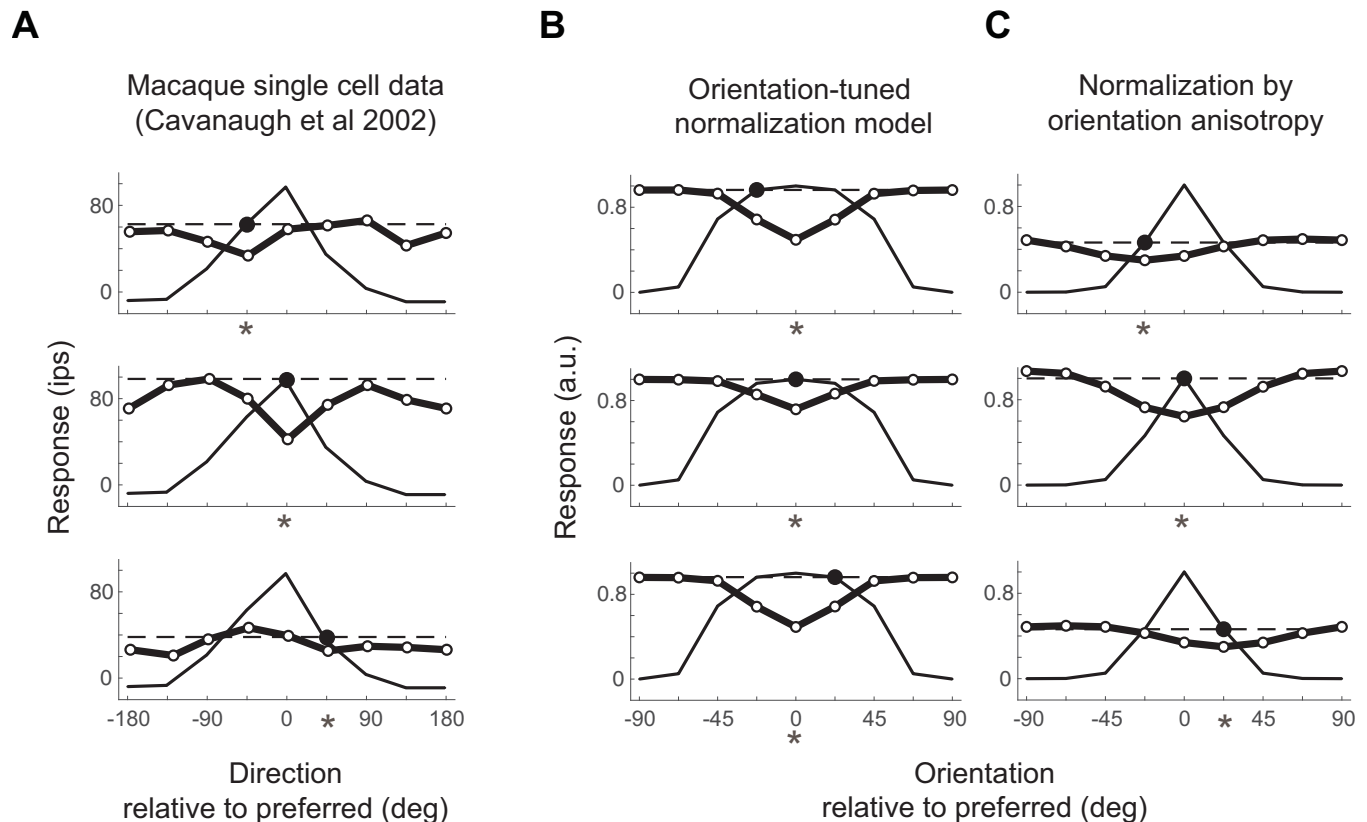


Fig 10. Center surround interactions in single units and computational models. (A) Data replotted from three example cells from macaque V1 [35]. The thin lines are tuning curves for drifting gratings in the cell's RF center, with 0 deg indicating the preferred direction. The thick lines and open circles show responses to stimuli including a center and surround. The center direction is indicated by the large black dot, and the surround direction is indicated by the values in the x-axis. The important observation is that the greatest suppression (indicated by the asterisks) occurs when the surround direction matches the center direction, not when the surround direction matches the center's preferred direction. (B) Simulations for the orientation-tuned normalization model for stimuli similar to those in panel A (but static rather than drifting). Because the surround suppression is matched to the center's preferred tuning, the largest suppression occurs at 0 deg, differing from the single unit data. (C) Simulations for the normalization by orientation anisotropy model for the same stimuli as in panel B. Here, the greatest suppression is when the surround orientation matches the center orientation, similar to the single unit data. Data and simulations plotted from the function `s4_visualize('figure 10')` in the [code repository](#).

<https://doi.org/10.1371/journal.pcbi.1011704.g010>

used in the center—suppression was often stronger for a given center stimulus when the parameters of the surround grating matched the parameters of the center grating even when the center grating was not itself of the optimal direction or orientation.” This pattern has been found in multiple studies, indicating that for macaque and cat V1 cells, surround suppression is often maximal when the surround and center stimulus orientation match, independent of the orientation preference of the cell [37,54,55]. Put another way, the tuning of the surround changes as the center orientation or direction changes. As shown by simulation, our normalization by anisotropy model can capture this observation from single units, but our tuned normalization model cannot (Fig 10). In this regard, the Normalization by orientation anisotropy model is more like Coen-Cagli et al's [22] model of single-cell data than is the Orientation-tuned surround model model.

The two models we implemented leave us in an unusual position. We have one model that has slightly higher prediction accuracy and is more in line with standard normalization models (orientation-tuned normalization), and a second model with slightly lower (but still high) prediction accuracy, but a closer fit to some data from single units (normalization by orientation anisotropy). The simplest conclusion is that a contrast energy model incorporating some form

of orientation-dependent normalization greatly outperforms models that do not have tuned normalization. Given the diversity of cells in visual cortex, it is not likely that a single, simplified model will be sufficient to capture the behavior of all cells or cell populations for all stimuli. A neural network model with recurrence, temporal dynamics, and with feedforward, feedback and lateral connections might provide insight into the specific ways in which the surround modulation emerges [54,56], particularly given evidence that normalization changes as the response to a stimuli unfolds [57–60].

Interestingly, the denominator in the normalization by anisotropy model (the normalizer) is almost the same as the numerator of an orientation variance model used to predict the amplitude of gamma oscillations measured with ECoG electrodes in human subjects [13]. The fact that the same term is used in the numerator to predict gamma oscillations and the denominator to predict BOLD is consistent with the empirical observation that many of the stimuli that are most effective for driving large gamma oscillations (high contrast luminance gratings) are relatively ineffective for eliciting BOLD signals [61,62], and multiunit action potentials [63,64]. The link may be that gamma oscillations, rather than being the fundamental mechanism for perception and long-range cortical communication [65,66], are rather a result of the normalization process [13,67]. An open question is whether gamma oscillations might be predicted as accurately, or perhaps more accurately, by the normalization pool used in our orientation-tuned normalization model than the normalization pool of the normalization by anisotropy model.

3.4 Orientation dependent normalization and image statistics

The normalization by orientation anisotropy model and the orientation-tuned model showed larger advantages over untuned models for V2 and V3 compared to V1. This pattern is consistent with findings from electrophysiology and computational modeling of behavior. First, evidence from the non-human primate visual system suggests that the orientation tuning of surround suppression in V1 arises in large part from feedback from extrastriate areas, especially for the more spatially distant effects of surround suppression [55,68,69]. If the tuned suppressive effects depend on computations in extrastriate regions, then these regions may also exhibit more tuned suppression than V1.

More generally, tuned suppression, either in the form of our tuned normalization model or the anisotropy model, reflects sensitivity to higher-order image statistics (correlations over space between filter outputs). Models of extrastriate neural responses, especially V2, also tend to be more sensitive to higher-order image statistics than to contrast energy *per se*, such as models based explicitly on texture statistics [28–30,70,71], or models with multiple subunits that may exhibit heterogeneous feature tuning [72–74].

We consider one specific way that our normalization by anisotropy model might be linked to the V2 models proposed by Simoncelli and colleagues. Suppose the V1 population computes contrast energy localized in orientation and space. (For simplicity, we ignore spatial frequency, as we did here experimentally by band-pass filtering our stimuli.) We then suppose that the V2 cells compute various weighted sums of the V1 outputs. Specifically, we assume that, for each spatial location, the weights among V2 cells form a Fourier basis set on the V1 outputs (across orientation). If these weights are arranged in pairs (similar to the odd and even V1 filters), and the outputs are squared and summed across phase (similar to the V1 energy model), then the summed V2 population output will be proportional to the variance in the V1 response across orientation (David Heeger, personal communication). If this population output is the normalization pool for V2, then we get a normalization term like that in the anisotropy model, which normalizes by the standard deviation across orientation channels. A V2 model based on

these principles would also need a similar term in the numerator, which we do not include. Hence the link is only to the denominator of the model.

3.5 Are typical contrast energy V1 models missing something important?

No model is complete. The standard normalized contrast energy model of V1, when fit with appropriate parameters, can capture substantial variance in V1 responses [75], but not all. Our results indicate that the model fails for at least some relatively simple stimulus classes, and that the failure can be large. But whether a more complex model is needed, such as a gated normalization model [22] or either of the orientation dependent normalization models we implemented, will depend on the stimulus set tested and the purpose of modeling. It is reasonably likely that the computations in both the anisotropy and the tuned normalization models are more connected to computations in extrastriate visual maps than in V1, but may also influence the response in V1 via feedback. [55,76].

More generally, since the development of divisive normalization models in the early 1990s, there has not yet been a generally agreed-upon description of exactly which cell populations contribute to normalization, and with what weights. Some attempts have been made, by assuming efficient coding of natural images [46,77] or by fitting model parameters to neural data [78]. And while there is a large literature on divisive normalization, including its tuning, many issues remain unresolved, including whether normalization within receptive field centers is orientation-tuned (in addition to the extra-classical surround being orientation-tuned) [78]. In this sense, a standard model of V1, with parameters set, that is downloadable and executable on arbitrary input images, does not yet really exist.

Both the anisotropy model and the tuned normalization model we presented make some advance but also have some important limits, most notably the lack of spatial receptive fields, as well as a lack of sensitivity to second-order contrast. Hence, they do not supersede other models, but rather provide compact computational summaries of response patterns that are not well captured by other models. Integrating better computational tools for validation [79], model-based stimulus development [80,81], high-quality standardized data sets [82], and theory [83], may offer the best path toward more complete understanding of neural circuits in visual cortex.

4. Methods

4.1 Ethics statement

Participants provided written informed consent. The experimental protocol was in compliance with the safety guidelines for MRI research and was approved by the University Committee on Activities Involving Human Subjects at New York University (IRB-FY2016-363).

We analyzed and modeled four fMRI data sets for this paper. Data sets 1 and 2 were collected at NYU. Data sets 3 and 4 are re-analyzed from a previous paper [11], for which the fMRI data and stimuli are freely available online (<http://kendrickkay.net/socmodel/>).

4.2 Participants

The two NYU participants were both experienced MRI subjects (female; 24, 29 yo). Data were collected at NYU's Center for Brain Imaging. The experimental protocol was approved by the University Committee on Activities Involving Human Subjects, and informed written consent was obtained from the participants before the study. Both participants had corrected-to-normal vision. The subjects participated in two separate scanning sessions, one for retinotopic mapping and one for the main study on encoding of textures.

4.3 Stimuli

Publicly accessible links to the stimuli, the names of the stimulus classes, and the correspondence between our naming convention and those in the Kay et al., 2013 [11] paper are described in Tables A, B, and C in [S1 Appendix](#). The total spectral power of all images can be visualized using function *visualizeStimulusPowerSpectrum* in the [code repository](#) and is shown for target stimuli from Dataset 2 of Fig A in [S1 Appendix](#).

4.3.1 Stimuli for data sets 3 and 4. Data sets 3 and 4 correspond to subject 1 and subject 2, respectively, in Kay et al., 2013 [11]. The stimuli for the two subjects were the same and are referred to as “Stimulus set 2” on the website (<http://kendrickkay.net/socmodel/>). The publicly available stimulus set includes 156 stimuli, 39 of which were used for this paper (**Tables B and C in S1 Appendix**). The reason we use only a subset of the stimuli is that we modeled how visual areas respond to textures ignoring retinotopic preference, and many of the stimuli used in the original paper varied systematically over space in order to map spatial receptive fields. We used only the large-field textures, i.e., the subset of stimuli whose patterns were similar across the whole circular aperture.

The 39 stimuli are organized into 7 groups, each of which we describe with two terms, one term for the type of texture and one term for the way in which the stimuli within the group vary. For example, GRATINGS (contrast) are stimuli which come from the grating family and vary from low to high contrast. GRATINGS (density) stimuli come from the same family but have uniform contrast and vary in the spacing between the contours. The correspondence between how we refer to the stimuli and how Kay et al., 2013 [11] referred to them is in **Table C in S1 Appendix**. Below we describe the general stimulus characteristics and the 7 specific classes used for this paper. Most of the text is duplicated from Kay et al., 2013 [11] (p. 11), indicated by italics.

General stimulus characteristics. Stimuli were constructed at a resolution of 256 pixels × 256 pixels and were upsampled to 800 pixels × 800 pixels for display purposes. All stimuli were presented within a circular aperture filling the height of the display; the rest of the display was filled with neutral gray. The outer 0.5 deg of the circular aperture was smoothly blended into the background using a half-cosine function.

Stimuli consisted of grayscale images restricted to a band-pass range of spatial frequencies centered at 3 cycles per degree. To enforce this restriction, a custom band-pass filter was used in the generation of some of the stimuli. The filter was a zero-mean isotropic 2D Difference-of-Gaussians filter whose amplitude spectrum peaks at 3 cycles per degree and drops to half-maximum at 1.4 and 4.7 cycles per degree. Restricting the spatial frequency content of the stimuli avoids the complications of building multiscale models and helps constrain the scope of the modeling endeavor. Even with the spatial frequency restriction, it is possible to construct a rich diversity of stimuli including objects and other naturalistic stimuli. ... Each stimulus consisted of nine distinct images that were presented in quick succession. The purpose of this design was to take advantage of the slow dynamics of the BOLD response and average over stimulus dimensions of no interest (e.g., using sinusoidal gratings differing in phase to average over phase).

A key motivating observation for this paper is that the response to gratings was lower than to curved stimuli. Four groups of stimuli, two groups of gratings and two groups of snakes, were studied first and are referred to throughout the paper as *target stimuli*. These are described first.

Target stimuli. SNAKES (contrast, 10 stimuli). Kay et al. refers to these stimuli as noise patterns: *Noise patterns were created by low-pass filtering white noise at a cutoff frequency of 0.5 cycles per degree, thresholding the result, performing edge detection using derivative filters, inverting image polarity such that edges are black, and applying the custom band-pass filter*

(described previously). We generated nine distinct noise patterns and scaled the contrast of the patterns to fill the full luminance range. [The contrast stimuli were then] constructed by varying the contrast of the noise patterns. . . . Ten different contrast levels were used: 1%, 2%, 3%, 4%, 6%, 9%, 14%, 21%, 32%, and 50%. These contrast levels are relative to the contrast of the patterns used in SPACE [not used in this study], which is taken to be 100%.

SNAKES (density, 5 stimuli). These stimuli used the same type of noise patterns as SPACE [not used here] but varied the amount of separation between contours. We generated noise patterns using cutoff frequencies of 2.8, 1.6, 0.9, 0.5, and 0.3 cycles per degree, and numbered these from 1 (smallest separation) to 5 (largest separation). The noise patterns used in SPACE correspond to separation 4; thus, we only constructed stimuli for the remaining separations 1, 2, 3, and 5. The noise patterns occupied the full stimulus extent (no aperture masking)

GRATINGS (contrast, 4 stimuli). These stimuli consisted of horizontal sinusoidal gratings at 2%, 4%, 9%, and 20% Michelson contrast. The spatial frequency of the gratings was fixed at 3 cycles per degree.

GRATINGS (density, 5 stimuli). The highest density stimulus in this group is similar to the horizontally oriented stimulus in **GRATINGS** (orientation), i.e., similar to a horizontal high-contrast grating, but it is not precisely a sinusoidal grating. It is made by convolving equally spaced horizontal lines with the custom band-pass filter (described previously). When the gratings are spaced appropriately ($\frac{1}{3}$ deg spacing) and filtered by a band-pass filter centered at 3 cycles per deg, the result is close to a sinusoidal grating at 3 cycles per deg. When the spacing is larger, there are several parallel band-pass contours with uniform gray between them. The spacing between parallel lines for the 5 stimuli varied in powers of 2, as $1/3$ deg \times 1, 2, 4, 8, or 16, from densest to sparsest. Because the grating and snakes stimuli were both constructed by convolving lines with the same band-pass filter, they have some similar properties. They differ in that the lines here were straight whereas the lines used for constructing the snakes stimuli were curved.

Additional stimuli. **GRATINGS** (orientation, 8 stimuli). These stimuli consisted of full-contrast sinusoidal gratings at eight different orientations. The spatial frequency of the gratings was fixed at 3 cycles per degree. Each [of the 9 exemplars per stimulus] consisted of gratings with the same orientation but nine different phases (equally spaced from 0 to 2π).

PLAID (contrast, 4 stimuli). These stimuli consisted of plaids at 2%, 4%, 9%, and 20% contrast (defined below). Each condition comprised nine plaids, and each plaid was constructed as the sum of a horizontal and a vertical sinusoidal grating (spatial frequency 3 cycles per degree, random phase). The plaids were scaled in contrast to match the root-mean-square (RMS) contrast of the **GRATING** stimuli. For example, the plaids in the 9% condition were scaled such that the average RMS contrast of the plaids is identical to the average RMS contrast of the gratings in the 9% **GRATING** stimulus.

CIRCULAR (contrast, 4 stimuli). These stimuli were identical to the **PLAID** stimuli except that sixteen different orientations were used instead of two.

4.3.2 Stimuli for Data set 1. Data set 1 was collected at NYU. The data set was designed to replicate some of the effects observed from data sets 3 and 4 (the greater response to snakes than gratings), but also to extend the measurements to new stimulus classes. The general stimulus characteristics were the same as those used in data sets 3 and 4. However, because the display size differed, the image resolution in pixels also differed (400×400 here, vs 800×800 above), and there were slight differences in the bandpass filter. The stimulus size in degrees of visual angle was the same (12.5 deg diameter). A total of 50 stimuli were tested. (The numbers below total more than 50 because some stimuli belong to more than one group, as indicated in **Table C in S1 Appendix**).

Target stimuli. **GRATINGS** (contrast, 5 stimuli). These stimuli are horizontal gratings (but not quite sinusoids), with a similar spatial pattern to the middle stimulus in the

GRATINGS (density) stimuli from data sets 3 and 4. They were made by convolving horizontal lines spaced every 1.75 deg with a custom band-pass filter. The images were scaled to yield 5 different contrasts of 3%, 10%, 25%, 50% and 100%.

GRATINGS (density; 5 stimuli). These stimuli are similar to horizontal gratings, made by convolving equally spaced horizontal lines with the custom band-pass filter. The 5 stimuli differed in the spacing of the horizontal lines, spaced every 3, 2.5, 1.75, 1, 0.33 deg. The contrast of all stimuli was 25%. The middle stimulus in this sequence was the same as the middle stimulus in the contrast sequence (spacing 1.75 deg, 25% contrast). The highest density (0.33 deg spacing) is close to a sinusoidal grating, as the spacing is the inverse of the peak spatial frequency of the band-pass filter (3 cycles per degree).

SNAKES (contrast, 5 stimuli). The spatial pattern is similar to the snakes stimuli in data set 3 and data set 4. Contrasts matched the grating contrasts (3%, 10%, 25%, 50% and 100%).

SNAKES (density, 5 stimuli). The spatial pattern is similar to the snakes stimuli in data set 3 and data set 4, but with 5 different densities of the contours. The contrast for all stimuli was 25% (lower than the contrast of the corresponding stimuli in data set 3 and data set 4). The range of densities used here was also lower than the range used in data set 3 and data set 4, with the densest pattern here similar to the middle stimuli in data set 3 and data set 4.

Additional stimuli. GRATINGS (orientation, 4 stimuli). These stimuli are the same as the third stimulus in the GRATINGS (density) group (25% contrast, 1.75 deg spacing between contours), except that they are rotated by 0, 45, 90, or 135 deg. Because the 0 deg rotation does not change the image, a new stimulus was not created; for visualization of results, the BOLD measurements and model predictions are plotted for both groups.

GRATINGS (cross, 4 stimuli). These stimuli contain horizontal contours similar to two of the stimuli in the GRATINGS (density) sequence, except that they have periodic vertical blank regions which interrupt the contours. For two of the stimuli, the spacing of the horizontal contours matches the densest stimuli in the density sequence (spacing of 0.33 deg) and for two of the stimuli, the spacing matched the middle stimulus in the density sequence (1.75 deg spacing). In all 4 images, the horizontal contours are interrupted by vertical blanks spaced every 1.75 deg. The vertical blanks are either thick (50% duty cycle; 1st and 3rd stimulus) or thin (25% duty cycle; 2nd and 4th stimuli).

NOISE BARS (density, 5 stimuli). These stimuli have the same contrast apertures as the GRATINGS (density) stimuli. Specifically, there are horizontal bands containing contrast patterns, spaced the same as the grating stimuli (bands every 3, 2.5, 1.75, 1, or 0.33 deg). These stimuli differ from the gratings in that each band contains band-pass filtered noise, equal in power across orientations, rather than horizontal contours.

NOISE BARS (contrast, 5 stimuli). These stimuli are matched in spatial pattern to the middle density of the NOISE BARS (density) stimuli (horizontal lines, spacing 1.75 deg), but scaled in contrast similar to the grating stimuli (3%, 10%, 25%, 50%, 100%).

NOISE BARS (orientation, 4 stimuli). The orientation sequence rotated the middle stimulus of the NOISE BARS (density) group (spacing 1.75 deg, contrast 25%) by 0, 45, 90, or 135 deg.

WAVES (density, 6 stimuli). These are identical to the snakes (density) stimuli, except that they have been filtered by orientation, such that they only contain power at or near the horizontal.

WAVES (contrast, 5 stimuli). These are identical to the snakes (contrast) stimuli, except that they have been filtered by orientation, such that they only contain power at or near the horizontal.

WAVES (orientation, 4 stimuli). These are identical to the densest stimulus in the snakes (density) group, except that they have been filtered by orientation, with filter centered at either 0, 45, 90, or 135 deg.

4.3.3 Stimuli for data set 2. Data set 2 was collected at NYU. The stimuli were nearly identical to those in data set 1, differing only in the following ways. First, the stimuli were 50% larger (18.75×18.75 deg and 600×600 pixels, rather than 12.5×12.5 deg and 400×400 pixels). The difference in size did not entail a difference in spatial frequency: The spatial frequency was matched between the two data sets (meaning that the stimuli were re-made with a larger aperture rather than by re-scaling). Second, those stimuli which were oriented were oriented vertically rather than horizontally. This applies to all GRATING stimuli, as well as NOISE BARS and WAVES. Third, the WAVES (density) stimuli had only 4 densities rather than 6. We reduced the number of stimuli to slightly shorten the MRI scans.

4.4 MRI

The methods for MRI acquisition and preprocessing for data sets 3 and 4 are described in Kay et al., 2013 [11]. In brief, each data set comes from one subject, who viewed a variety of stimuli in an event-related fMRI design. Data set 3 was collected over two scan sessions and each stimulus was presented 6 times. Data set 4 was collected over one scan session and each stimulus was presented 3 times. (Note that both in this paper and the Kay website (<http://kendrickkay.net/socmodel/>), these two data sets are referred to as data sets 3 and 4. However, the Kay website refers to the stimuli for these two data sets as stimulus set 2 and the subjects themselves as “subject B” and “subject C”. We do not adopt these latter two conventions.)

After preprocessing the data (slice-time correction, co-registration, spatial unwarping), a general linear model was applied using the GLMdenoise toolbox [11]. The output of this algorithm includes a coefficient (beta weight) for each stimulus for each voxel solved from the whole fMRI session, as well as 30 bootstrapped estimates of each beta weight (bootstrapping across fMRI runs). The publicly available data (<http://kendrickkay.net/socmodel/>) are already pre-processed, denoised, and organized by ROI. Specifically, the data we used are in the files called “data set03.mat” and “data set04.mat” (http://kendrickkay.net/socmodel/data_set03.mat, http://kendrickkay.net/socmodel/data_set04.mat). The data sets are described on the website as “Data set 3 (subject B)” and “Data set 4 (subject C),” respectively. Within the MATLAB files, we used the stored 3D array called “betas” (voxels \times stimuli \times bootstraps), limited to V1, V2, V3 as indicated in the grouping variables “roi” and “roilabels”, and limited to the 39 stimuli indicated in Table B in S1 Appendix. Visual areas were identified by retinotopic mapping in a separate session.

4.4.1 Acquisition of data sets 1 and 2. Data sets 1 and 2 were acquired in one scanning session each. Each scanning session had 12 fMRI runs of 249 s each (data set 1) or 241.5 s each (data set 2). For each data set, half of the stimuli were assigned to odd fMRI runs and half to even runs, so that each stimulus was shown 6 times in the session. The stimulus events were 3 s long, consisting of 9 alternations between stimulus exemplar and blank, $\frac{1}{6}$ s each. Trial onsets were every 7.5 s (so 4.5 s blank between trials). To help estimate the hemodynamic response function, there were 12 s of blank at the beginning and end of each run, as well as 5 additional trials randomly interspersed with no stimulus (meaning that 5 times during the scan, trials were separated by 15 s instead of 7.5 s). Thus, each complete run consisted of either (25 stimuli + 5 blanks) \times 7.5 s + 24 s = 249 s (data set 1) or (24 stimuli + 5 blanks) \times 7.5 s + 24 s = 241.5 s (data set 2).

All MRI data were acquired at New York University Center for Brain Imaging using a Siemens Allegra 3T head-only scanner with a Nova Medical phased array, 8-channel receive surface coil (NMSC072). For each participant, we collected functional images (single shot echo planar images, 1500 ms TR, 30 ms TE, and 72° flip angle). Voxels were 2.0 mm^3 isotropic, with 24 slices, with an inplane sampling of 104×80 voxels ($208 \text{ mm A/P} \times 160 \text{ mm L/R}$). The slice

prescription covered most of the occipital lobe, and the posterior part of both the temporal and parietal lobes. Images were corrected for B0 field inhomogeneity using a calibration scan and Center for Brain Imaging algorithms during offline image reconstruction.

We also acquired 1 or 2 T1-weighted whole-brain anatomical scans (MPRAGE sequence; 1mm³), as well as a T1-weighted “inplane” image with the same slice prescription as the functional scans. This scan had an inplane resolution of 1.25 × 1.25 mm and a slice thickness of 2.5 mm, and was collected to aid alignment of the functional images to the high-resolution T1 weighted anatomical images.

In a separate session, retinotopy scans were collected and analyzed using a pRF model as implemented in the Vistasoft software tool (<https://github.com/vistalab/vistasoft>). The methods for acquisition and analysis of the retinotopy data are identical to that described by Zhou et al. 2018 [84].

4.4.2 Data preprocessing and analysis. *Data preprocessing.* Processing of the fMRI data was identical to that described by Zhou et al. 2018 [84]:

We coregistered and segmented the T1 weighted whole-brain anatomical images into gray and white matter voxels using FreeSurfer’s autosegmentation algorithm (<http://surfer.nmr.mgh.harvard.edu>). Using custom software Vistasoft (<https://github.com/vistalab/vistasoft>), the functional data were slice-time corrected by resampling the time series in each slice to the center of each 1.5 s volume. Data were then motion-corrected by coregistering all volumes of all scans to the first volume of the first scan. The first 8 volumes (12 s) of each scan were discarded for analysis to allow longitudinal magnetization and stabilized hemodynamic response.

GLM. The preprocessed fMRI data were then fit by a general linear model, GLMDenoise [11]. This algorithm denoises the data by projecting out nuisance regressors derived in a data-driven manner, and estimates coefficients for each of the 48 or 50 stimuli for each voxel in the functional images. The algorithm bootstraps the data over fMRI runs. For data sets 1 and 2, we generated 50 bootstraps for data set 1 and 100 bootstraps for data set 2. The publicly available data from Kay et al, 2013 [11], included 30 bootstraps per subject. The algorithm also estimated a hemodynamic impulse response function as a finite impulse response function, with 35 time points (52.5 s) per subject.

ROIs. Regions of interest for V1, V2, V3 were delineated manually using the Vistasoft (<https://github.com/vistalab/vistasoft>) graphical user interface to visualize the results of the pRF models. These methods for identifying these boundaries are well established, as described in many publications [85,86, summarized by 87]. The ROIs for V1, V2 and V3 were identified on the cortical surface and then projected to the functional images. For purposes of data summary and model fitting, we took the average signal from each ROI. We did this by averaging the beta weight across voxels within an ROI separately for each stimulus, after voxel selection (Table A in S1 Appendix). Because noise can be correlated across voxels, but should not be correlated across scans, when we bootstrapped the data, we average across voxels within an ROI for each bootstrap. For the purposes of model fitting, each of the 4 data sets comprised two matrices, one for the means and one for the standard deviation across bootstraps, each of which had a size equal to the number of stimuli by number of ROIs.

4.5 Model equations

In the Results, we compared the accuracy of four models fit to the data, three of which are based on existing models or empirical findings—a contrast energy model, a untuned

normalization model, and an orientation-tuned normalization model—and one new model, which computes normalization by orientation anisotropy. In this section we describe the computation that comprises each model.

All four models consist of three primary steps: (1) computation of oriented contrast energy, (2) pooling across orientation and space, and (3) a power-law nonlinearity. Steps 1 and 3 are identical for all models. Step 2, spatial pooling, varies between models.

1. Contrast energy. We denote by $I(x,y)$ the value of the pre-processed input image at coordinates (x,y) . The pre-processing causes the image values to have mean 0 and range from -0.5 to 0.5. The image is projected onto a set of 128 Gabor filters, which comprise 8 orientations θ , spaced every 22.5 deg; 8 spatial frequencies f , with peak spatial frequency log spaced from 0.75 cpd to 6 cpd; and 2 phases ϕ , separated by 90° (i.e., “quadrature”). $F(x,y,\theta,f,\phi)$ indicates the Gabor filter at a spatial location (x,y) , orientation θ , spatial frequency f , and phase ϕ . Each filter comprised a cosine or sine function of 4 cycles, windowed by a Gaussian with SD of 1 cycle. The outputs over the two phases are squared and summed to compute the contrast energy and summed across spatial frequencies. Finally, the contrast energy as a function of spatial position (x,y) and orientation θ becomes

$$E(x, y, \theta) = \sum_{\phi, f} \left(\sum_{x', y'} I(x - x', y - y') F(x', y', \theta, f, \phi) \right)^2 \quad (1)$$

We convolve the image I and the filter F . The computation of contrast energy has no free parameters. Prior to convolution, stimuli were padded with uniform gray (mean luminance) on all sides by the width of the largest filter. After convolution, all energy images were down-sampled to 12 pixels per degree for computational efficiency. We note that for simplicity, we summed over spatial frequency channels with uniform weights. If one were to fit separate parameters for each voxel, then one might expect spatial frequency tuning to vary with eccentricity. Nonetheless, the simplification of uniform weighting is reasonable given that the spatial frequency content of our stimuli is concentrated in a single octave (~2–4 cpd), and fMRI studies of spatial frequency tuning find a wide bandwidth at the voxel level, std of 2.2 octaves, or full width at half max of 5.1 octaves [88].

2. Spatial pooling. Each model differs in how contrast energy is pooled to yield a scalar value, s :

$$s = \Phi[E] \quad (2)$$

where we use square brackets to indicate a function of a function (also called a *functional*). We describe the pooling functional Φ , for each model below.

3. Power-law nonlinearity. Finally, the scalar is passed through a power-law nonlinearity to predict the BOLD amplitude r in units of percent signal change:

$$r = g \cdot s^\alpha \quad (3)$$

Where g is the gain and α is the exponent parameter. These are free parameters fit to the fMRI data. The power-law nonlinearity is similar to divisive normalization in the case where each unit in a population is normalized by the same pool [10].

4.5.1 Pooling functional for contrast energy model. In the contrast energy model, the contrast energy is summed over orientations and space to yield a scalar output, s .

$$s = \frac{1}{N_{\text{ori}} \cdot N_{\text{pixels}}} \sum_{x,y,\theta} E(x, y, \theta) \quad (4)$$

N_{ori} is the number of orientation channels (always 8) and N_{pixels} is the number of pixels per stimulus in the padded images (344^2 , 419^2 , or 342^2). There are no free parameters in this

pooling functional for contrast energy, so the complete contrast energy model has only two free parameters, g and α , both from the power-law nonlinearity step (Eq 3).

4.5.2 Pooling functional for divisive normalization model. The contrast energy in the untuned divisive normalization model is normalized before it is summed: Each (x, y, θ) element in the energy image is normalized by a weighted sum of elements at (x', y', θ') . The weighting is a Gaussian function of distance from location (x, y) and can thus be expressed as a convolution of the contrast energy E , with a Gaussian, G :

$$Z(x, y, \theta) = \sum_{x', y', \theta'} E(x - x', y - y', \theta) G(x', y', \theta'; \theta) \quad (5)$$

The standard deviation of G is 4% of the padded image size, which is approximately 1 deg. G is identical across the 8 orientations.

The normalized contrast energy is

$$d(x, y, \theta) = \frac{E(x, y, \theta)}{\sigma + Z(x, y, \theta)} \quad (6)$$

where σ is a parameter to control the strength of normalization. When σ is large, the normalization is low, and the overall expression approximates the contrast energy model. When σ approaches 0, there is strong normalization. We then sum d across space and orientation to result in the scalar, s .

$$s = \frac{1}{N_{\text{ori}} \cdot N_{\text{pixels}}} \sum_{x, y, \theta} d(x, y, \theta) \quad (7)$$

The pooling functional for divisive normalization introduces one free parameter, σ . As with the contrast energy model, a power-law nonlinearity is applied to s to predict the BOLD response in percent signal change (Eq 3). Hence the complete model has three free parameters.

We note that the complete divisive normalization model has two similar non-linearities, one in the pooling functional (divisive normalization) and one on the final output (power-law). This is consistent with prior work showing that two stages of normalization (a cascade model) improved model accuracy [11].

4.5.3 Pooling functional for orientation-tuned normalization model. We implement the orientation-tuned normalization (OTN) model identically to the divisive normalization model (Eqs 5–7), except that the contrast energy normalizer $Z(x, y, \theta)$ is now orientation-tuned. When the orientation channel of the image and filter matches, $\theta' = \theta$, the 2-D filter of the channel is a 2D Gaussian identical to the untuned normalization (4.4.2). At all other orientations (i.e., $\theta' \neq \theta$), the filter is a symmetric 2D Gaussian distribution with a much small standard deviation, effectively just one pixel (Fig 4). This is akin to summing two forms of normalization, cross-orientation suppression (same location, other orientations), and an orientation-tuned surround (same orientation, other locations). As with the untuned normalization model, the pooling functional introduces only one free parameter, σ . The complete model, including the power-law non-linearity (Eq 3) has 3 free parameters.

4.5.4 Pooling functional for normalization by orientation anisotropy. In the normalization by orientation anisotropy (NOA) model, the pooling step first sums the contrast energy

across space within an orientation band, resulting in one value per orientation band, $E_{\text{ori}}(\theta)$:

$$E_{\text{ori}}(\theta) = \frac{1}{N_{\text{pixels}}} \sum_{x,y} E(x, y, \theta) \quad (8)$$

E_{ori} indicates the oriented energy. This energy at each orientation is then normalized by the standard deviation across the 8 orientations, and then summed to produce a scalar.

$$s = \frac{1}{N_{\text{ori}}} \sum_{\theta} \frac{E_{\text{ori}}(\theta)}{\sigma + \text{Std}(E_{\text{ori}})} \quad (9)$$

Where $\text{Std}(E_{\text{ori}}) = \sqrt{\sum_{\theta} (E_{\text{ori}}(\theta) - \bar{E}_{\text{ori}})^2}$ calculates the standard deviation of the oriented energy and a non-negative parameter w controls the strength of the normalization. When σ is large, the normalization is low, and the overall expression approximates the contrast energy model. When σ approaches 0, there is strong normalization by the standard deviation across orientation channel outputs. Calculating the standard deviation of oriented energy involves a squaring operation. To keep the parameters comparable across different models, we also square the numerator and the parameter σ . The NOA pooling functional has one free parameter, σ . The complete model, including the power-law non-linearity (Eq 3) has 3 free parameters.

4.6 Optimization

In each model, we fitted the model free parameters using the MATLAB optimization tool *fmincon* by minimizing the squared error between the model prediction and the corresponding BOLD amplitude. Because each stimulus consisted of 9 exemplars shown to the subject in rapid succession, the model prediction for each stimulus was obtained by averaging the model predictions across the exemplars. To avoid getting stuck in the local minima of the nonconvex landscape, we ran the optimization algorithm with 40 different parameter initializations. Each initialized value was picked randomly. All parameters were unbounded in the search, minimizing human interference in the fitting. Parameter α is passed through a sigmoid function to ensure its value is between 0 and 1.

4.7 Cross-validation scheme

All models were fit using an leave-one-out cross-validation scheme, where n is the number of stimuli. Thus, the BOLD signal prediction for each stimulus was generated by a model fit to all stimuli except that one. Under this scheme, the models are less likely to overfit data sets.

4.8 Accuracy metric

The model accuracy was quantified as the percentage of the explained variance (R^2) in the human BOLD data by the cross-validated model predictions,

$$R^2 = 1 - \frac{\sum_{i=1}^M (r_i - \hat{r}_i)^2}{\sum_{i=1}^M (r_i - \bar{r})^2} \quad (10)$$

where r_i represents the BOLD amplitude to the i^{th} stimulus, \hat{r}_i represents the corresponding model prediction, and \bar{r} is the mean response across stimuli. We can understand this metric as the extra uncertainty reduction brought by the model beyond describing the BOLD data by its mean.

Supporting information

S1 Appendix. Data Set and Stimulus Properties.
(PDF)

S2 Appendix. Stimulus Images.
(PDF)

S3 Appendix. Model Variance Explained.
(PDF)

S4 Appendix. Model Parameter Estimates.
(PDF)

S5 Appendix. Plots of Model Fits.
(PDF)

Author Contributions

Conceptualization: Zeming Fang, Ilona M. Bloem, Catherine Olsson, Wei Ji Ma, Jonathan Winawer.

Data curation: Zeming Fang, Ilona M. Bloem, Catherine Olsson, Jonathan Winawer.

Formal analysis: Zeming Fang, Ilona M. Bloem, Catherine Olsson, Wei Ji Ma, Jonathan Winawer.

Funding acquisition: Jonathan Winawer.

Investigation: Zeming Fang, Ilona M. Bloem, Catherine Olsson, Jonathan Winawer.

Methodology: Zeming Fang, Ilona M. Bloem, Wei Ji Ma, Jonathan Winawer.

Project administration: Wei Ji Ma, Jonathan Winawer.

Resources: Jonathan Winawer.

Software: Zeming Fang, Ilona M. Bloem, Jonathan Winawer.

Supervision: Wei Ji Ma, Jonathan Winawer.

Validation: Zeming Fang, Ilona M. Bloem, Jonathan Winawer.

Visualization: Zeming Fang, Ilona M. Bloem, Jonathan Winawer.

Writing – original draft: Zeming Fang, Jonathan Winawer.

Writing – review & editing: Zeming Fang, Ilona M. Bloem, Wei Ji Ma, Jonathan Winawer.

References

1. Heeger DJ. Modeling simple-cell direction selectivity with normalized, half-squared, linear operators. *J Neurophysiol.* 1993; 70(5):1885–98. <https://doi.org/10.1152/jn.1993.70.5.1885> PMID: 8294961.
2. Wandell BA. Foundations of vision. Sunderland, Mass.: Sinauer Associates; 1995. xvi, 476 p., [4] p. of plates p.
3. Hubel DH, Wiesel TN. Receptive fields, binocular interaction and functional architecture in the cat's visual cortex. *J Physiol.* 1962; 160:106–54. Epub 1962/01/01. <https://doi.org/10.1113/jphysiol.1962.sp006837> PMID: 14449617; PubMed Central PMCID: PMC1359523.
4. Adelson EH, Bergen JR. Spatiotemporal energy models for the perception of motion. *J Opt Soc Am A.* 1985; 2(2):284–99. <https://doi.org/10.1364/josaa.2.000284> PMID: 3973762.
5. Pollen DA, Ronner SF. Visual cortical neurons as localized spatial frequency filters. *IEEE Transactions on Systems, Man, and Cybernetics.* 1983;(5):907–16.

6. Heeger DJ. Half-squaring in responses of cat striate cells. *Visual neuroscience*. 1992; 9(5):427–43. <https://doi.org/10.1017/s095252380001124x> PMID: 1450099.
7. Heeger DJ. Normalization of cell responses in cat striate cortex. *Visual neuroscience*. 1992; 9(2):181–97. <https://doi.org/10.1017/s0952523800009640> PMID: 1504027.
8. Carandini M, Demb JB, Mante V, Tolhurst DJ, Dan Y, Olshausen BA, et al. Do we know what the early visual system does? *The Journal of neuroscience: the official journal of the Society for Neuroscience*. 2005; 25(46):10577–97. <https://doi.org/10.1523/JNEUROSCI.3726-05.2005> PMID: 16291931.
9. Kay KN, Naselaris T, Prenger RJ, Gallant JL. Identifying natural images from human brain activity. *Nature*. 2008; 452(7185):352–5. <https://doi.org/10.1038/nature06713> PMID: 18322462; PubMed Central PMCID: PMC3556484.
10. Kay KN, Winawer J, Mezer A, Wandell BA. Compressive spatial summation in human visual cortex. *J Neurophysiol*. 2013; 110(2):481–94. <https://doi.org/10.1152/jn.00105.2013> PMID: 23615546; PubMed Central PMCID: PMC3727075.
11. Kay KN, Winawer J, Rokem A, Mezer A, Wandell BA. A two-stage cascade model of BOLD responses in human visual cortex. *PLoS computational biology*. 2013; 9(5):e1003079. <https://doi.org/10.1371/journal.pcbi.1003079> PMID: 23737741; PubMed Central PMCID: PMC3667759.
12. Brouwer GJ, Heeger DJ. Cross-orientation suppression in human visual cortex. *J Neurophysiol*. 2011; 106(5):2108–19. <https://doi.org/10.1152/jn.00540.2011> PMID: 21775720; PubMed Central PMCID: PMC3214101.
13. Hermes D, Petridou N, Kay KN, Winawer J. An image-computable model for the stimulus selectivity of gamma oscillations. *Elife*. 2019;8. Epub 2019/11/09. <https://doi.org/10.7554/eLife.47035> PMID: 31702552.
14. Winawer J, Kay KN, Foster BL, Rauschecker AM, Parvizi J, Wandell BA. Asynchronous broadband signals are the principal source of the BOLD response in human visual cortex. *Current biology: CB*. 2013; 23(13):1145–53. <https://doi.org/10.1016/j.cub.2013.05.001> PMID: 23770184; PubMed Central PMCID: PMC3710543.
15. David SV, Vinje WE, Gallant JL. Natural stimulus statistics alter the receptive field structure of v1 neurons. *The Journal of neuroscience: the official journal of the Society for Neuroscience*. 2004; 24(31):6991–7006. Epub 2004/08/06. <https://doi.org/10.1523/JNEUROSCI.1422-04.2004> PMID: 15295035; PubMed Central PMCID: PMC6729594.
16. Vinje WE, Gallant JL. Sparse coding and decorrelation in primary visual cortex during natural vision. *Science*. 2000; 287(5456):1273–6. Epub 2000/02/26. <https://doi.org/10.1126/science.287.5456.1273> PMID: 10678835.
17. Olshausen BA, Field DJ. How close are we to understanding v1? *Neural Comput*. 2005; 17(8):1665–99. <https://doi.org/10.1162/0899766054026639> PMID: 15969914.
18. Rust NC, Movshon JA. In praise of artifice. *Nat Neurosci*. 2005; 8(12):1647–50. Epub 2005/11/25. <https://doi.org/10.1038/nn1606> PMID: 16306892.
19. Hubel DH, Wiesel TN. Receptive fields and functional architecture of monkey striate cortex. *J Physiol*. 1968; 195(1):215–43. <https://doi.org/10.1113/jphysiol.1968.sp008455> PMID: 4966457; PubMed Central PMCID: PMC1557912.
20. Hubel DH, Wiesel TN. Receptive Fields and Functional Architecture in Two Nonstriate Visual Areas (18 and 19) of the Cat. *J Neurophysiol*. 1965; 28:229–89. <https://doi.org/10.1152/jn.1965.28.2.229> PMID: 14283058.
21. Tang S, Lee TS, Li M, Zhang Y, Xu Y, Liu F, et al. Complex Pattern Selectivity in Macaque Primary Visual Cortex Revealed by Large-Scale Two-Photon Imaging. *Current biology: CB*. 2018; 28(1):38–48 e3. Epub 2017/12/19. <https://doi.org/10.1016/j.cub.2017.11.039> PMID: 29249660.
22. Coen-Cagli R, Kohn A, Schwartz O. Flexible gating of contextual influences in natural vision. *Nat Neurosci*. 2015; 18(11):1648–55. <https://doi.org/10.1038/nn.4128> PMID: 26436902; PubMed Central PMCID: PMC4624479.
23. Walker EY, Sinz FH, Cobos E, Muhammad T, Froudarakis E, Fahey PG, et al. Inception loops discover what excites neurons most using deep predictive models. *Nat Neurosci*. 2019; 22(12):2060–5. Epub 2019/11/04. <https://doi.org/10.1038/s41593-019-0517-x> PMID: 31686023.
24. Mannion DJ, Kersten DJ, Olman CA. Scene coherence can affect the local response to natural images in human V1. *Eur J Neurosci*. 2015; 42(11):2895–903. Epub 2015/09/24. <https://doi.org/10.1111/ejn.13082> PMID: 26390850; PubMed Central PMCID: PMC4715660.
25. Qiu C, Burton PC, Kersten D, Olman CA. Responses in early visual areas to contour integration are context dependent. *J Vis*. 2016; 16(8):19. Epub 2016/07/02. <https://doi.org/10.1167/16.8.19> PMID: 27366994; PubMed Central PMCID: PMC4946811.

26. Bloem IM, Ling S. Normalization governs attentional modulation within human visual cortex. *Nat Commun*. 2019; 10(1):5660. Epub 2019/12/13. <https://doi.org/10.1038/s41467-019-13597-1> PMID: [31827078](#); PubMed Central PMCID: PMC6906520.
27. Klimova M, Bloem IM, Ling S. The specificity of orientation-tuned normalization within human early visual cortex. *J Neurophysiol*. 2021; 126(5):1536–46. Epub 2021/09/23. <https://doi.org/10.1152/jn.00203.2021> PMID: [34550028](#); PubMed Central PMCID: PMC8794056.
28. Freeman J, Ziemba CM, Heeger DJ, Simoncelli EP, Movshon JA. A functional and perceptual signature of the second visual area in primates. *Nat Neurosci*. 2013; 16(7):974–81. <https://doi.org/10.1038/nn.3402> PMID: [23685719](#); PubMed Central PMCID: PMC3710454.
29. Movshon JA, Simoncelli EP. Representation of Naturalistic Image Structure in the Primate Visual Cortex. *Cold Spring Harb Symp Quant Biol*. 2015. <https://doi.org/10.1101/sqb.2014.79.024844> PMID: [25943766](#).
30. Okazawa G, Tajima S, Komatsu H. Gradual Development of Visual Texture-Selective Properties Between Macaque Areas V2 and V4. *Cereb Cortex*. 2017; 27(10):4867–80. Epub 2016/09/23. <https://doi.org/10.1093/cercor/bhw282> PMID: [27655929](#).
31. Sanchez-Giraldo LG, Laskar MNU, Schwartz O. Normalization and pooling in hierarchical models of natural images. *Curr Opin Neurobiol*. 2019; 55:65–72. Epub 2019/02/18. <https://doi.org/10.1016/j.conb.2019.01.008> PMID: [30785005](#).
32. Cadieu CF, Hong H, Yamins DL, Pinto N, Ardila D, Solomon EA, et al. Deep neural networks rival the representation of primate IT cortex for core visual object recognition. *PLoS computational biology*. 2014; 10(12):e1003963. Epub 2014/12/18. <https://doi.org/10.1371/journal.pcbi.1003963> PMID: [25521294](#); PubMed Central PMCID: PMC4270441.
33. Yamins DL, Hong H, Cadieu CF, Solomon EA, Seibert D, DiCarlo JJ. Performance-optimized hierarchical models predict neural responses in higher visual cortex. *Proc Natl Acad Sci U S A*. 2014; 111(23):8619–24. Epub 2014/05/08. <https://doi.org/10.1073/pnas.1403112111> PMID: [24812127](#); PubMed Central PMCID: PMC4060707.
34. Zhuang C, Yan S, Nayebi A, Schrimpf M, Frank MC, DiCarlo JJ, et al. Unsupervised neural network models of the ventral visual stream. *Proc Natl Acad Sci U S A*. 2021; 118(3). <https://doi.org/10.1073/pnas.2014196118> PMID: [33431673](#); PubMed Central PMCID: PMC7826371.
35. Cavanaugh JR, Bair W, Movshon JA. Selectivity and spatial distribution of signals from the receptive field surround in macaque V1 neurons. *J Neurophysiol*. 2002; 88(5):2547–56. <https://doi.org/10.1152/jn.00693.2001> PMID: [12424293](#).
36. DeAngelis GC, Freeman RD, Ohzawa I. Length and width tuning of neurons in the cat's primary visual cortex. *J Neurophysiol*. 1994; 71(1):347–74. <https://doi.org/10.1152/jn.1994.71.1.347> PMID: [8158236](#).
37. Sillito AM, Grieve KL, Jones HE, Cudeiro J, Davis J. Visual cortical mechanisms detecting focal orientation discontinuities. *Nature*. 1995; 378(6556):492–6. <https://doi.org/10.1038/378492a0> PMID: [7477405](#).
38. Shushruth S, Nurminen L, Bijanzadeh M, Ichida JM, Vanni S, Angelucci A. Different orientation tuning of near- and far-surround suppression in macaque primary visual cortex mirrors their tuning in human perception. *The Journal of neuroscience: the official journal of the Society for Neuroscience*. 2013; 33(1):106–19. <https://doi.org/10.1523/JNEUROSCI.2518-12.2013> PMID: [23283326](#); PubMed Central PMCID: PMC3711542.
39. Trott AR, Born RT. Input-gain control produces feature-specific surround suppression. *The Journal of neuroscience: the official journal of the Society for Neuroscience*. 2015; 35(12):4973–82. <https://doi.org/10.1523/JNEUROSCI.4000-14.2015> PMID: [25810527](#); PubMed Central PMCID: PMC4389596.
40. Xing J, Heeger DJ. Measurement and modeling of center-surround suppression and enhancement. *Vision Res*. 2001; 41(5):571–83. [https://doi.org/10.1016/s0042-6989\(00\)00270-4](https://doi.org/10.1016/s0042-6989(00)00270-4) PMID: [11226503](#).
41. Wang HX, Heeger DJ, Landy MS. Responses to second-order texture modulations undergo surround suppression. *Vision Res*. 2012; 62:192–200. <https://doi.org/10.1016/j.visres.2012.03.008> PMID: [22811987](#); PubMed Central PMCID: PMC3477815.
42. Solomon JA, Sperling G, Chubb C. The lateral inhibition of perceived contrast is indifferent to on-center/off-center segregation, but specific to orientation. *Vision Res*. 1993; 33(18):2671–83. Epub 1993/12/01. [https://doi.org/10.1016/0042-6989\(93\)90227-n](https://doi.org/10.1016/0042-6989(93)90227-n) PMID: [8296464](#).
43. Zenger-Landolt B, Heeger DJ. Response suppression in v1 agrees with psychophysics of surround masking. *The Journal of neuroscience: the official journal of the Society for Neuroscience*. 2003; 23(17):6884–93. <https://doi.org/10.1523/JNEUROSCI.23-17-06884.2003> PMID: [12890783](#); PubMed Central PMCID: PMC2275204.
44. McDonald JS, Seymour KJ, Schira MM, Spehar B, Clifford CW. Orientation-specific contextual modulation of the fMRI BOLD response to luminance and chromatic gratings in human visual cortex. *Vision Res*. 2009; 49(11):1397–405. Epub 2009/01/22. <https://doi.org/10.1016/j.visres.2008.12.014> PMID: [19167419](#).

45. Joo SJ, Boynton GM, Murray SO. Long-range, pattern-dependent contextual effects in early human visual cortex. *Current biology: CB*. 2012; 22(9):781–6. Epub 20120412. <https://doi.org/10.1016/j.cub.2012.02.067> PMID: 22503498; PubMed Central PMCID: PMC3350565.
46. Schwartz O, Simoncelli EP. Natural signal statistics and sensory gain control. *Nat Neurosci*. 2001; 4(8):819–25. <https://doi.org/10.1038/90526> PMID: 11477428.
47. Morrone MC, Burr DC, Maffei L. Functional implications of cross-orientation inhibition of cortical visual cells. I. Neurophysiological evidence. *Proc R Soc Lond B Biol Sci*. 1982; 216(1204):335–54. Epub 1982/10/22. <https://doi.org/10.1098/rspb.1982.0078> PMID: 6129633.
48. Vinck M, Bosman CA. More Gamma More Predictions: Gamma-Synchronization as a Key Mechanism for Efficient Integration of Classical Receptive Field Inputs with Surround Predictions. *Front Syst Neurosci*. 2016; 10:35. Epub 2016/05/21. <https://doi.org/10.3389/fnsys.2016.00035> PMID: 27199684; PubMed Central PMCID: PMC4842768.
49. Sceniak MP, Ringach DL, Hawken MJ, Shapley R. Contrast's effect on spatial summation by macaque V1 neurons. *Nat Neurosci*. 1999; 2(8):733–9. <https://doi.org/10.1038/11197> PMID: 10412063.
50. Cavanaugh JR, Bair W, Movshon JA. Nature and interaction of signals from the receptive field center and surround in macaque V1 neurons. *J Neurophysiol*. 2002; 88(5):2530–46. <https://doi.org/10.1152/jn.00692.2001> PMID: 12424292.
51. Li CY, Creutzfeldt O. The representation of contrast and other stimulus parameters by single neurons in area 17 of the cat. *Pflügers Arch*. 1984; 401(3):304–14. <https://doi.org/10.1007/BF00582601> PMID: 6473083.
52. Albrecht DG, Hamilton DB. Striate cortex of monkey and cat: contrast response function. *J Neurophysiol*. 1982; 48(1):217–37. Epub 1982/07/01. <https://doi.org/10.1152/jn.1982.48.1.217> PMID: 7119846.
53. Sclar G, Freeman RD. Orientation selectivity in the cat's striate cortex is invariant with stimulus contrast. *Exp Brain Res*. 1982; 46(3):457–61. <https://doi.org/10.1007/BF00238641> PMID: 7095050.
54. Shushruth S, Mangapathy P, Ichida JM, Bressloff PC, Schwabe L, Angelucci A. Strong recurrent networks compute the orientation tuning of surround modulation in the primate primary visual cortex. *The Journal of neuroscience: the official journal of the Society for Neuroscience*. 2012; 32(1):308–21. <https://doi.org/10.1523/JNEUROSCI.3789-11.2012> PMID: 22219292; PubMed Central PMCID: PMC3711470.
55. Angelucci A, Bijanzadeh M, Nurminen L, Federer F, Merlin S, Bressloff PC. Circuits and Mechanisms for Surround Modulation in Visual Cortex. *Annu Rev Neurosci*. 2017; 40:425–51. Epub 20170503. <https://doi.org/10.1146/annurev-neuro-072116-031418> PMID: 28471714; PubMed Central PMCID: PMC5697758.
56. Heeger DJ, Zemlianova KO. A recurrent circuit implements normalization, simulating the dynamics of V1 activity. *Proc Natl Acad Sci U S A*. 2020. Epub 2020/08/28. <https://doi.org/10.1073/pnas.2005417117> PMID: 32843341.
57. Webb BS, Dhruv NT, Solomon SG, Tailby C, Lennie P. Early and late mechanisms of surround suppression in striate cortex of macaque. *The Journal of neuroscience: the official journal of the Society for Neuroscience*. 2005; 25(50):11666–75. <https://doi.org/10.1523/JNEUROSCI.3414-05.2005> PMID: 16354925; PubMed Central PMCID: PMC6726034.
58. Schallmo MP, Kale AM, Murray SO. The time course of different surround suppression mechanisms. *J Vis*. 2019; 19(4):12. <https://doi.org/10.1167/19.4.12> PMID: 30952163; PubMed Central PMCID: PMC6464404.
59. Petrov Y, McKee SP. The time course of contrast masking reveals two distinct mechanisms of human surround suppression. *J Vis*. 2009; 9(1):21 1–11. Epub 20090120. <https://doi.org/10.1167/9.1.21> PMID: 19271891; PubMed Central PMCID: PMC2842926.
60. Self MW, Lorteije JA, Vangeneugden J, van Beest EH, Grigore ME, Levelt CN, et al. Orientation-tuned surround suppression in mouse visual cortex. *The Journal of neuroscience: the official journal of the Society for Neuroscience*. 2014; 34(28):9290–304. <https://doi.org/10.1523/JNEUROSCI.5051-13.2014> PMID: 25009262; PubMed Central PMCID: PMC6608354.
61. Butler R, Bernier PM, Lefebvre J, Gilbert G, Whittingstall K. Decorrelated Input Dissociates Narrow Band gamma Power and BOLD in Human Visual Cortex. *The Journal of neuroscience: the official journal of the Society for Neuroscience*. 2017; 37(22):5408–18. Epub 2017/04/30. <https://doi.org/10.1523/JNEUROSCI.3938-16.2017> PMID: 28455370.
62. Hermes D, Nguyen M, Winawer J. Neuronal synchrony and the relation between the blood-oxygen-level dependent response and the local field potential. *PLoS Biol*. 2017; 15(7):e2001461. <https://doi.org/10.1371/journal.pbio.2001461> PMID: 28742093.
63. Hermes D, Kasteleijn-Nolst Trenite DGA, Winawer J. Gamma oscillations and photosensitive epilepsy. *Current biology: CB*. 2017; 27(9):R336–R8. <https://doi.org/10.1016/j.cub.2017.03.076> PMID: 28486114; PubMed Central PMCID: PMC5438467.

64. Ray S, Maunsell JH. Different origins of gamma rhythm and high-gamma activity in macaque visual cortex. *PLoS biology*. 2011; 9(4):e1000610. Epub 2011/05/03. <https://doi.org/10.1371/journal.pbio.1000610> PMID: 21532743; PubMed Central PMCID: PMC3075230.
65. Engel AK, Singer W. Temporal binding and the neural correlates of sensory awareness. *Trends Cogn Sci*. 2001; 5(1):16–25. Epub 2001/02/13. [https://doi.org/10.1016/s1364-6613\(00\)01568-0](https://doi.org/10.1016/s1364-6613(00)01568-0) PMID: 11164732.
66. Fries P. A mechanism for cognitive dynamics: neuronal communication through neuronal coherence. *Trends Cogn Sci*. 2005; 9(10):474–80. Epub 2005/09/10. <https://doi.org/10.1016/j.tics.2005.08.011> PMID: 16150631.
67. Ray S, Ni AM, Maunsell JH. Strength of gamma rhythm depends on normalization. *PLoS Biol*. 2013; 11(2):e1001477. Epub 2013/02/09. <https://doi.org/10.1371/journal.pbio.1001477> PMID: 23393427; PubMed Central PMCID: PMC3564761.
68. Angelucci A, Levitt JB, Walton EJ, Hupe JM, Bullier J, Lund JS. Circuits for local and global signal integration in primary visual cortex. *The Journal of neuroscience: the official journal of the Society for Neuroscience*. 2002; 22(19):8633–46. <https://doi.org/10.1523/JNEUROSCI.22-19-08633.2002> PMID: 12351737; PubMed Central PMCID: PMC6757772.
69. Shmuel A, Korman M, Sterkin A, Harel M, Ullman S, Malach R, et al. Retinotopic axis specificity and selective clustering of feedback projections from V2 to V1 in the owl monkey. *The Journal of neuroscience: the official journal of the Society for Neuroscience*. 2005; 25(8):2117–31. <https://doi.org/10.1523/JNEUROSCI.4137-04.2005> PMID: 15728852; PubMed Central PMCID: PMC6726055.
70. Freeman J, Simoncelli EP. Metamers of the ventral stream. *Nat Neurosci*. 2011; 14(9):1195–201. <https://doi.org/10.1038/nn.2889> PMID: 21841776; PubMed Central PMCID: PMC3164938.
71. Ziemba CM, Simoncelli EP. Opposing effects of selectivity and invariance in peripheral vision. *Nat Commun*. 2021; 12(1):4597. Epub 20210728. <https://doi.org/10.1038/s41467-021-24880-5> PMID: 34321483; PubMed Central PMCID: PMC8319169.
72. Liu L, She L, Chen M, Liu T, Lu HD, Dan Y, et al. Spatial structure of neuronal receptive field in awake monkey secondary visual cortex (V2). *Proc Natl Acad Sci U S A*. 2016; 113(7):1913–8. Epub 20160202. <https://doi.org/10.1073/pnas.1525505113> PMID: 26839410; PubMed Central PMCID: PMC4763736.
73. Willmore BD, Prenger RJ, Gallant JL. Neural representation of natural images in visual area V2. *The Journal of neuroscience: the official journal of the Society for Neuroscience*. 2010; 30(6):2102–14. <https://doi.org/10.1523/JNEUROSCI.4099-09.2010> PMID: 20147538; PubMed Central PMCID: PMC2994536.
74. Tao X, Zhang B, Smith EL, 3rd, Nishimoto S, Ohzawa I, Chino YM. Local sensitivity to stimulus orientation and spatial frequency within the receptive fields of neurons in visual area 2 of macaque monkeys. *J Neurophysiol*. 2012; 107(4):1094–110. Epub 20111123. <https://doi.org/10.1152/jn.00640.2011> PMID: 22114163; PubMed Central PMCID: PMC3289454.
75. Carandini M, Heeger DJ, Movshon JA. Linearity and normalization in simple cells of the macaque primary visual cortex. *The Journal of neuroscience: the official journal of the Society for Neuroscience*. 1997; 17(21):8621–44. <https://doi.org/10.1523/JNEUROSCI.17-21-08621.1997> PMID: 9334433.
76. Henry CA, Jazayeri M, Shapley RM, Hawken MJ. Distinct spatiotemporal mechanisms underlie extra-classical receptive field modulation in macaque V1 microcircuits. *Elife*. 2020;9. Epub 20200527. <https://doi.org/10.7554/eLife.54264> PMID: 32458798; PubMed Central PMCID: PMC7253173.
77. Coen-Cagli R, Dayan P, Schwartz O. Cortical Surround Interactions and Perceptual Saliency via Natural Scene Statistics. *PLoS computational biology*. 2012; 8(3):e1002405. Epub 20120301. <https://doi.org/10.1371/journal.pcbi.1002405> PMID: 22396635; PubMed Central PMCID: PMC3291533.
78. Burg MF, Cadena SA, Denfield GH, Walker EY, Tolias AS, Bethge M, et al. Learning divisive normalization in primary visual cortex. *PLoS computational biology*. 2021; 17(6):e1009028. Epub 2021/06/08. <https://doi.org/10.1371/journal.pcbi.1009028> PMID: 34097695; PubMed Central PMCID: PMC8211272.
79. Lerma-Usabiaga G, Benson N, Winawer J, Wandell BA. A validation framework for neuroimaging software: The case of population receptive fields. *PLoS computational biology*. 2020; 16(6):e1007924. Epub 2020/06/26. <https://doi.org/10.1371/journal.pcbi.1007924> PMID: 32584808; PubMed Central PMCID: PMC7343185.
80. Golan T, Raju PC, Kriegeskorte N. Controversial stimuli: Pitting neural networks against each other as models of human cognition. *Proc Natl Acad Sci U S A*. 2020; 117(47):29330–7. Epub 2020/11/25. <https://doi.org/10.1073/pnas.1912334117> PMID: 33229549; PubMed Central PMCID: PMC7703564.
81. Wang Z, Simoncelli EP. Maximum differentiation (MAD) competition: a methodology for comparing computational models of perceptual quantities. *J Vis*. 2008; 8(12):8 1–13. Epub 2008/10/04. <https://doi.org/10.1167/8.12.8> PMID: 18831621; PubMed Central PMCID: PMC4143340.

82. Allen EJ, St-Yves G, Wu Y, Breedlove JL, Prince JS, Dowdle LT, et al. A massive 7T fMRI dataset to bridge cognitive neuroscience and artificial intelligence. *Nat Neurosci*. 2021. Epub 2021/12/18. <https://doi.org/10.1038/s41593-021-00962-x> PMID: 34916659.
83. Simoncelli EP, Olshausen BA. Natural image statistics and neural representation. *Annu Rev Neurosci*. 2001; 24:1193–216. Epub 2001/08/25. <https://doi.org/10.1146/annurev.neuro.24.1.1193> PMID: 11520932.
84. Zhou J, Benson NC, Kay KN, Winawer J. Compressive Temporal Summation in Human Visual Cortex. *The Journal of neuroscience: the official journal of the Society for Neuroscience*. 2018; 38(3):691–709. Epub 2017/11/30. <https://doi.org/10.1523/JNEUROSCI.1724-17.2017> PMID: 29192127; PubMed Central PMCID: PMC5777115.
85. Engel SA, Glover GH, Wandell BA. Retinotopic organization in human visual cortex and the spatial precision of functional MRI. *Cereb Cortex*. 1997; 7(2):181–92. <https://doi.org/10.1093/cercor/7.2.181> PMID: 9087826.
86. Sereno MI, Dale AM, Reppas JB, Kwong KK, Belliveau JW, Brady TJ, et al. Borders of Multiple Visual Areas in Humans Revealed by Functional Magnetic Resonance Imaging. *Science*. 1995; 268(5212):889–93. <https://doi.org/10.1126/science.7754376> PMID: 7754376
87. Wandell BA, Winawer J. Imaging retinotopic maps in the human brain. *Vision Res*. 2011; 51(7):718–37. Epub 2010/08/06. <https://doi.org/10.1016/j.visres.2010.08.004> PMID: 20692278; PubMed Central PMCID: PMC3030662.
88. Broderick WF, Simoncelli EP, Winawer J. Mapping spatial frequency preferences across human primary visual cortex. *J Vis*. 2022; 22(4):3. Epub 2022/03/11. <https://doi.org/10.1167/jov.22.4.3> PMID: 35266962; PubMed Central PMCID: PMC8934567.