

A Model for Image Segmentation in Retina

Christopher Warner^{1,2} and Friedrich T. Sommer^{1,3,4}

¹Redwood Center for Theoretical Neuroscience

²Biophysics Graduate Group

University of California, Berkeley

Berkeley, CA 94720

³Helen Wills Neuroscience Institute

University of California, Berkeley

Berkeley, CA 94720

⁴Intel Labs, Santa Clara, CA 95054-1549

cwarner@berkeley.edu, fsommer@berkeley.edu

May 7, 2020

Abstract

While traditional filter models can reproduce the rate responses of retinal ganglion neurons to simple stimuli, they cannot explain why synchrony between spikes is much higher than expected by Poisson firing [6], and can be sometimes rhythmic [25, 16]. Here we investigate the hypothesis that synchrony in periodic retinal spike trains could convey contextual information of the visual input, which is extracted by computations in the retinal network. We propose a computational model for image segmentation consisting of a Kuramoto model of coupled oscillators whose phases model the timing of individual retinal spikes. The phase couplings between oscillators are shaped by the stimulus structure, causing cells to synchronize if the local contrast in their receptive fields is similar. In essence, relaxation in the oscillator network solves a graph clustering problem with the graph representing feature similarity between different points in the image. We tested different model versions on the Berkeley Image Segmentation Data Set (BSDS). Networks with phase interactions set by standard representations of the feature graph (adjacency matrix, Graph Laplacian or modularity) failed to exhibit segmentation performance significantly over the baseline, a model of independent sensors. In contrast, a network with phase interactions that takes into account not only feature similarities but also geometric distances between receptive fields exhibited segmentation performance significantly above baseline.

1 Introduction

For decades the commonly accepted view of retinal processing has been that it provides a bank of independent, linear filters that decorrelate stimulus features in space and time, reducing the redundancy in the retina's representation [5]. Linear spatio-temporal filters factorized into center-surround spatial and biphasic temporal components followed by pointwise non-linearities encode local stimulus features in the spike rates of retinal ganglion cells (RGC) [17]. There remain, however, severe puzzles, unexplained by the textbook view of retina.

First, for retinal ganglion cells it would be inefficient to use spikes exclusively in a rate code with rather long integration window. This assumption is in conflict not only with theoretical principles, such as the efficient coding hypothesis [4], but with experimental observations. For example, it has been shown that time to first spike in RGCs can be very reliable, containing nearly as much information about the stimulus as spike rates [10].

Second, the circuitry in the anatomical retinal network is exquisitely complex, consisting of >60 distinct neuron types stratified into at least 12 parallel and interconnected circuits providing roughly 20 diverse representations of the visual world, discussed at length in [21], [22], [36], [11]. Simple linear spatio-temporal filtering requires only a handful of cell types in the outer retina, leaving the rest of the network unexplained. By "Occam's razor", the textbook view must be at least incomplete.

Third, the textbook model of retina fails to account for complex phenomena such as precise spiking of RGCs relative to the phase of network oscillations in the gamma range (50-80Hz) [25, 16]. Although the function of retinal oscillations is yet unknown in mammals, they have been observed in mouse [23], cat [24] and primate [27]. Further, gamma-band retinal oscillations have been causally connected to the perception of spatially extended stimuli in the frog [14]. Specifically, it has been observed that neurons in the cat lateral geniculate nucleus (LGN) often receive periodic retinal spike trains in the gamma band. Estimates of information rate in LGN spike trains suggest that in cells with periodic inputs, the spike train could multiplex two different types of information. While rate modulation in a courser time window encodes local stimulus contrast, a significant fraction of the total information is encoded by spike timing at a fine time scale, conveying the phase of the gamma frequency in the neurons input [16].

Fourth, computational models reflecting the text book view, such as the linear nonlinear Poisson (LNP) model and generalized linear models (GLM), predict RGC responses to a simple white noise stimulus [34] with reasonable accuracy. However, looking more closely, one observes pairwise correlations in retinal activity, even in the absence of stimulus (correlations) [33]. Taking into account these pairwise activity correlations improves decoding of retinal responses to white noise [29] – but does not explain why the retina introduces such correlations to begin with. The situation with ecologically relevant natural movie stimuli, in which pixels possess dependencies across space and time, is even more puzzling. The model prediction by independent encoding models

becomes rather poor [34], and even encoding models that include second-order correlations fail to replicate responses to natural movie stimuli [13]. We suggest to take these mismatches between retina and its current computational models as an encouragement to design and investigate novel computational models of retina.

Here we approach the challenge to design better retina models from a computational perspective and ask: "What type of image analysis could be computed in an array simple sensors with access to (center surround) image features, like found in retina, above and beyond independent sensors proposed in the textbook model?" Specifically, we follow the lead suggested in the discussion of experimental work [14] and investigate whether, in addition to encoding local image features, the retinal network can also extract spatially extended visual features and multiplex the extracted information into the retinal output using phase synchrony in periodic spike trains [16].

To concretely design a sensor network model with this function, we build on contributions provided in various streams of earlier work, the insight that image segmentation (IS) can be cast as a graph clustering problem [35], and the insight that, in addition to spectral methods, graph clustering can be efficiently solved in networks of phase-coupled oscillators [3]. To evaluate the performance of the model, the Berkeley Image Segmentation Dataset (BSDS) was essential. While the motivation for this work is to model a computation in retina, it should be noted that the network model we propose is still quite abstract. The model aims to serve as a proof of principle that the network computation could be efficiently performed by biological retinas, it is not intended as a neurobiologically detailed circuit model.

A coarse overview of the model is given in Fig. 1. The firing rate r_i in a coarse time window represents the local image contrast in the classical receptive field of neuron i . The similarity between pairs of local features in the image determines the strength of phase interaction between the periodic structure in the spike trains. Phase diffusion through the phase couplings does not change firing rates but produces sets of neurons with similar spike times on a fine time scale. These sets of synchronous neurons represent spatially extended image features, image segments. The resulting spike trains multiplex two types of information, local contrast in individual spike rates, and image segments in sets of neurons that fire nearly synchronously [16]. In our example, two image segments are represented by groups of neurons with different phases. Note that in this study, we only consider models of the phase dynamics, omitting aspects of spikes and spike rates.

The remainder of this paper is structured as follows. The Methods section describes prerequisites for our study from the literature. Section 2.1 concisely defines the putative computation of our retina model, image segmentation (IS) using simple image features available in retina, local contrast values or local center surround image features. The evaluation pipeline proposed in the BSDS image segmentation database [19] is explained, which is essential to quantitatively compare the performances of different models. Following [35], section 2.2 describes how image segmentation can be cast as a graph clustering

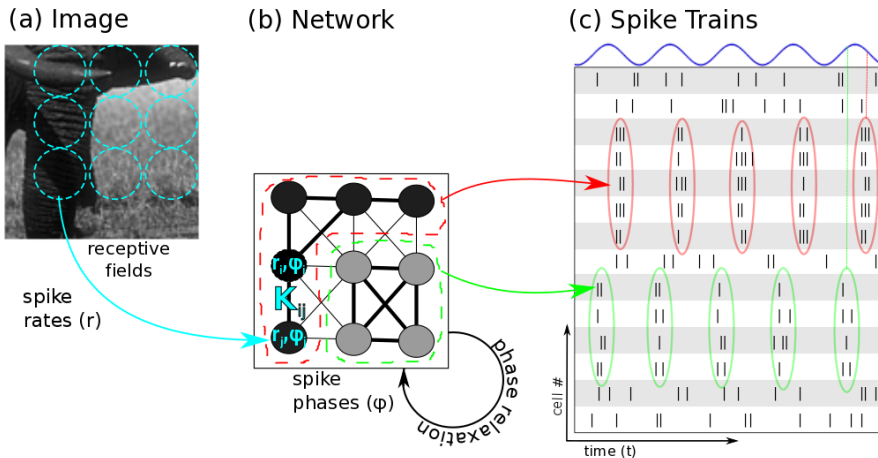


Figure 1: **Image Segmentation Model:** (a) Input image with superimposed retinal receptive fields (dashed cyan circles). (b) Network of retinal neurons. The neural firing rates r_i represent local contrast in the receptive fields. The phase interactions K_{ij} are displayed by the links between neurons. Line thickness represents the strength of the interaction which is set by the similarity of local features. Recurrent propagation in the network produces the phase structure ϕ_i of the periodic spike trains. (c) Resulting spike trains. Information about local features is represented in firing rates and segmentation is represented in phase structure.

problem, and how an adjacency graph is constructed for a particular image. Section 2.3 describes three common graph clustering methods from the literature, average association, graph Laplacian and modularity, that we will compare in our image segmentation experiments. Section 2.4 describes how, as an alternative to computing eigenvalues of the graph representation matrix, relaxation of phase-coupled oscillators can be used to solve graph clustering problems. This step is critical in mapping the computation of image segmentation to the network model in Fig. 1b.

The Results section contains original contributions of our study. Section 3.1 describes topographic modularity, a novel graph-clustering method based on modularity [26] that we propose for clustering multigraphs. Image segmentation can be understood as clustering of a multigraph, in which one type of edges represent feature similarity and the other geometric vicinity of the features in the image plane. Section 3.2 compares the performance of image segmentation of commonly used eigenvector-based "spectral methods" [8] for graph clustering to the method of phase relaxation [3]. We find that phase relaxation generally outperforms spectral methods, independent of the choice of a particular image graph or receptive field structure. Thus, our further experiments focus on phase relaxation, the method that also has the advantage of being easily implementable as an oscillation-based computation [15]. The central experimental results of our study are described in section 3.3. We compare segmentation

performances of different network models to a baseline segmentation algorithm based on thresholding image feature histograms, a computation which does not require a network. While the standard graph clustering methods are not able to significantly outperform histogram thresholding, one model stands out significantly, the network implementing topographic modularity. Section 3.5 describes experiments to elucidate why the network with topographic modularity outperforms the competitor models.

In the Discussion section we delineate the various implications of the presented results. We describe the predictions our model makes for future neuroscience experiments and its potential for applications of image processing with coupled sensors.

2 Methods

2.1 Berkeley Segmentation Data Set

Image segmentation is a challenging and important problem in computer vision and the Berkeley Segmentation Data Set (BSDS) is a standard benchmarking data set for many computer vision image segmentation algorithms [19, 1]. It consists of 500 large ($\sim 400 \times 300$ pixels) color images each with multiple (~ 5) human drawn boundary contours (green box in Fig. 2), as well as code provided for standard benchmarking and comparison of algorithms. Since image segmentation is closely related to boundary detection and quantification of boundary detection performance is more straightforward than that of image segmentation, segments in images are often recast as boundaries for benchmarking. Binary boundary pixel locations are compared to human drawn boundaries using the precision, recall, f-measure framework. In this context, "Precision" is the proportion of image pixels hypothesized by a method to belong to segment boundaries that agree with the ground truth. "Recall" is the percentage of ground truth boundary pixels that are found by a method. F-measure is the harmonic mean of Precision and Recall.

In order to leverage the BSDS resource, we must first convert the output of a segmentation model - a phase, spectral or feature activation map (blue box in Fig. 2) - into binary boundaries. Intuitively, a good segmentation of an image has been achieved if the model output map has very similar values within segments and large discontinuities at boundaries. We compute spatial derivatives ($\delta/\delta r$) in the output map and normalize the values between 0 and 1, allowing us to interpret resulting probabilistic boundary (pb) as the algorithm's confidence that there is a boundary between segments at a particular image location. We can threshold pb's at multiple values and compare each resulting binary boundary map (bb) to each human drawn ground-truth boundary map (gT), generating a pixel match set by a logical AND operation. Because human drawn boundaries are not precise down to the pixel, we allow small misalignment between gT and bb pixel including a pair in the match set if they are within d_t pixels of one another.

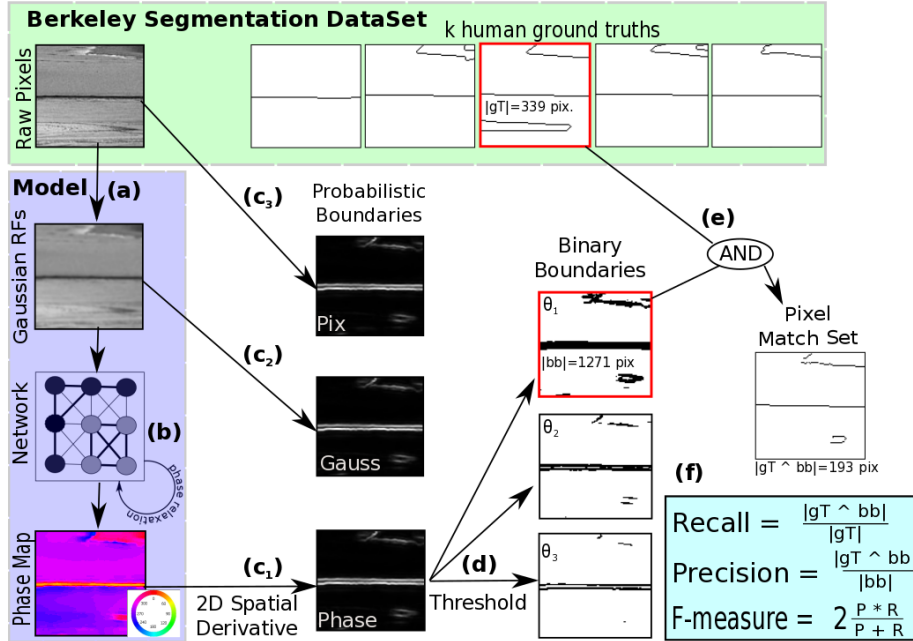


Figure 2: **Performance and benchmarking:** Input image patch and associated human drawn ground truth boundaries (gT) provided by BSDS is displayed in the green box. The operations performed by model are displayed in the blue box. Other steps of model evaluation are illustrated in the remainder. (a) Filtering the raw image patch with a Gaussian kernel ($\sigma = 1$). (b) Phase relaxation in the network (Fig. 1b) produces a phase map. (c₁) Spatial gradient operation ($\delta/\delta r$) and normalization resulting in probabilistic boundary map ($pb \in [0, 1]$). (d) Thresholding pb map at several values yielded binary boundary (bb) maps. (e) Match set was computed for each bb-gT pair at different distance tolerances, d_t . (f) Precision, recall and F-measure were computed by ratios of boundary pixel sets. (c_{2,3}) To assess the performance of network models relative to baselines, we repeated steps (c) - (f) on Gaussian RF and image pixels independent sensors models, comparing F-measures by subtraction.

We compared the ability of different phase coupled oscillator models to segment images from the Berkeley Segmentation Dataset (BSDS) [19]. The models differed in the phase couplings. One baseline model contained isotropic couplings, while the couplings in the other models AA, GL, M and TM were the transformations of the adjacency of features described above. We set the parameters σ_f , σ_d and σ_ω to adequate common values and performed for each method a parameter grid search in neighborhood connectivity R_M and scale K_s to maximize the average ΔF_b across 500 image patches. The oscillator frequency was 60 Hz (typical for retinal gamma oscillations) and we gave the networks 300ms to relax the phases, corresponding to an interval between saccades.

2.2 Image segmentation as a problem of graph clustering

Within a stage of visual processing, in which a set of local visual features is extracted, image segmentation can be viewed as a graph clustering problem [35]. Consider an image and its corresponding neural representation in retina or LGN, in which the activity in individual cells represent the strength of local center-surround features. Image segmentation consists of clustering sets of local image features that share properties and thus likely correspond to larger objects in the image. However, much more efficient than clustering pixel values, is to apply clustering on more sophisticated local image features, e.g., center-surround, edges, multi-scale textures, as done in state-of-the-art segmentation methods [35, 32].

The problem of graph clustering, that is finding “cliques” of strongly connected nodes in a graph, is obviously related to finding pixel sets with similar local features. To recast image segmentation in terms of graph clustering, one first uses kernels to construct an adjacency matrix, in which an element is large whenever two features have similar values and lie nearby each other in the image plane [35, 32]. The segmentation of the image corresponds to finding the communities (subsets of nodes) that are strongly interconnected within the community, and well separated from nodes outside. The goal then is to find non-trivial subsets of nodes that can be separated from one another by cutting through the minimum weight of edges, know as the "mincut" problem. Though a brute force, optimal solution to this problem would be combinatorically intractable, approximate solutions can be found efficiently by leveraging the machinery of "spectral graph theory" [8].

Following [35], we define a graph for segmenting an image by the *adjacency* matrix:

$$\mathbf{A}_{ij} = e^{-\frac{(f_i - f_j)^2}{2\sigma_f^2}} \cdot e^{-\frac{(r_i - r_j)^2}{2\sigma_r^2}} \cdot \left(1 - H(\sqrt{(r_i - r_j)^2} - R_M)\right) \quad (1)$$

with $H(x)$ the Heaviside step function. The first factor reflects the dissimilarity of the local features f_i and f_j , in our case local contrasts. It was found experimentally that $\sigma_f = 0.2$ provides reasonable dynamic range in adjacency weight distribution. The second and third terms reflect the distance between the local features in the image plane. Since we are interested how well segmentation can be performed in networks with local neighborhood connectivity and for simplicity, we null out the second term by setting $\sigma_r = \infty$ and add the third term, a binary rectangular Heaviside function $1 - H(\sqrt{(r_i - r_j)^2} - R_M)$ that is 1 within a maximum radius, R_M , and 0 outside. We explored R_M values of 1,3,5 and 10.

2.3 Three common graph clustering methods

The simplest strategy of graph clustering, referred to as *average association* (AA), is to analyze the adjacency matrix directly [32]. Eigenvalues of the adjacency quantify the amount of correlated structure and the associated eigenvectors

characterize the location of the correlated structure in the image. Other methods of graph clustering utilize transformations of the adjacency matrix, often incorporating the node "degree", $d_i = \sum_j A_{ij}$, which captures the total weight of connections to each node from all other nodes in the network. One such transformation we considered is the normalized *graph Laplacian* (GL) or Kirchhoff matrix: $L = D^{-1/2}(D - A)D^{-1/2}$ with diagonal matrix $D_{ij} = \delta_{ij} \sum_k A_{kj}$, δ_{ij} the Kronecker symbol. This strategy, combined with more sophisticated image features, forms the basis of a very successful image segmentation algorithm, the "Normalized Cut" [35]. The eigenvectors and associated smallest eigenvalues of the Laplacian matrix find divisions in the input characterized by large feature differences.

A third transformation of the association matrix we considered is *modularity* (M) [26], which has successfully discovered community structure in social and information networks, outperforming the graph Laplacian in these tasks. The modularity matrix can be written as

$$Q = A - N \quad \text{with} \quad N_{ij} = D_i D_j \quad \text{and} \quad D_i := \frac{d_i}{\sqrt{2m}} \quad (2)$$

where D_i and D_j denote the "degree" of nodes i and j respectively, normalized by the total weight of edges in the graph, $2m = \sum_k d_k$. Importantly, the null model matrix, N , contains the expectation of the weight value between each node pair N_{ij} based on the strength of connectivity of both nodes. In this way, an expected graph is constructed by assuming an otherwise random graph with node degrees constrained (an Erdos-Renyi random graph). Comparing the observed adjacency graph to the null model by subtraction reveals graph structure beyond what could expectedly be introduced by heterogeneous node degrees. In section 5, we discuss modularity further and introduce an extension, called *topographic modularity* (TM), with null model adapted for graphs embedded in space.

Once an associated matrix representing a graph is constructed, spectral methods have been predominantly used within the graph clustering community to find clusters within because eigenvalues and eigenvectors efficiently find an approximate solution to the combinatorially intractable "mincut" problem. It has been observed on simple networks that the eigenvalue spectrum of an associated matrix resembles the temporal progression of clusters discernible from phases of nodes in a Kuramoto network [2], this time evolution of clusters forming a hierarchical clustering of a network. Given this observation, we compute the time evolution of a phase coupled oscillator network dynamical system as an alternative to eigenvector-based graph clustering methods.

2.4 Kuramoto Phase Relaxation Model

The described graph clustering methods in 2.2 compute the eigenvectors of the associated matrices [8] which, in essence, is assessing anisotropic diffusion in these networks. This process has also been related to the path a random walker would take through the graph where edge weights represent transition probabilities and the distribution of electrical potentials on nodes in a resistor

network where an edge weight represents the conductance of a particular resistor [12]. A further parallel has been between eigenvector based methods for graph clustering and the "fundamental mode(s) of a spring-mass system" [35]. To rigorously investigate this last claim, we simulate phase relaxation in a network of Kuramoto coupled oscillators [18] with networks defined by methods described in 2.2.

Here we followed [2] and assessed diffusion properties by relaxing a network of phase-coupled oscillators :

$$\Delta\phi_i = \omega_i + \sum_j K_{ij} \sin(\phi_i - \phi_j), \quad K_{ij} = k_s M_{ij} \quad (3)$$

with each node's natural frequency $\omega_i = 60Hz$ and where M_{ij} is one of the graph matrices mentioned above. For intuition, Eq. 3 loosely simulates a lattice of oscillating masses connected by different size and signed springs. The lattice is shaken at initialization and through the relaxation dynamics, masses connected by strong positive springs are attracted in phase while strong negative springs repel one another. In the original Kuramoto model [18], couplings K were set to be uniform, supporting isotropic diffusion. As a baseline, we also investigated the effects of isotropic diffusion (ISO) for image segmentation. Unlike the uniform network, a network with heterogeneous weights relax to stable states containing multiple distinct clusters of phase aligned oscillators.

In the implementation of the model, the overall positive scaling factor k_s is critical. If coupling weights are too large, phasers will spin wildly in response to even small phase differences. Conversely, if too small, oscillators will adjust their phase too slowly and the relaxation will not converge in time. Importantly, the phase relaxation was limited to 300ms or 20 periods of the 60Hz signal, which is the average duration of fixation before a saccade brings the eye's gaze to a new point, refreshing the input and beginning the computation once again. The value for the k_s parameter was set for each graph individually based on mathematical considerations in equation 3. A middle value k_s^{mid} was chosen so that the phase change of the node with largest degree $D_{k_{max}}$ is limited to $\pi/2$ radians in one full period of the 60Hz signal when all its neighbors are aligned $\pi/2$ radians away and exerting maximal pull.

$$k_s^{mid} = 60Hz \cdot \frac{2\pi}{\pi/2} \cdot D_{k_{max}} \quad (4)$$

We then bracketed that value above and below by an order of magnitude.

The final result of the phase relaxation simulation is a phase map with a phase value, $\phi_i \in [0, 2\pi]$, associated with every node, i , in the network and corresponding location i in the image. Spectral methods also yield a value associate with each location, i , in the image with $v_i \in [-\infty, \infty]$. In order to compare our results to other algorithms using the BSDS resources, we convert these maps to probabilistic boundaries and recast the image segmentation problem as a boundary detection one as discussed in Section 2.1 and illustrated in Fig. 2.

In practice, two meta parameters, r_M defining the neighborhood structure of the Adjacency graph and k_s defining an overall scaling on the strength of phase interactions in the network, impacted image segmentation performance. They were optimized for each method and results shown are with optimized parameters, shown in Fig. 11. To optimize parameters for each method, we performed segmentation of 500 image patches with four r_M values ranging from 1 to 10 and bracketing k_s as discussed above and chose the parameter settings with best average performance across all images and across d_t . Fig. 3 illustrates the procedure for one particular method. It shows average performance across k_s values for optimal r_M on the left and performance across r_M values for optimal k_s on the right. Fig. 4 shows the effect of different parameter settings on one example image patch.

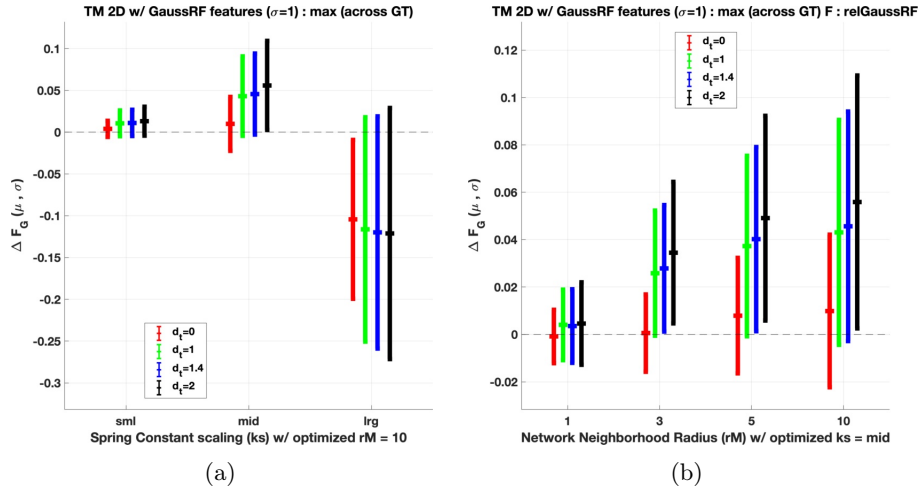


Figure 3: **Hyper-parameter optimization:** Network neighborhood graph structure r_M and coupling spring-constant scaling k_s are important meta parameters of the algorithm, discussed in Sections 2.3 and 2.4 respectively. We plot mean and standard deviation across 500 image patches of ΔF -measure relative to Gaussian RF independent sensors for the 2D topographic modularity network. Colors indicate pixel distance tolerances d_t (see Fig. 2 for explanation). *Left* panel shows performance at three k_s values, with r_M fixed at optimal. *Right* panel shows performance at four r_M values, with k_s fixed at optimal. Fig. 4 shows the effects of the different parameters on a single example image patch.

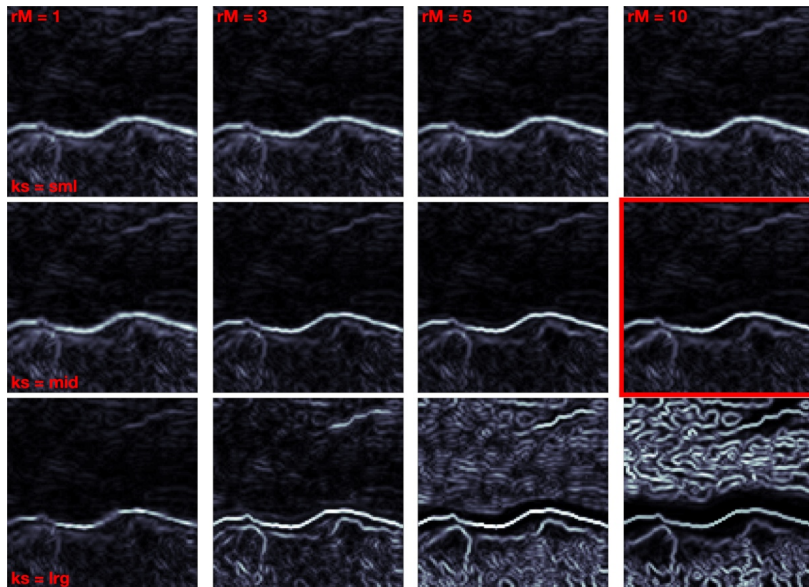


Figure 4: **Effect of hyper-parameters, single image patch example:** Probabilistic boundary maps shown for resulting phase distribution from TM 2D method for combinations of 3 k_s (rows) and 4 r_M (columns) hyper-parameters.

3 Results

3.1 Modularity null models for images

An image can be described by a multi-graph, in which pixels or local image features are represented by nodes and each pair of pixels has two different types of edges connecting them. One edge type represents geometric distance in the image plane and the other edge type represents feature differences. The two types of edges are given by adjacency matrices, resulting from the two types of distances and corresponding kernel functions (like a Gaussian kernel), as in Eq. 1. Shi and Malik [35] proposed a way to collapse this multi-graph of an image to an ordinary graph by forming the Hadamard product of the two adjacency matrices. An entry in the resulting single adjacency matrix A represents the two distinct similarities between pixels, geometric proximity and feature similarity by a single number. Specifically, an entry in A can only be large, if both, distance and feature differences are small in the corresponding pair of pixels. In order to find image segments, researchers then used "spectral" graph clustering methods on the matrix A [32, 35].

For some graph clustering methods, such as modularity [26], the collapsing of the multi-graph into an ordinary graph destroys information, which is critical for segmenting images. The modularity matrix consists of the difference of the adjacency matrix and a null model. The null model represents an average

adjacency value. In the standard modularity method, Eq. 2, the average is computed from the degrees of the two nodes involved, the row and column sum of the collapsed graph. However, in natural images, the average feature similarity of a pair of pixels is a function of geometric distance [31], see also Fig. 20 in supplemental section B.6. Thus, an appropriate null model for images should also depend on the geometric adjacency matrix.

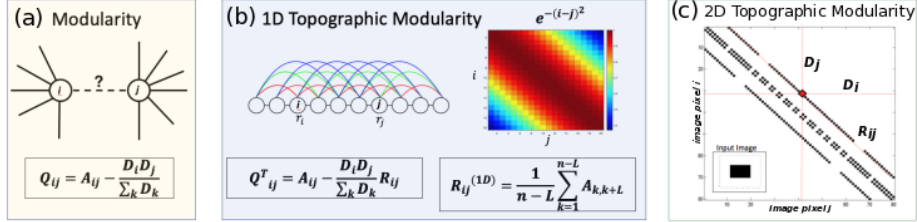


Figure 5: **Modularity null models & space:** In the null model of Newman’s modularity [26] (*panel a*) the average weight between nodes i and j is proportional to the product of their node degrees ($D_i \cdot D_j$). The topographic modularity’s null model (*panels b & c*) additionally includes a distance-dependent factor, R_{ij} , which is the average edge weight between all node pairs in the graph separated by the same distance that separates nodes i and j . *Panel b* illustrates R_{ij} for a schematized 1D graph, shown with edges colored based on distance between the nodes they connect. Inset plot shows geometric factor in the topographic null model. Each term in $R_{ij}^{(1D)}$ is an off-diagonal sum in the adjacency matrix. *Panel c* shows the mask associated with a single geometric distance in a 2D image. Here $R_{ij}^{(2D)}$ at 1 pixel separation has a complex structure in the Adjacency matrix for even the simple binary image shown in the inset.

To address this issue, we devised a novel graph clustering method called *topographic modularity* (TM) in which the null model takes topographic distance in the image plane into account. Like the standard modularity [26], see Fig. 5a and Eq. 2, an entry of the topographic modularity matrix, Q^T , is the difference between the entry A_{ij} and the expected connectivity, captured by the null model, N_{ij} . Here, the topographic null model N_T accounts for distance dependent factors in feature similarity with the R_{ij} term in addition to node degree heterogeneity.

$$Q_{ij}^T = A_{ij} - N_{ij}^T \quad \text{where} \quad N_{ij}^T = D_i \cdot D_j \cdot R_{ij} \quad (5)$$

The R_{ij} factor represents the average connectivity between all node pairs that are separated by the same geometric distance as the nodes i and j . For a network in space along a 1D line, Fig. 5b, the distance dependent contribution to the null model can be written mathematically as

$$R_{ij}^{(1D)} = \left(\frac{1}{n-L} \right) \sum_{k=1}^{n-L} A_{k,k+L} \quad (6)$$

where L is the distance separating nodes i and j (i.e., $L = r_i - r_j$) and n is the total number of nodes (or pixels or features in the image). In 1D, the average connectivity of all nodes separated by a distance L is equal to the mean along the L 'th diagonal.

For networks with 2D grid-like geometry, like Adjacency graphs constructed from images, the computation of $R_{ij}^{(2D)}$ is more involved, yet the interpretation is the same. Reshaping a 2D image into a 1D vector so that similarity relationships can be represented in a 2D matrix introduces discontinuities in spatial relationships between entries in the matrix. Weights between nodes separated by a particular distance can be labeled by a mask specific to the dimensions of a particular image. Fig. 5c shows the weights between all neighboring nodes ($L = 1$) in the network derived from the 11 x 11 binary image in the inset. For completeness, we show $R_{ij}^{(2D)}$ masks for other pixel separations in supplement section B Fig. 19.

Before comparing the different null models in an image segmentation we compare how well they capture the structure of an adjacency matrix of an image. The null model in Newman's modularity, by construction, is a "consistent" estimator of node degrees [7], ensuring that $\sum_j N_{ij} = \sum_j A_{ij}$ (blue line in Fig. 6 middle). However, it is clearly the wrong null model for natural image Adjacency graphs for two reasons. First, the null model incorrectly contains positive diagonal weights in proportion to D_i^2 , although the diagonal elements of the adjacency matrix are zero. Second, it does not capture the distance dependence of the adjacencies, thereby underestimating average adjacency between proximal nodes and overestimating it for distant node pairs. Both problems manifest in the difference between the blue and the dashed lines in Fig. 6 bottom.

While the TM-1D and TM-2D null models are not strictly consistent in node degree or distance dependence, they are nearly so (green and red lines respectively in Fig. 6). Introducing distance-dependent statistics into the TM-1D null model corrects for the spatial "inconsistency", vastly improving estimates of edge-weight over M . TM-2D offers improvement over TM-1D due to further refinement of its null model, see Eq. 6 and surrounding text. Further discussion of null model consistency and bias in the supplemental section ??.

Null Model Consistency

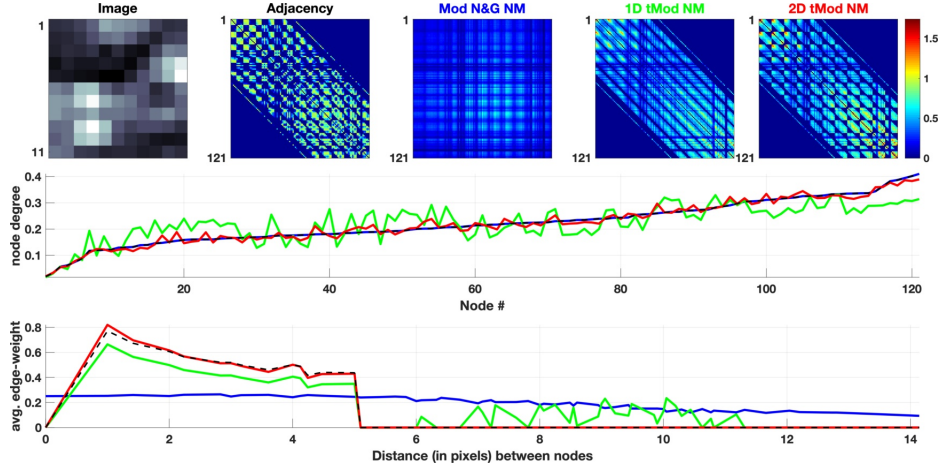


Figure 6: **Null model consistency:** *Top row* from left to right shows image patch, the adjacency (black) constructed from the patch with $r_{max} = 5$, and null models for modularity (blue), 1D topographic modularity (green) and 2D topographic modularity (red), with colorbar indicating edge weight. Models represented by line colors in plots as well. *Center plot* shows average node degree (row sums in each matrix) sorted by strength in adjacency. *Bottom plot* shows average edge weight as a function of distance in the image plane.

Importantly, the difference between an adjacency value and its average in the modularity can become negative. In a Kuramoto net relaxation, these negative weights mediate phase repulsion and introduce targeted phase desynchronization, see Sec. 2.4, at boundaries in an image where gross image statistics change. In contrast, if the modularity value between a node pair is positive, it contributes to phase synchronization. Fig. 7 illustrates image segmentation performance before and after phase relaxation through connections defined by M (in blue), TM-1D (in green) and TM-2D (in red). While M does not significantly change image segmentation performance over Gaussian RF independent sensors, TM-1D does so (p-value ~ 0.004) and TM-2D does so even more (p-value $\sim 4 \cdot 10^{-7}$). With TM-2D, we see improvement for $\sim 460/500$ image patches.

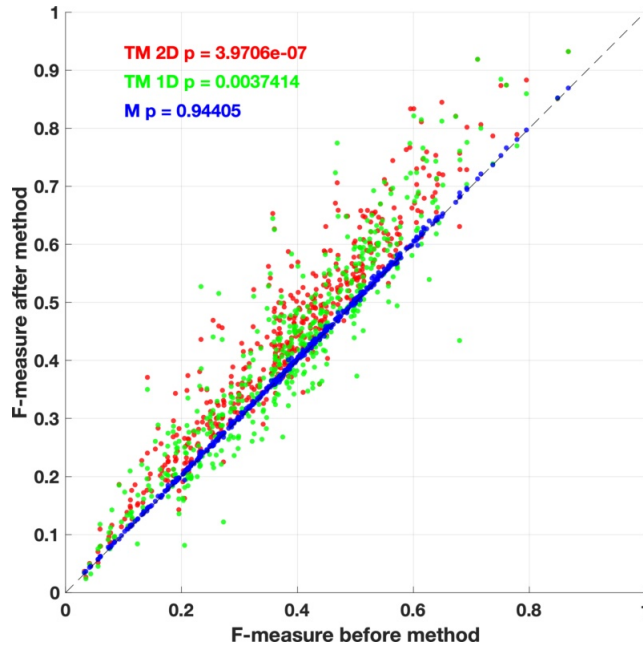


Figure 7: **Modularity performance comparison:** Each scatter point represents one image patch. Newman’s modularity (M) in blue, 1-dimensional topographic modularity (TM 1D) in green and 2-dimensional topographic modularity (TM 2D) in red. Points above the unity line indicate image patches with improved image segmentation with network phase relaxation over-and-above Gaussian RF independent sensors. P-values quantify the difference between F-measure distribution across 500 image patches before and after network computation.

3.2 Broad comparison of models on image segmentation

Following [15], we investigate the idea whether a phase-coupled network of simple sensors of local image features, similar to those in the retina, could at the same time represent local and contextual image features in its output. Specifically, phase interactions mediated through heterogeneous network edges which are influenced by local features similarities can segment an image, grouping regions within a segment into the same relative phase and introducing phase breaks at segment boundaries. In a biological system, the contextual image information encoded by phase can be represented by the timing of spikes and be multiplexed into spike trains, whose rates represent the local features Fig. 1.

This idea is tested on images provided in the Berkeley Segmentation Data set (Sec. 2.1). For an image patch, we construct a graph based on local features in the image (Sec. 2.2) and segment the image by either computing eigenvectors or by performing anisotropic phase diffusion in a Kuramoto net. Computing a spatial derivative on either eigenvectors or the final phase distributions and normalizing values between 0 and 1 converts the output into probabilistic boundaries, which

can be quantitatively compared to assess relative performance of different image segmentation methods.

We ask whether a phase-coupled sensor array can add to an image segmentation that can be done based on the independent sensor measurements alone. Thus, the network computation must outperform two baseline methods. The first method computes normalized spatial gradients on the raw image pixels (magenta, RawPix). In the second method the image pixels are first convolved with Gaussian receptive fields, roughly similar to those measured in retina (cyan, GaussRF). As a third baseline method, we include isotropic diffusion in a network with homogeneous phase couplings between nearest-neighbor nodes (black, ISO).

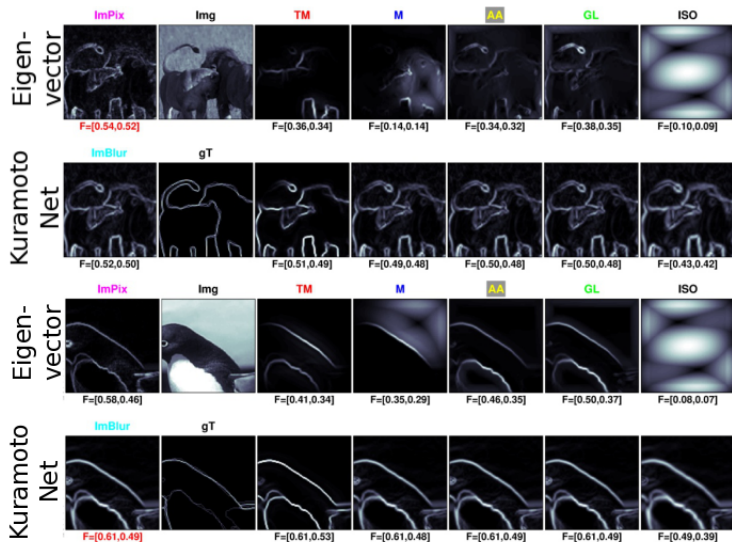


Figure 8: **Spectral methods vs. Kuramoto Net Examples:** Two example image patches (top two rows and bottom two rows) show probabilistic boundaries found by different network (TM, M, AA, GL) and baseline models (ImPix, GaussRF, ISO). Network models are segmented using eigen-methods (1st and 3rd row) and Kuramoto Net phase relaxation (2nd and 4th row). Qualitatively, boundaries found with spectral methods are less crisp and more localized than those found with Kuramoto Net phase relaxation.

Probabilistic boundaries (pb) can be interpreted as the algorithm’s confidence that a boundary exists between two segments at a particular location in the image. Fig. 8 shows pb’s resulting from the segmentation of different networks constructed from the same image patch, either by computing eigenvectors and by performing Kuramoto net relaxation. Qualitatively, we observe that eigenvectors seem to focus a spotlight on a region of the image patch while information propagated through the Kuramoto Net covers all parts of the image patch. Regardless of the network method used, boundaries found with the Kuramoto net are crisper and extend further across the image patch than do those found

by computing eigenvectors.

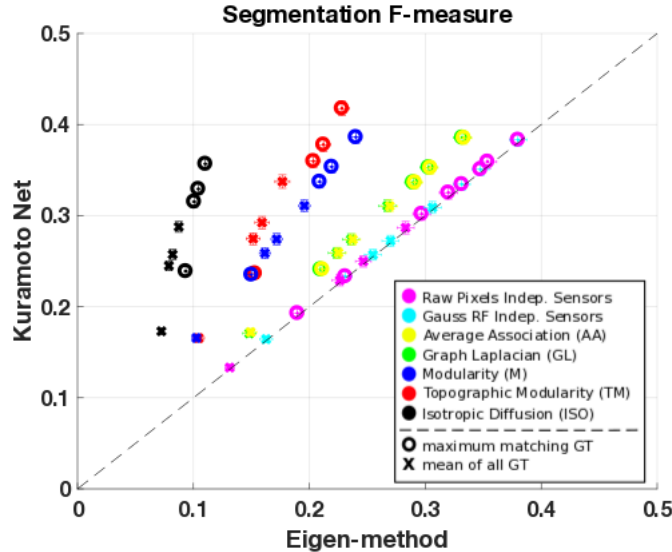


Figure 9: **Spectral methods vs. Kuramoto Net Statistics:** F-measure computed across 500 image patches, mean and standard error errorbars. Colors indicating different network and baseline models are used consistently throughout this paper. Circles indicates that F-measure for each image patch taken for maximum matching GT and x's shows mean value across all GT's. Network models built with Gaussian RF features are segmented by the best combination of the top 3 eigenvectors on the x-axis and by the phase distribution after Kuramoto Net relaxation on the y-axis. The dashed unity line indicates equal performance and the independent sensors baseline models (magenta and cyan) do not deviate from it.

To assess whether this trend in image segmentation performance is statistically significant, we calculate Precision, Recall and F-measure across 500 image patches, shown in Fig. 9. Plotting F-measure statistics for network and baseline models segmented by Kuramoto-net and Eigen-methods, we find that segmentation without network computation (magenta and cyan) outperforms the results from the best combination of the top 3 eigenvectors, regardless of the model. We also find that all scatter points lie above the unity line, indicating superior image segmentation performance of anisotropic phase diffusion in a Kuramoto net verses the spectral clustering methods. As a consequence of this observation, we focus in the reminder on the superior methods based on Kuramoto Nets.

3.3 Influence of receptive fields choice

We further observe that the features from which networks are constructed influence segmentation performance achieved. This comes as no surprise since

state-of-the-art image segmentation algorithms rely on a combination of sophisticated spatially-extended features.

We constrain our investigation to the relatively simple and local stimulus features that retina is supposed to have access to. Specifically, we investigate the difference in segmentation caused by switching between raw pixels and Gaussian receptive fields with different radii. Again, we compare the segmentation performance of networks with phase relaxation to baseline models representing independent sensors, and a model with isotropic diffusion through a homogeneous neighbor connections. We find that Gaussian receptive field features provide better segmentation than raw image pixels both when used as independent sensors and to construct phase interaction networks. Fig. 10 shows the segmentation performance (F-measure and the change in F-measure relative to the independent sensors image pixels baseline model) as a function of pixel match distance tolerance (d_t).

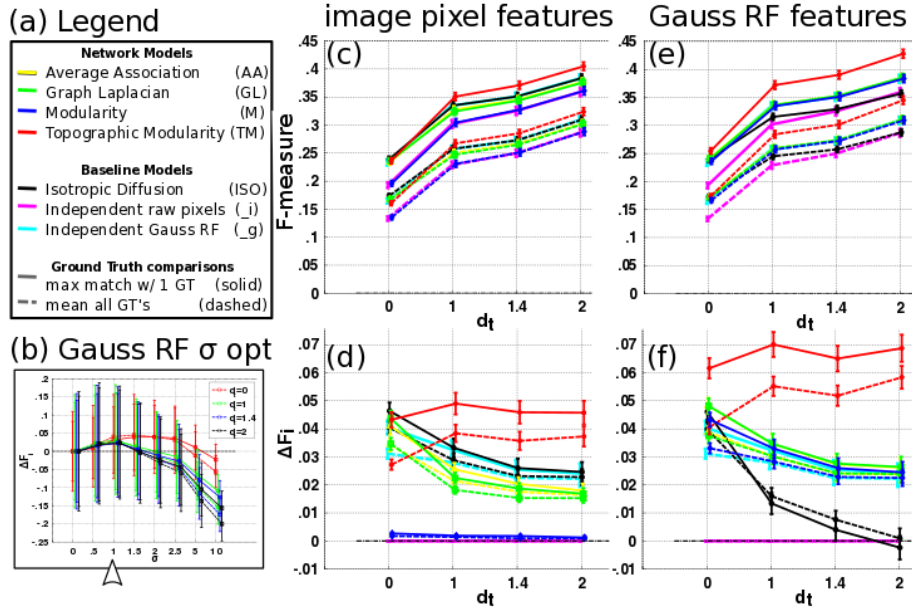


Figure 10: **Gaussian RFs improves segmentation:** Performances of 4 anisotropic diffusion and 3 baseline models are compared using raw image pixel features and Gaussian RF features, *center and right columns respectively*, lines representing average and bars standard error across 500 BSDS image patches. (*a*) Colors indicate different models and line styles indicate ground truth comparison as in Fig. 9. (*b*) Optimal spread, σ , of Gaussian RF's chosen by maximizing change in F-measure relative to the independent raw pixels baseline model, ΔF_i , averaged across all image patches. Recall d_t is the "distance tolerance" when computing the pixel match set, Fig. 2. Optimal performance for all d_t values obtained for Gaussian RF $\sigma = 1$. (*c*) F-measure and (*d*) ΔF_i when models receive raw image pixels as features. (*e*) F-measure and (*f*) ΔF_i when models receive Gauss RF activation as input features.

For small tolerances d_t in the F-measure (see section 2.2) the simpler isotropic phase diffusion model was a surprisingly strong competitor, even beating some of the anisotropic networks (black lines in Fig. 10c and d). Isotropic diffusion with optimized parameters provides mild smoothing of image structure, which operates indiscriminately within and across segments. To introduce the effect of smoothing in other models, we introduced Gaussian RF features. The filters corresponded to optical blur and the extended (centers of) receptive fields in retinal ganglion cells. Fig. 10b shows segmentation performance as a function of the width of the Gaussian filter, σ , and tolerance parameter, d_t . We find that Gaussian RF features with $\sigma = 1$ were beneficial and near optimal across different tolerance values. Interestingly, the size of the optimal Gaussian coincides with the size of retinal ganglion cell receptive fields measured in primate retina [9]. See supporting information A for further discussion.

Fig. 10e and 10f show the improvement in segmentation performance using Gaussian features above using image pixel features. In particular the method TM displayed a significant increase in ΔF_i which became more prominent for larger pixel match distance tolerances d_t . Among all methods TM was able to improve segmentation performance the most, compared to that achievable with the Gaussian RF independent sensors model.

3.4 Detailed model comparison between most promising models

To assess the overall performance of different models on the diverse input images, each model was run with optimized parameters. Fig. 11 shows image segmentation performance improvement from Gaussian RF independent sensors. Here the models TM-1D, TM-2D and ISO were significantly different from the three other methods that stayed near baseline $\Delta F = 0$. ISO stayed below baseline because the input kernels provide near optimal blur and therefore additional isotropic blurring deteriorated the segmentation performance. TM-1D performs well too, but not as well as TM-2D. This is because the null models are increasingly accurate, section 3.1. Shown results are with best matching ground truth. Results hold with average across all human drawn ground truths, though less pronounced.

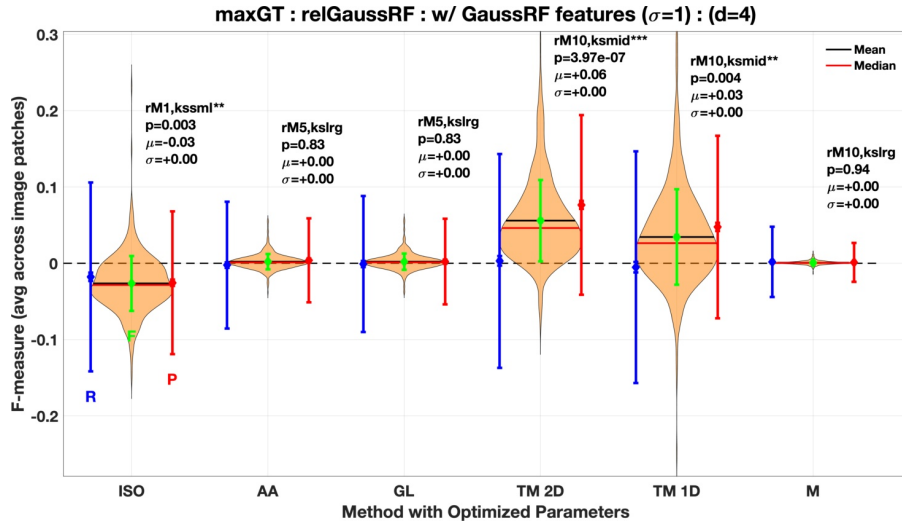


Figure 11: **ΔF -measure model comparison:** Violin plots show ΔF_G distribution with moments of Δ Precision, Δ Recall and ΔF -measure distributions across 500 image patches in red, blue and green, respectively. ΔF relative to Gaussian RF Independent Sensors model. Optimal hyper-parameters(r_M, k_s), statistical significance, p-values and distribution moments indicated above each method. Performance of ISO, TM-1D and TM-2D models relative to Gauss RF are statistically significant, as determined by Mann-Whitney U (aka rank-sum) test.

Segmentation performance via anisotropic phase diffusion in a Kuramoto net depends critically on the structure of the phase couplings. Kuramoto nets using the graph Laplacian, average association or Newman’s modularity as the phase couplings do not improve segmentation performance significantly over the independent sensors Gaussian RF baseline model. Only the Kuramoto model with the topographic modularity as phase couplings increases segmentation information over baseline independent sensors, homogeneous network and competitor heterogeneous network models, as quantified by the F-measure.

3.5 Why is the Kuramoto model with topographic modularity superior?

The F-measure combines the performance measures Precision and Recall, each with intuitive interpretations described in section 2.1. To analyze the differences between our different models, we separately plot the precision and recall distributions in Fig. 12. Note the position of curves for each network method relative to the independent sensors Gaussian RF baseline model (cyan dashed curve). Focusing first on the F-measure, in panel a, three of the network models (AA in yellow, GL in green, M in blue) did not show significant differences. The ISO model (black) degraded segmentation performance while the TM models (red

& magenta) improved relative to the Gaussian RF baseline. In panels b and c, the precision distribution of both TM models shifts significantly to higher values while the recall distribution shifts only slightly to lower values. Thus, the performance improvement of the TM model is mainly caused by increased precision, reflecting superior ability to suppress spurious boundaries, texture or "noise" in the probabilistic boundary maps.

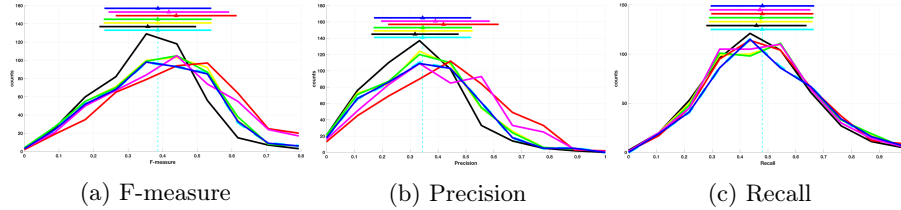


Figure 12: **Precision & Recall model comparison:** (a) F-measure, (b) precision and (c) recall across 1000 image patches for Gaussian RF independent sensors baseline model and 4 network models with optimized parameters and $d_t = 2$. Distribution μ and σ denoted above. Note colors same as in Figs. 9&10.

To better understand the computation in the TM-2D model, we visualize changes to Precision and Recall together for individual image patches in Fig. 13.

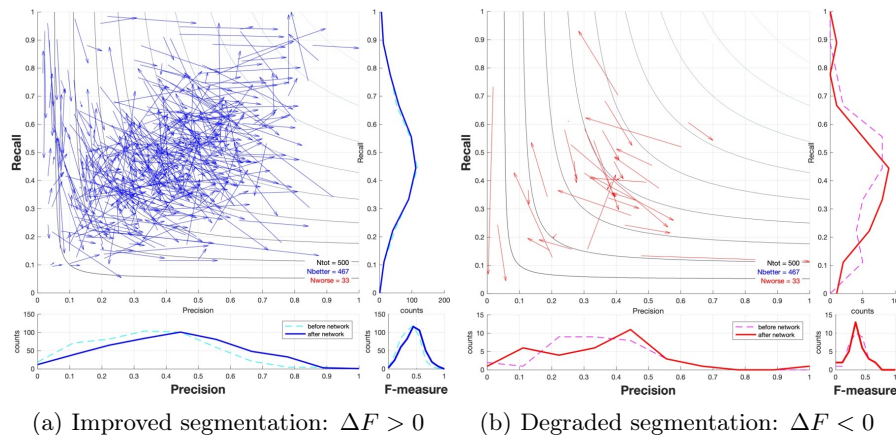


Figure 13: **Δ Precision and Δ Recall with TM-2D model:** In *panel a*, arrows show change in P & R for 467 image patches where network increased F-measure. Arrow tails indicate values before network relaxation and heads values after. Surrounding are distributions showing P,R,F before network (in cyan) and after (in blue). *Panel b* shows the same for 33 image patches where network decreased F-measure. Distributions before network in magenta.

The TM-2D network relaxation improved segmentation for $\sim 93\%$ of all image patches, in blue, panel a. Clear positive shifts in the precision and F-measure distributions can be observed from the independent sensors Gaussian

RF model (dashed cyan) to the phase output from the TM-2D network relaxation (solid blue). No clear trend emerges for the recall distribution with improved images. No clear trend exists for images where the network relaxation decreased performance. For some precision increased, and recall decreased. For others, vice versa.

3.6 Visual assessment of model performances

Finally, to provide some intuition what a ΔF value means for individual images, some examples are shown in Fig. 14. Compared to the results from other methods, the TM model produces probabilistic boundaries (pb's) that are often thinner and cleaner.

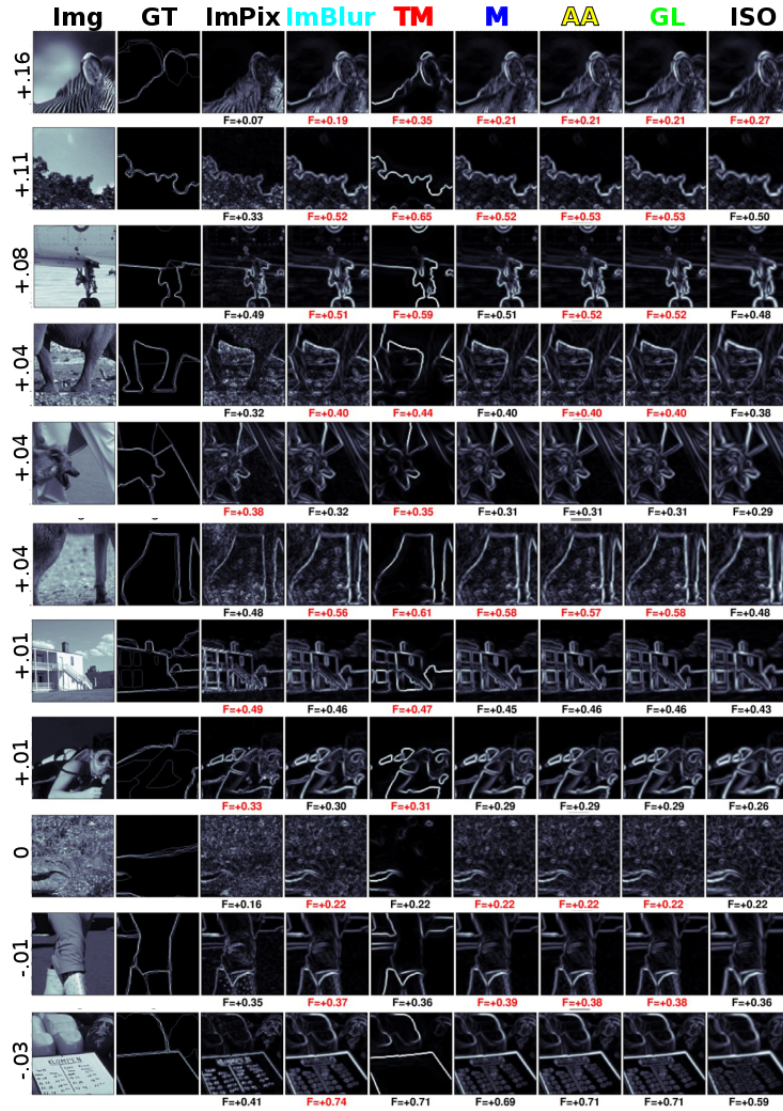


Figure 14: **Representative image patches:** Each row shows one example image patch ordered by change in F-measure between Gaussian RF independent sensors baseline and TM models (indicated on left). Columns show image pixels, gT boundaries and pb maps obtained from raw pixels, Gaussian RF and 5 network models. Mean F-measure value across all gT's noted below each pane is red if $\Delta F_G > 0$.

Further, in Fig. 15, we show samples of image patches with varying image segmentation performance relative to the Gaussian RF independent sensors model.

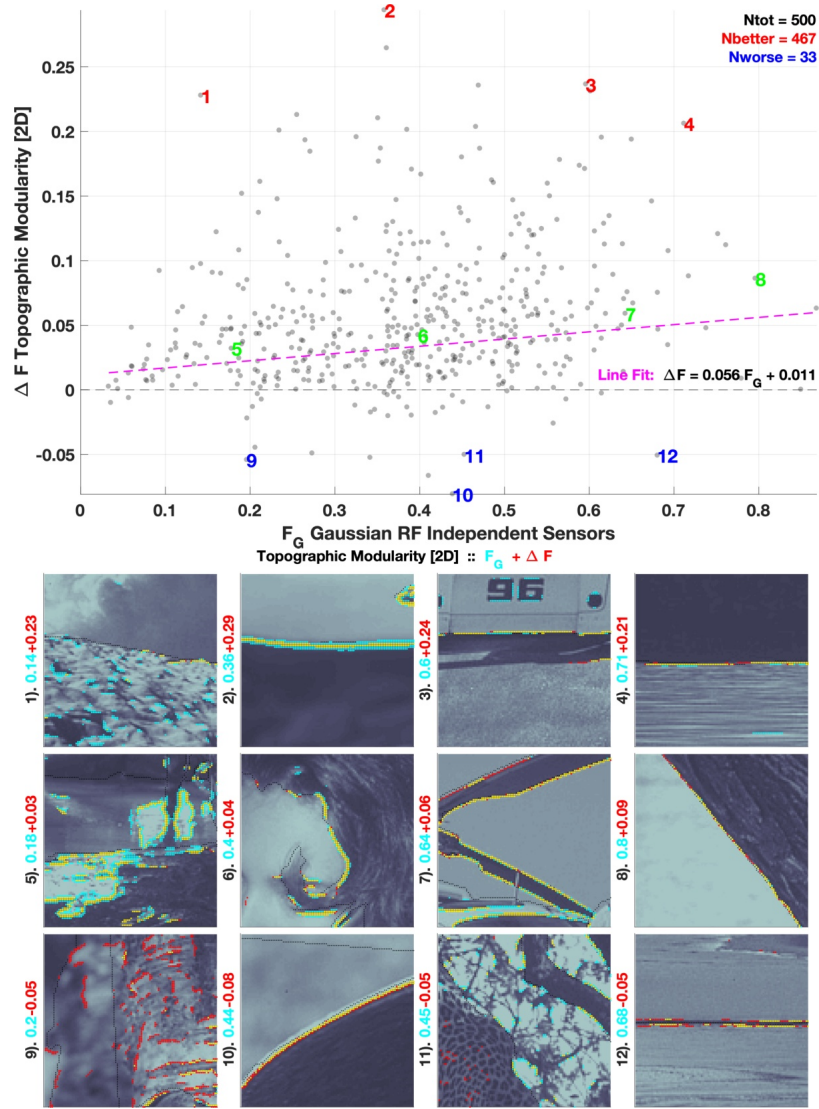


Figure 15: **Examples of TM 2D model performance:** *Top* panel scatters F-measure in Gaussian RF independent sensors model vs. ΔF after 2D topographic modularity network phase relaxation. Out of 500 total image patches, 467 show positive improvement. Best fit line to scatter points in magenta. Colored numbers indicate randomly sampled image patches (shown in bottom panel) where ΔF performance is best (#1-4), average (#5-8) and worst (#9-12). *Bottom* panel shows image patches with best matching ground truth boundaries, in black. Yellow points indicate pixels found to be boundaries both by the Gaussian RF independent sensors model and the topographic modularity network model. Cyan number and points indicate F-measure under Gaussian RF model and boundaries found only by it. Red number and points indicate ΔF after TM 2D network phase diffusion and boundaries found only by TM-2D. Note that image patches are shown at $1/2$ contrast to highlight boundaries found.

4 Discussion

Here we have shown that phase relaxation in coupled oscillators receiving inputs from simple image sensors (with unoriented Gaussian receptive fields) can provide image segmentation performance above and beyond the baseline, the segmentation performance that can be achieved by just using local contrast measurements. First, we have demonstrated that the type of graph clustering matters, the common spectral methods do not perform as well as relaxation in a Kuramoto model [2]. Second, we have demonstrated that the graph derived from the image structure matters. Specifically, we introduced topographic modularity, a modularity matrix that can capture the distance dependence in the statistics of image features. We find that a Kuramoto model using the topographic modularity matrix as phase couplings was the only network model that significantly outperformed the baseline.

A critical ingredient in the successful model is negative phase coupling weights, which introduces phase desynchronization at segment boundaries. Interestingly, we saw the best segmentation results with Gaussian receptive fields sizes similar to those measured in retina [9]. In essence, the successful segmentation model provides a "cartoonization" [37] of images - smoothing texture and variation within segments while maintaining crisp segment boundaries. Examples of phase relaxation results on two sample images are shown in Fig 16. Note the halos at the base of the lizard tail and surrounding the elk, where low contrast segment boundaries have been accentuated.



Figure 16: **Two examples of cartoonization:** Original images on left and resulting phase of TM-2D network computation on right

We quantify performance on the BSDS and show that anisotropic phase diffusion through the TM-2D improves F-measure significantly above baseline performance. This improvement is obtained by increased Precision with only slightly decreased Recall. However, there are some caveats with benchmarking our retina model on the BSDS data. First, BSDS is designed for state-of-the-art image segmentation methods that require combinations of sophisticated image filters, etc. However, context extraction in the retina can only use the simple image features of retinal cells, such as center surround features. Second, human image segmenters that provide the ground truth in the BSDS database can take advantage of the full image in color, while our model has only access to a 100×100 pixel image patch in greyscale. Third, humans segmenting images use consciously and unconsciously high-level semantic information to draw boundaries, while our algorithm just uses information from the image patch.

The model presented in this work is abstract and does not directly map to the biological features of retina. But some experimental evidence supports the plausibility of a computation in retina, as proposed by our model. Ganglion cell spike trains have been observed to be periodic in the Gamma frequency range [25] and the phase of this periodicity is transmitted with high precision through thalamus spikes to cortex [16]. The time to first spike in ganglion cells is quite

precise [10] and provides a possible mechanism for phase initialization following global suppression during eye saccades [30]. The phase coupling (without amplitude coupling) in our model could result simply from weak interactions between retinal cells, that slightly advance or delay spikes without adding or removing them. Both, phase synchronization and desynchronization through positive and negative weights in the model can be mapped onto excitation, inhibition and inhibition-of-inhibition circuits in retina. The spatial null model’s distance dependent term, R_{ij} term in Eq. 2, which requires global knowledge in the model could be implemented in retina via sampling through long distance inhibitory interactions from polyaxonal amacrine cells [28] or through eye movements implementing a temporal null model based on comparing feature similarity at a current stimulus location to feature similarity at a previous fixation. However, one central feature in our model still lacks experimental support. The model requires a mechanism for fast adaptation of the phase couplings to a particular stimulus.

Our modeling results suggest that, in principle, a coarse image segmentation or grouping/clustering of image features could be computed at the first stage of visual processing, in retina. While individual cell spike rates encode local stimulus contrast features through Gaussian-like receptive fields of ganglion cells, fine-time spike synchrony across the cell population encode extra-classical receptive field features, such as extended image segments. Fine-time correlations are multiplexed into ganglion cell spike-trains alongside with the rate-coded local stimulus features [15].

Acknowledgements: CW has been supported by the Systems On Nanoscale Information fabriCs (SONIC) program, FTS has been partly supported by grant 1R01EB026955-01 from the National Institute of Health. Both researchers have been partially supported by NIH award 1R25MH109070-01.

References

- [1] P. Arbelaez, M. Maire, C. Fowlkes, and J. Malik. Contour detection and hierarchical image segmentation. *IEEE transactions on pattern analysis and machine intelligence*, 33(5):898–916, 2011.
- [2] A. Arenas and A. Diaz-Guilera. Synchronization and modularity in complex networks. *The European Physical Journal Special Topics*, 143(1):19–25, 2007.
- [3] A. Arenas, A. Díaz-Guilera, and C. J. Pérez-Vicente. Synchronization reveals topological scales in complex networks. *Phys. Rev. Lett.*, 96:114102, Mar 2006.
- [4] J. J. Atick and A. N. Redlich. What does the retina know about natural scenes? *Neural computation*, 4(2):196–210, 1992.
- [5] H. B. Barlow et al. Possible principles underlying the transformation of sensory messages. *Sensory communication*, 1:217–234, 1961.
- [6] I. H. Brivanlou, D. K. Warland, and M. Meister. Mechanisms of concerted firing among retinal ganglion cells. *Neuron*, 20(3):527–539, 1998.
- [7] Y.-T. Chang, R. M. Leahy, and D. Pantazis. Modularity-based graph partitioning using conditional expected models. *Physical Review E*, 85(1):016109, 2012.
- [8] F. R. Chung. *Spectral graph theory*, volume 92. American Mathematical Soc., 1997.
- [9] L. J. Croner and E. Kaplan. Receptive fields of p and m ganglion cells across the primate retina. *Vision research*, 35(1):7–24, 1995.
- [10] T. Gollisch and M. Meister. Rapid neural coding in the retina with relative spike latencies. *science*, 319(5866):1108–1111, 2008.
- [11] T. Gollisch and M. Meister. Eye smarter than scientists believed: neural computations in circuits of the retina. *Neuron*, 65(2):150–164, 2010.
- [12] L. Grady. Random walks for image segmentation. *IEEE transactions on pattern analysis and machine intelligence*, 28(11):1768–1783, 2006.
- [13] A. Heitman, N. Brackbill, M. Greschner, A. Sher, A. M. Litke, and E. Chichilnisky. Testing pseudo-linear models of responses to natural scenes in primate retina. *bioRxiv*, page 045336, 2016.
- [14] H. Ishikane, M. Gangi, S. Honda, and M. Tachibana. Synchronized retinal oscillations encode essential information for escape behavior in frogs. *Nature neuroscience*, 8(8):1087–1095, 2005.

- [15] K. Koepsell, X. Wang, J. Hirsch, and F. T. Sommer. Exploring the function of neural oscillations in early sensory systems. *Frontiers in neuroscience*, 3:10, 2010.
- [16] K. Koepsell, X. Wang, V. Vaingankar, Y. Wei, Q. Wang, D. L. Rathbun, M. W. Usrey, J. Hirsch, and F. T. Sommer. Retinal oscillations carry visual information to cortex. *Frontiers in Systems Neuroscience*, 3:4, 2009.
- [17] S. W. Kuffler. Discharge patterns and functional organization of mammalian retina. *Journal of neurophysiology*, 16(1):37–68, 1953.
- [18] Y. Kuramoto. Chemical oscillations, waves, and turbulence, ser. *Springer Series in Synergetics*. Berlin/Heidelberg, Germany: Springer-Verlag, 19, 1984.
- [19] D. Martin, C. Fowlkes, D. Tal, and J. Malik. A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics. 2:416–423, 2001.
- [20] L. M. Martinez, M. Molano-Mazón, X. Wang, F. T. Sommer, and J. A. Hirsch. Statistical wiring of thalamic receptive fields optimizes spatial sampling of the retinal image. *Neuron*, 81(4):943–956, 2014.
- [21] R. H. Masland. Cell populations of the retina: the proctor lecture. *Investigative ophthalmology & visual science*, 52(7):4581–4591, 2011.
- [22] R. H. Masland. The tasks of amacrine cells. *Visual neuroscience*, 29(01):3–9, 2012.
- [23] J. Menzler and G. Zeck. Network oscillations in rod-degenerated mouse retinas. *Journal of Neuroscience*, 31(6):2280–2291, 2011.
- [24] S. Neuenschwander, M. Castelo-Branco, and W. Singer. Synchronous oscillations in the cat retina. *Vision research*, 39(15):2485–2497, 1999.
- [25] S. Neuenschwander and W. Singer. Long-range synchronization of oscillatory light responses in the cat retina and lateral geniculate nucleus. *Nature*, 379(6567):728–733, 1996.
- [26] M. E. Newman. Finding community structure in networks using the eigenvectors of matrices. *Physical review E*, 74(3):036104, 2006.
- [27] T. E. Ogden. The oscillatory waves of the primate electroretinogram. *Vision research*, 13(6):1059–IN4, 1973.
- [28] B. P. Ölveczky, S. A. Baccus, and M. Meister. Segregation of object and background motion in the retina. *Nature*, 423(6938):401–408, 2003.
- [29] J. W. Pillow, J. Shlens, L. Paninski, A. Sher, A. M. Litke, E. Chichilnisky, and E. P. Simoncelli. Spatio-temporal correlations and visual signalling in a complete neuronal population. *Nature*, 454(7207):995–999, 2008.

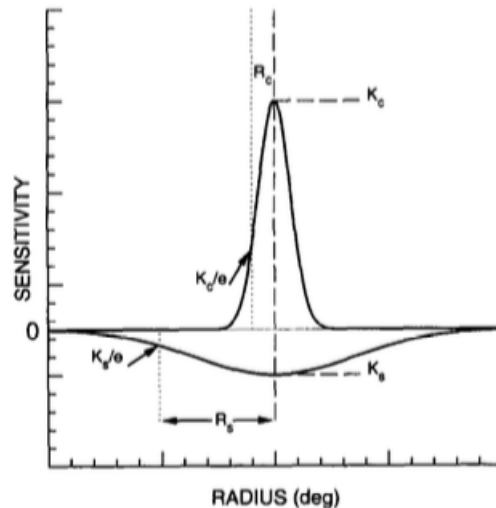
- [30] B. Roska and F. Werblin. Rapid global shifts in natural scenes block spiking in specific ganglion cell types. *Nature neuroscience*, 6(6):600–608, 2003.
- [31] D. L. Ruderman. The statistics of natural images. *Network: computation in neural systems*, 5(4):517–548, 1994.
- [32] S. Sarkar and K. L. Boyer. Quantitative measures of change based on feature organization: Eigenvalues and eigenvectors. *Computer vision and image understanding*, 71(1):110–136, 1998.
- [33] E. Schneidman, M. J. Berry, R. Segev, and W. Bialek. Weak pairwise correlations imply strongly correlated network states in a neural population. *Nature*, 440(7087):1007–1012, 2006.
- [34] O. Schwartz, J. W. Pillow, N. C. Rust, and E. P. Simoncelli. Spike-triggered neural characterization. *Journal of vision*, 6(4):13–13, 2006.
- [35] J. Shi and J. Malik. Normalized cuts and image segmentation. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 22(8):888–905, 2000.
- [36] F. S. Werblin. The retinal hypercircuit: a repeating synaptic interactive motif underlying visual function. *The Journal of physiology*, 589(15):3691–3702, 2011.
- [37] W. Yin, D. Goldfarb, and S. Osher. Total variation based image cartoon-texture decomposition. Technical report, DTIC Document, 2005.
- [38] F. Zenke, W. Gerstner, and S. Ganguli. The temporal paradox of hebbian learning and homeostatic plasticity. *Current opinion in neurobiology*, 43:166–176, 2017.

A Optimal Gaussian RF size

There are multiple independent sensor models to which we could compare the network models. We constrain our sensors to have access to relatively simple image features similar to those which the retina would encode. For comparison, we compute image segmentation using two independent sensor null models. The first uses raw image pixels and the second passes image pixels through Gaussian filters that mimic retinal ganglion cell (RGC) receptive fields (RFs). Center-surround RGC RFs are modelled by a difference-of-Gaussian filter with an excitatory center and inhibitory surround. Gaussian filters fit to the centers and surrounds of primate midget and parasol ganglion cells were observed to be strongly center dominant [9]. Thus the receptive field of an RGC can reasonably be modelled by a single excitatory Gaussian center to first approximation and the optimal Gaussian RF size reasonably matches average RGC RF sizes measured in primate retina.

In our simulations, the phase initialization of each individual oscillator as well as the connectivity strength between oscillators are both determined by the cell’s activation - that is, how closely incoming stimulus matches the filter that is defined as a cell’s receptive field. We began with the simplest receptive field model, each cell responding to the greyscale pixel intensity value at its location. Then, motivated by the biological fact that retinal receptive fields are spatially extended, we extended the receptive field model for each oscillating cell to be a localized Gaussian RF kernel. To determine the best Gaussian RF size (σ), we numerically explored a range of spread values and kept the one that provided best average segmentation performance across 500 image patches in the Berkeley Segmentation Dataset (BSDS) [19]. Segmentation performance was determined by F-measure calculated on the match between spatial gradients in phase maps output by network models and ground truth boundaries drawn by human subjects. Interestingly, we determined that a Gaussian RF kernel with $\sigma = 1$ pixel performed best empirically, improving the F-measure value by a modest but statistically significant 0.04 points over raw image pixels.

Motivated further by the excitatory and inhibitory center-surround nature of biological receptive fields in retina, we employ difference of gaussian (DoG) filters with parameters based on retinal physiology [9]. The Croner paper provides parameters fit to DoG receptive fields for M and P cells in primate retina for eccentricities ranging from 0 – 40° in its Table 1. In contrast with LGN center-surround cells [20], retinal receptive fields have very weak surrounds ($\sim 1/100^{th}$) compared with the strength of the center portion. From the many receptive field parameters fit to different cell types at different eccentricities in the primate retina, we distilled out 4 clusters that were different enough to test via simulations. In our simulations using DoG filters with P-avg and M-avg parameter values, we did not see image segmentation improvement over simple Gaussian filter with $\sigma = 1$.



| | R_c | R_s | K_s/K_c |
|-------|-------|-------|-----------|
| P-avg | 1 | 8 | 0.01 |
| P-40° | 3 | 13 | 0.06 |
| M-avg | 3 | 14.5 | 0.01 |
| M-40° | 5 | 12.5 | 0.025 |

Figure 17: **Primate center-surround RFs:** modeled as difference-of-Gaussians. Note: R_c and R_s in image pixels. Values are given for magnocellular projecting (P) and parvocellular projecting (M) cells averaged across all eccentricities (avg) and at the visual periphery (-40°) Image of measured retinal RF size from Croner 1995 [9]

Using a simple back-of-the-envelope visual angle calculation, illustrated in Fig. 18, and a few reasonable assumptions we approximate the size of retinal receptive field centers and surrounds in terms of image pixels for our models. The calculation goes as follows: Full images in the BSDS are 321×481 pixels and we assume that the displayed image size is $8.5'' \times 11''$. Given these assumptions, an image pixel is approximately $0.02''$ on a side. Next, we assume that the projection screen is placed $24''$ away from the eye. Then, the angle that a single pixel subtends on the retina is approximately 0.05° . Using this relation, we convert numbers provided in the Croner paper for retinal receptive field sizes into pixels and provide them in Fig. 17.

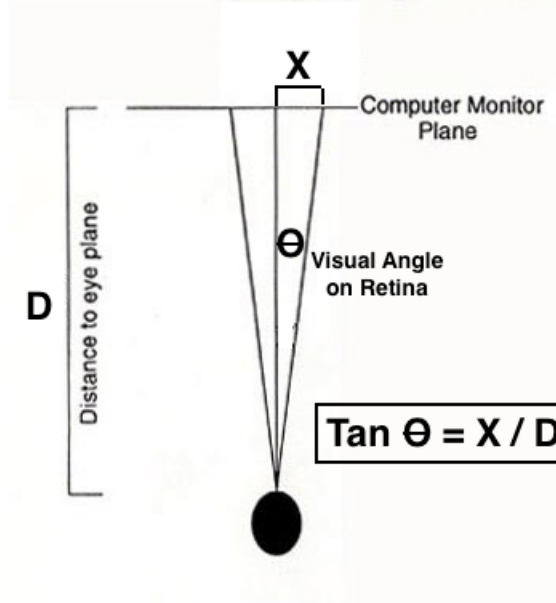


Figure 18: Visual angle calculation schematic

B Motivating modularity

B.1 Homogenization and Null Model as Expected Value of Weight

Most generally, an entry in the *modularity* matrix (Q_{ij}) is defined as the difference in weight between a pair of nodes in the actual network, characterized in the *adjacency matrix* (A_{ij}), and the expected value of that weight ($\mathbf{E}[A_{ij}]$) in a “homogenized network”, with connections between nodes made to reflect gross statistics of the network’s connectivity.

$$Q_{ij} = A_{ij} - \mathbf{E}[A_{ij}] \quad \text{with} \quad \mathbf{E}[A_{ij}] = \int A_{ij} p(A_{ij}) dA_{ij} \quad (7)$$

The expected value of weights is parameterized in the *null model* (N_{ij}) which is chosen to reflect the modeller’s knowledge of network structure and connectivity.

$$Q_{ij} = A_{ij} - N_{ij} \quad \text{where} \quad N_{ij} = \mathbf{E}[A_{ij}] \quad (8)$$

The null model is constrained only by two considerations. First, because the networks considered have undirected edges, both adjacency and null model matrices are symmetric, with $N_{ij} = N_{ji}$ and $A_{ij} = A_{ji}$. Second, it is axiomatically required that the total weight of edges in the null model are equal to the total weight of edges in the actual network because $Q = 0$ when all the vertices are

placed in the same partition. This leads to a normalizing constraint on the null model matrix,

$$\Sigma = \sum_{ij} A_{ij} = \sum_{ij} N_{ij} \quad (9)$$

where Σ is twice the total weight of edges in the network to account for double counting in the double sum over vertices (Note: $\sum_{ij} := \sum_i \sum_j$). Beyond these basic requirements, we are free to choose from many possible null models, each one containing a different number of parameters, requiring a different number of computations and capturing the expectation of edge weights at different levels of homogeneity by calculating different statistics on the adjacency matrix.

B.2 I.I.D. or Homogeneous Random Graph

The simplest null model, based on a Bernoulli or Erdos-Renyi random graph with weights allowed to take real values (i.e. are not constrained to be binary), assigns a single uniform expectation weight to all edges in the network, $\bar{A} = \frac{\Sigma}{n^2 - n}$, which is the average edge weight in the actual network. Note that n is the number of nodes in the network and $\binom{n}{2} = \frac{n^2 - n}{2}$ is the number of possible undirected edges that connect them with all-to-all connectivity, barring self-loops.

$$\mathbf{E}[A_{ij} | \frac{\Sigma}{n^2 - n}] = \int A_{ij} \cdot p(A_{ij} | \frac{\Sigma}{n^2 - n}) dA_{ij} = \int A_{ij} \delta(A_{ij} - c \frac{\Sigma}{n^2 - n}) dA_{ij} \quad (10)$$

$$N_{ij} = \mathbf{E}[A_{ij} | \frac{\Sigma}{n^2 - n}] = c \frac{\Sigma}{n^2 - n} \quad (11)$$

Solving for c by equation 9, we find

$$c = \frac{n - 1}{n}. \quad (12)$$

Combining the I.I.D. edge weight assumption with the constraint on total weight strength, we derive that the null model which assumes Bernoulli random graph connectivity patterns expects each weight in the network to take the following value.

$$N_{ij} = \frac{\Sigma}{n^2} \quad (13)$$

This is a very simple representation of the network which requires only a single number - the average edge weight across the entire network (\bar{A}), however it is inadequate to capture the structure in all but the simplest networks.

B.3 Independent-Vertex or Inhomogeneous Random Graph (N&G Modularity)

Relaxing the “identical” assumption of the I.I.D. graph null model, the “Independent-Vertex” model allows the expected value of each weight in the null model network to be different (inhomogeneous). The expected value of a weight between two nodes is the product of the degree of each of those nodes. This null model capture the expectation that two strongly connected nodes are more likely to be connected to one another and two nodes which are generally weakly connected are unlikely to be connected to one another. Specifically,

$$\mathbf{E}[A_{ij} | \frac{d_i}{n}, \frac{d_j}{n}] = \int A_{ij} p(A_{ij} | \frac{d_i}{n}, \frac{d_j}{n}) dA_{ij} = \int A_{ij} \delta(A_{ij} - c \frac{d_i}{n} \frac{d_j}{n}) dA_{ij} \quad (14)$$

$$N_{ij} = \mathbf{E}[A_{ij} | \frac{d_i}{n}, \frac{d_j}{n}] = c \frac{d_i}{n} \frac{d_j}{n} \quad \text{where} \quad d_i = \sum_{i=1}^n A_{ij} \quad (15)$$

where n is the number of vertices and d_i is the “degree” of node i or strength of connectivity from node i to all other nodes in the network, defined as the row (or equivalently column) sums of the adjacency matrix. Solving for c by equation 9, we find

$$c = \frac{n^2}{\Sigma} \quad (16)$$

making the full null model

$$N_{ij} = \frac{d_i d_j}{\Sigma}. \quad (17)$$

This requires n numbers or statistics calculated from the network to characterize the null model, namely the degree of each node. This is the model used by Newman [26] and works well finding community structure in networks with no inherent spatial layout or topography.

B.4 Line-Distance Dependent, Independent-Vertex Random Graph in 1D (Mod SKH Adj)

In networks with 1D spatial relationships, where each vertex is more likely or more strongly connected to nearby vertices than to distant vertices, the independent-vertex null model which just considers vertex degrees fails to capture this spatial structure and the modularity’s ability to find communities in such topographical networks suffers. The simplest spatial arrangement of nodes in a network is along a line in one dimension. Here, we can expand the vertex-independent null model to include a line-distance dependent ($b_{|i-j|}$) term which characterizes the expectation of a weight between nodes separated by a distance ($|i - j|$).

$$\begin{aligned} \mathbf{E}[A_{ij} | \frac{d_i}{n}, \frac{d_j}{n}, \frac{b_{|i-j|}}{n-|i-j|}] = & \\ & \int A_{ij} p(A_{ij} | \frac{d_i}{n}, \frac{d_j}{n}, \frac{b_{|i-j|}}{n-|i-j|}) dA_{ij} = \\ & \int A_{ij} \delta(A_{ij} - c \frac{d_i}{n} \frac{d_j}{n} \frac{b_{|i-j|}}{n-|i-j|}) dA_{ij} \end{aligned} \quad (18)$$

$$N_{ij} = \mathbf{E}[A_{ij} | \frac{d_i}{n}, \frac{d_j}{n}, \frac{b_{|i-j|}}{n-|i-j|}] = c \frac{d_i}{n} \frac{d_j}{n} \frac{b_{|i-j|}}{n-|i-j|} \quad (19)$$

where

$$d_i = \sum_{i=1}^n A_{ij} \quad \text{and} \quad b_{|i-j|} = \sum_{k=1}^{n-|i-j|} A_{k, k+|i-j|} \quad (20)$$

Solving for c by equation 9 yeilds

$$c = \frac{n^2 \Sigma}{\sum_{ij} (d_i d_j \frac{b_{|i-j|}}{n-|i-j|})} \quad (21)$$

and the full null model is

$$N_{ij} = \frac{d_i d_j \frac{b_{|i-j|}}{n-|i-j|} \Sigma}{\sum_{ij} (d_i d_j \frac{b_{|i-j|}}{n-|i-j|})} \quad (22)$$

where $\frac{d_i}{n}$ is the average weight from node i to other nodes in the network, and $\frac{b_{|i-j|}}{n-|i-j|}$ is the average weight between a pair of nodes separated by the distance $|i-j|$. Since nodes are arranged along a line, their separation distance in 1 dimensional space directly translates into distance from the diagonal in the adjacency matrix. Namely, the first off-diagonal contains weights between nodes separated by one distance unit, the second off diagonal by two units, and so on. This method requires $2n$ values computed from A to characterize the null model, the n normalized row (or column) sums and the n normalized diagonal sums. Although it is not entirely correct for networks arranged on a 2D grid, it can be used and yeilds better performance than the Independent-Vertex null model.

B.5 Grid-Distance Dependent, Independent-Vertex Random Graph in 2D (Mod SKH Euc)

A more correct null model for networks constructed from images admits the arrangement of nodes in a 2D lattice. The setup follows very closely the construction discussed above in the Line-Distance Dependent case with independent contributions from node degrees and from the connectivity-distance relationship

across the entire network. When nodes are arranged in a two dimensional grid, however, the relationship between distance in the network and location in the adjacency matrix is no longer simple to express mathematically, as in diagonal sums of \mathbf{A} in the 1D case. Fig. 19 below shows entries in the adjacency matrix representing the collection of edges separating pairs of nodes by the distance indicated in each pane in an 11x11 image patch.

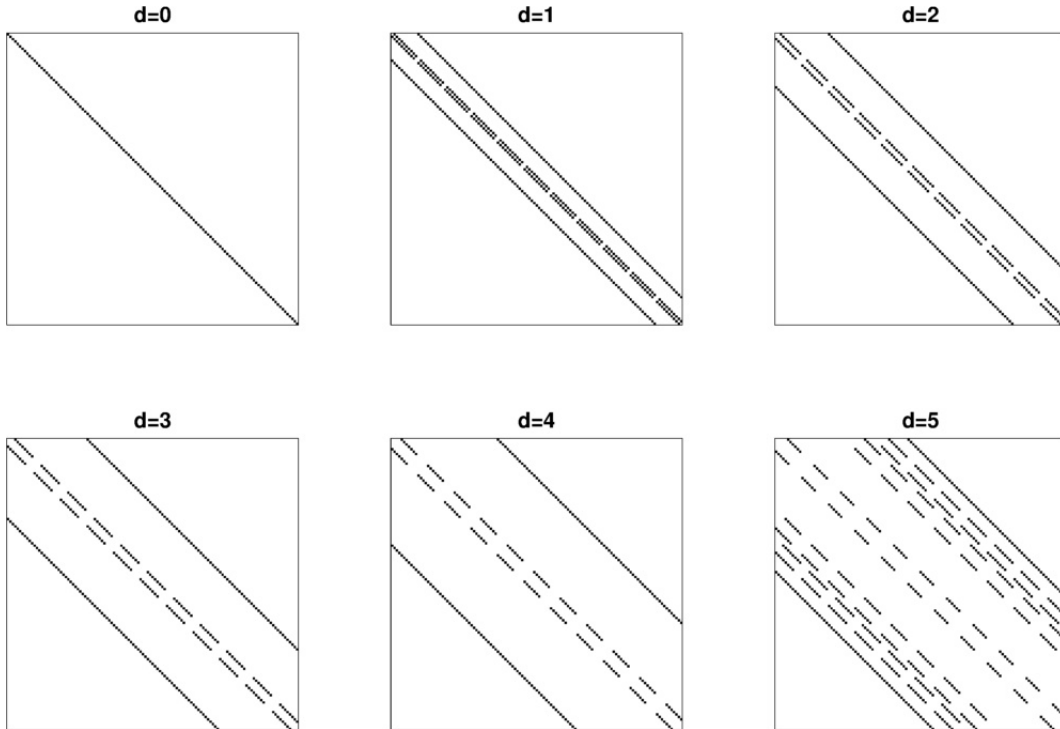


Figure 19: **Grid-Distance Dependence:** Distance mask in \mathbf{A} matrix: Elements within the adjacency matrix that are separated by distance $d = |r_i - r_j|$ in an 11x11 network arranged on a 2D lattice.

For all but $|r_i - r_j| = 0$, distances in the image plane translate into patterns in the adjacency matrix that are more complex than just off-diagonals. Note that each pattern includes some of the $|r_i - r_j|^{th}$ off-diagonal, with additional entries resulting from the way which the $n \times n$ image is rasterized to make to form the $n^2 \times n^2$ adjacency matrix. In our implementation, we do not attempt to express the $b_{|r_i - r_j|}$ term analytically, rather we algorithmically compute distances in the image plane and construct an adjacency matrix mask for each distance that we use to compute the distance-dependent average connectivity. Aside from difference in implementation, the motivation behind this model is identical to the 1D case. Here specifically,

$$\begin{aligned} \mathbf{E}[A_{ij} | \frac{d_i}{n}, \frac{d_j}{n}, \frac{b_{|r_i-r_j|}}{\#b_{|r_i-r_j|}}] = \\ \int A_{ij} p(A_{ij} | \frac{d_i}{n}, \frac{d_j}{n}, \frac{b_{|r_i-r_j|}}{\#b_{|r_i-r_j|}}) dA_{ij} = \\ \int A_{ij} \delta(A_{ij} - c \frac{d_i}{n} \frac{d_j}{n} \frac{b_{|r_i-r_j|}}{\#b_{|r_i-r_j|}}) dA_{ij} \end{aligned} \quad (23)$$

$$N_{ij} = \mathbf{E}[A_{ij} | \frac{d_i}{n}, \frac{d_j}{n}, \frac{b_{|r_i-r_j|}}{\#b_{|r_i-r_j|}}] = c \frac{d_i}{n} \frac{d_j}{n} \frac{b_{|r_i-r_j|}}{\#b_{|r_i-r_j|}} \quad (24)$$

where

$$d_i = \sum_{i=1}^n A_{ij} \quad (25)$$

and $b_{|r_i-r_j|}$ is implemented by masks illustrated in Fig. 19. Here, the $\#b_{|r_i-r_j|}$ term refers to the number of non-zero entries in the mask for the given distance. Since edges are undirected and \mathbf{A} is symmetric, the distance mask could also be implemented using the upper or lower triangular version of the adjacency matrix.

Solving for c by equation 9 yields

$$c = \frac{n^2 \Sigma}{\sum_{ij} (d_i d_j \frac{b_{|r_i-r_j|}}{\#b_{|r_i-r_j|}})} \quad (26)$$

and the full null model with the normalization constant is

$$N_{ij} = \frac{d_i d_j \frac{b_{|r_i-r_j|}}{\#b_{|r_i-r_j|}} \Sigma}{\sum_{ij} (d_i d_j \frac{b_{|r_i-r_j|}}{\#b_{|r_i-r_j|}})} \quad (27)$$

B.6 Temporal Modularity Null Model

While topographic modularity models are powerful tools for image segmentation, it is difficult to interpret how they could be implemented in retinal circuitry. The distance-dependent term $b_{|r_i-r_j|}$ requires that each edge in the network have access to global knowledge, namely the average edge weight across the entire network of all edges that span the same physical distance for the current input stimulus. However, the null model can be constructed with only local information if each neuron pair samples and stores the average edge weight between them over an ensemble of past stimuli. Hebbian plasticity in the ganglion-amacrine cell anatomical connectivity network could nicely account for such a computation.

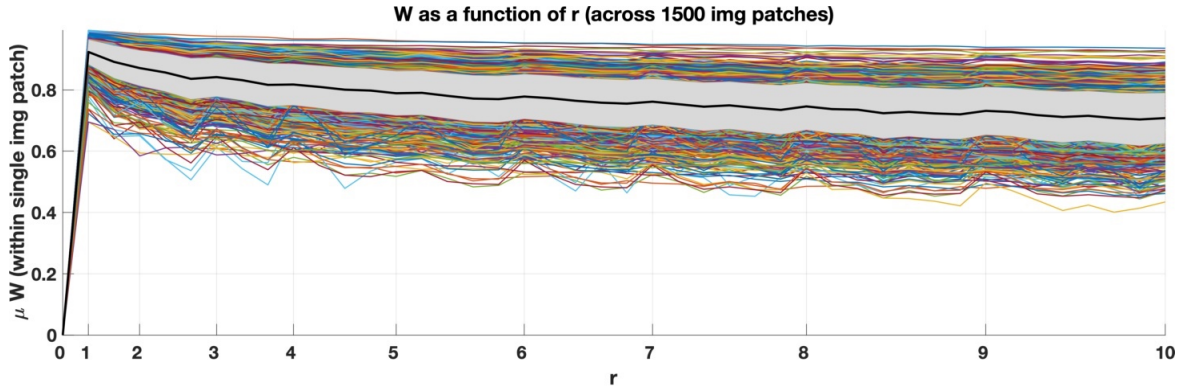


Figure 20: **Adjacency edge weight vs distance:** Average edge weight between node pairs in the adjacency matrix separated by distance r as a function of distance in image. Colored lines denote individual image patches and black line with grey error bars indicates μ and σ across 1500 image patches that are 50x50pixels.

Within a single scene or image, this spatial statistic can be converted to a local, temporal statistic via eye movements in a persistent scene if the timescale of plasticity is shorter than the scene duration [38]. For longer Hebbian timescales, the argument holds across an ensemble of natural scenes in so far as the distance-dependent feature similarity in single images is captured by an average across the ensemble. Pixel values in images of natural scenes have been shown to be much more highly correlated for nearby pairs of pixels than for distant pairs [4]. Fig. 20 shows the average weight in the Adjacency matrix across all node pairs i and j separated by a distance $r = |r_i - r_j|$ as a function of r , within single image patches as colored lines and the mean and standard deviation across an ensemble in black and grey.

A further advantage of a temporally sampled null model, beyond node degree and distance-dependence, is that *all* parameters describing the relationship between cells (such as cell types and direction) are trivially captured the cell pair itself is used to compute the null model. Thus the null model effectively controls for all influences to network connectivity other than image content, which is marginalized out over many samples across time. The temporal null model has not been explored in this work and is left for future development.