

---

# Unsupervised learning of features and object boundaries from local prediction

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 A visual system has to learn both which features to extract from images and how  
2 to group locations into (proto-)objects. Those two aspects are usually dealt with  
3 separately, although predictability is discussed as a cue for both. To incorporate  
4 features and boundaries into the same model, we model a layer of feature maps  
5 with a pairwise Markov random field model in which each factor is paired with an  
6 additional binary variable, which switches the factor on or off. Using one of two  
7 contrastive learning objectives, we can learn both the features and the parameters of  
8 the Markov random field factors from images without further supervision signals.  
9 The features learned by shallow neural networks based on this loss are local averages,  
10 opponent colors, and Gabor-like stripe patterns. Furthermore, we can infer  
11 connectivity between locations by inferring the switch variables. Contours inferred  
12 from this connectivity perform quite well on the Berkeley segmentation database  
13 (BSDS500) without any training on contours. Thus, computing predictions across  
14 space aids both segmentation and feature learning, and models trained to optimize  
15 these predictions show similarities to the human visual system. We speculate that  
16 retinotopic visual cortex might implement such predictions over space through  
17 lateral connections.

## 18 1 Introduction

19 A long-standing question about human vision is how representations initially be based on parallel  
20 processing of retinotopic feature maps can represent *objects* in a useful way. Most research on  
21 this topic has focused on computing later object-centered representations from the feature map  
22 representations. Psychology and neuroscience identified features that lead to objects being grouped  
23 together [37, 38], established feature integration into coherent objects as a sequential process [73], and  
24 developed solutions to the binding problem, i.e. ways how neurons could signal whether they represent  
25 parts of the same object [17, 57, 67, 72]. In computer vision, researchers also focused on how feature  
26 map representations could be turned into segmentations and object masks. Classically, segmentation  
27 algorithm were clustering algorithms operating on extracted feature spaces [2, 12, 13, 16, 66], and  
28 this approach is still explored with more complex mixture models today [74]. Since the advent of  
29 deep neural network models, the focus has shifted towards models that directly map to contour maps  
30 or semantic segmentation maps [21, 27, 39, 50, 65, 83], as reviewed in [54].

31 Diverse findings suggest that processing within the feature maps take object boundaries into account.  
32 For example, neurons appear to encode border ownership [34, 57, 63] and to fill in information across  
33 surfaces [40] and along illusory contours [23, 76]. Also, attention spreading through the feature  
34 maps seems to respect object boundaries [4, 59]. And selecting neurons that correspond to an object  
35 takes time, which scales with the distance between the points to be compared [35, 41]. Finally, a  
36 long history of psychophysical studies showed that changes in spatial frequency and orientation

37 content can define (texture) boundaries [e.g. 5, 45, 81]. In both human vision and computer vision,  
 38 relatively little attention has been given to these effects of grouping or segmentation on the feature  
 39 maps themselves.

40 Additionally, most theories for grouping and segmentation take the features in the original feature  
 41 maps as given. In human vision, these features are traditionally chosen by the experimenter [37,  
 42 73, 72] or are inferred based on other research [57, 63]. Similarly, computer vision algorithms used  
 43 off-the-shelf feature banks originally [2, 12, 13, 16, 66], and have recently moved towards deep neural  
 44 network representations trained for other tasks as a source for feature maps [21, 27, 39, 50, 65, 83].

45 Interestingly, predictability of visual inputs over space and time has been discussed as a solution for  
 46 both these limitations of earlier theories. Predictability has been used as a cue for segmentation since  
 47 the law of common fate of Gestalt psychology [37], and both lateral interactions in visual cortices and  
 48 contour integration respect the statistics of natural scenes [19, 20]. Among other signals like sparsity  
 49 [55] or reconstruction [36], predictability is also a well known signal for self-supervised learning  
 50 of features [80], which has been exploited by many recent contrastive learning [e.g. 15, 24, 29, 75]  
 51 and predictive coding schemes [e.g. 51, 52, 75] for self-supervised learning. However, these uses of  
 52 predictability for feature learning and for segmentation are usually studied separately.

53 Here, we propose a model that learns both features and segmentation without supervision. Predictions  
 54 between locations provide a self-supervised loss to learn the features, how to perform the prediction  
 55 and how to infer which locations should be grouped. Also, this view combines contrastive learning  
 56 [24, 75], a Markov random field model for the feature maps [46] and segmentation into a coherent  
 57 framework. We implement our model using some shallow architectures. The learned features  
 58 resemble early cortical responses and the object boundaries we infer from predictability align well  
 59 with human object contour reports from the Berkeley segmentation database (BSDS500 [2]). Thus,  
 60 retinotopic visual cortex might implement similar computational principles as we propose here.

## 61 2 Model

62 To explain our combined model of feature maps and their local segmentation information, we start  
 63 with a Gaussian Markov random field model [46] with pairwise factors. We then add a variable  
 64  $w \in \{0, 1\}$  to each factor that governs whether the factor enters the product or not. This yields a joint  
 65 distribution for the whole feature map and all  $w$ 's. Marginalizing out the  $w$ 's yields a Markov random  
 66 field with "robust" factors for the feature map, which we can use to predict feature vectors from the  
 67 vectors at neighboring positions. We find two contrastive losses based on these predictions that can  
 68 be used to optimize the feature extraction and the factors in the Markov random field model.

69 We model the distribution of  $k$ -dimensional feature maps  $\mathbf{f} \in \mathbb{R}^{k, m', n'}$  that are computed from input  
 70 images  $I \in \mathbb{R}^{c, m, n}$  with  $c = 3$  color channels (see Fig. 1 A & B). We use a Markov random field  
 71 model with pairwise factors, i.e. we define the probability of encountering a feature map  $\mathbf{f}$  with  
 72 entries  $f_i$  at locations  $i \in [1 \dots m'] \times [1 \dots n']$  as follows:

$$p(\mathbf{f}) \propto \prod_i \psi_i(f_i) \prod_{(i,j) \in N} \psi_{ij}(f_i, f_j), \quad (1)$$

73 where  $\psi_i$  is the local factor,  $N$  is the set of all neighboring pairs, and  $\psi_{ij}$  is the pairwise factor  
 74 between positions  $i$  and  $j$ <sup>1</sup>. We will additionally assume shift invariance, i.e. each point has the same  
 75 set of nearby relative positions in the map as neighbors,  $\psi_i$  is the same factor for each position, and  
 76 each factor  $\psi_{ij}$  depends only on the relative position of  $i$  and  $j$ .

77 We now add a binary variable  $w \in \{0, 1\}$  to each pairwise factor that encodes whether the factor  
 78 is 'active' ( $w = 1$ ) for that particular image (Fig. 1 C). To scale the probability of  $w = 1$  and  
 79  $w = 0$  relative to each other, we add a factor that scales them with constants  $p_{ij} \in [0, 1]$  and  $1 - p_{ij}$   
 80 respectively:

$$p(\mathbf{f}, \mathbf{w}) \propto \prod_i \psi_i(f_i) \prod_{(i,j) \in N} p_{ij}^{w_{ij}} (1 - p_{ij})^{1 - w_{ij}} \psi_{ij}(f_i, f_j)^{w_{ij}} \quad (2)$$

---

<sup>1</sup> $i$  and  $j$  thus have two entries each

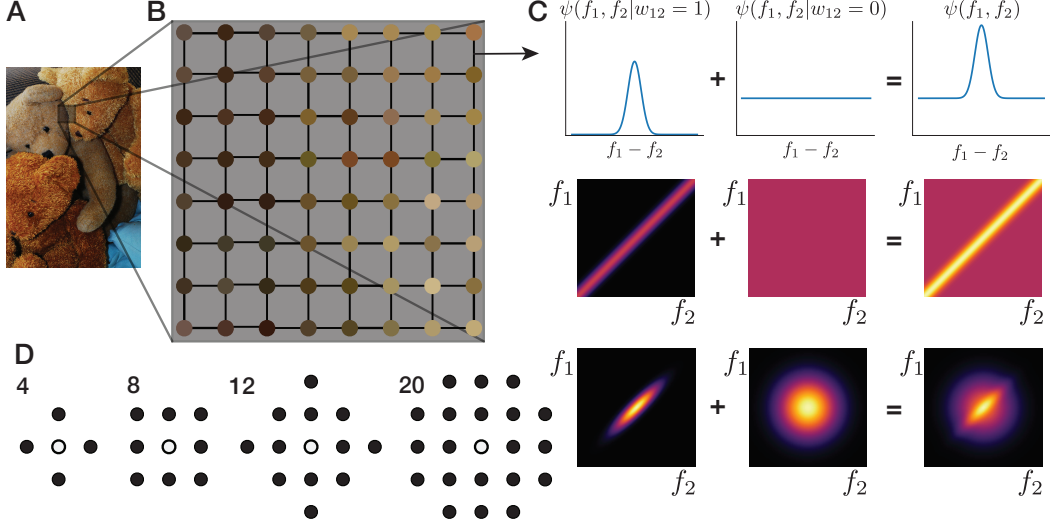


Figure 1: Illustration of our Markov random field model for the feature maps. **A**: An example input image. **B**: Feature map with 4 neighborhood connectivity and pixel color as the extracted feature. In the actual models, these feature maps are higher dimensional maps extracted by a convolutional neural network model. **C**: Illustration of the factor that links the feature vectors at two neighboring locations for a 1D feature. Top row: projection of the factor  $\psi_{ij}$  onto the difference between the features value  $f_i - f_j$ , showing the combination of a Gaussian around 0 and a constant function for the connection variable  $w_{ij}$  being 1 or 0 respectively. Middle row: 2D representation of the factor and its parts plotted against both feature values. Bottom row: Multiplication of the middle row with the standard normal factor for each position yielding the joint distribution of two isolated positions. **D**: Neighborhoods of different sizes used in the models, scaling from 4 to 20 neighbors for each location.

81 Finally, we assume that the factors are Gaussian and the feature vectors are originally normalized to  
 82 have mean 0 and variance 1:

$$p(\mathbf{f}, \mathbf{w}) = \frac{1}{Z_0} \mathcal{N}(\mathbf{f}, 0, \mathbf{I}) \prod_{(i,j) \in N} \frac{p_{ij}^{w_{ij}} (1 - p_{ij})^{1-w_{ij}}}{Z(w_{ij}, C_{ij})} \exp\left(-\frac{w_{ij}}{2} (f_i - f_j)^T C_{ij} (f_i - f_j)\right), \quad (3)$$

83 where  $Z_0$  is the overall normalization constant,  $N(\mathbf{f}, 0, \mathbf{I})$  is the density of a standard normal  
 84 distribution with  $k \times m' \times n'$  dimensions,  $C_{ij}$  governs the strength of the coupling in the form of a  
 85 precision matrix, which we will assume to be diagonal, and  $Z(w_{ij}, C_{ij})$  scales the distributions with  
 86  $w_{ij} = 0$  and  $w_{ij} = 1$  relative to each other.

87 We set  $Z(w_{ij}, C_{ij})$  to the normalization constant of the Gaussian with standard Gaussian factors  
 88 for  $f_i$  and  $f_j$  respectively. For  $w = 0$  this is just  $(2\pi)^{-k}$ , the normalization constant of a standard  
 89 Gaussian in  $2k$  dimensions. For  $w = 1$  we get:

$$Z(w_{ij} = 1, C_{ij}) = \int \int \exp\left(-\frac{1}{2} f_i^T f_i - \frac{1}{2} f_j^T f_j - \frac{1}{2} (f_i - f_j)^T C_{ij} (f_i - f_j)\right) df_i df_j \quad (4)$$

$$= (2\pi)^{-k} \det \begin{vmatrix} I + C_{ij} & C_{ij} \\ C_{ij} & I + C_{ij} \end{vmatrix}^{\frac{1}{2}} \quad (5)$$

$$= (2\pi)^{-k} \prod_l \sqrt{1 + 2c_{ll}} \quad (6)$$

90 which we get by computing the normalization constant of a Gaussian with the given precision and  
 91 then using the assumption that  $C_{ij}$  is a diagonal matrix with diagonal entries  $c_{ll}$ .

92 This normalization depends only on  $w$  and the coupling matrix  $C$  of the factor  $\psi_{ij}$  and thus induces a  
 93 valid probability distribution on the feature maps. Two points are notable about this normalization  
 94 though: First, once other factors also constrain  $f_i$  and/or  $f_j$ , this normalization will not guarantee  
 95  $p(w_{ij} = 1) = p_{ij}$ .<sup>2</sup> Second, the  $w_{ij}$  are not independent in the resulting distribution. For example, if  
 96 pairwise factors connect  $a$  to  $b$ ,  $b$  to  $c$  and  $a$  to  $c$  the corresponding  $w$  are dependent, because  $w_{ab} = 1$   
 97 and  $w_{bc} = 1$  already imply a smaller difference between  $f_a$  and  $f_c$  than if these factor were inactive,  
 98 which increases the probability for  $w_{ac} = 1$ .

## 99 2.1 Learning

100 To learn our model from data, we use a contrastive learning objective on the marginal likelihood  $p(\mathbf{f})$ .  
 101 To do so, we first need to marginalize out the  $w$ 's, which is fortunately simple, because each  $w$  affects  
 102 only a single factor:

$$p(\mathbf{f}) = \sum_{\mathbf{w}} p(\mathbf{f}, \mathbf{w}) = \frac{1}{Z_0} \mathcal{N}(\mathbf{f}, 0, \mathbf{I}) \prod_{(i,j) \in N} [p_{ij} \psi_{ij}(f_i, f_j) + (1 - p_{ij})] \quad (7)$$

103 Using this marginal likelihood directly for fitting is infeasible though, because computing  $Z_0$ , i.e.  
 104 normalizing this distribution is not computationally tractable.

105 We resort to contrastive learning to fit the unnormalized probability distribution [24], i.e. we optimize  
 106 discrimination from a noise distribution with the same support as the target distribution. Following  
 107 [75] we do not optimize the Markov random field directly, but optimize predictions based on the  
 108 model using features from other locations as the noise distribution. For this noise distribution, the  
 109 factors that depend only on a single location (the first product in (1)) will cancel. We thus ignore the  
 110  $\mathcal{N}(\mathbf{f}, 0, \mathbf{I})$  in our optimization and instead normalize the feature maps to mean 0 and unit variance  
 111 across each image. We define two alternative losses that make predictions for positions based on all  
 112 their neighbors or for a single factor respectively.

### 113 2.1.1 Position loss

114 The *position loss* optimizes the probability of the feature vector at each location relative to the  
 115 probability of randomly chosen other feature vectors from different locations and images:

$$l_{\text{pos}}(\mathbf{f}) = \sum_i \log \frac{p(f_i | f_j \forall j \in N(i))}{\sum_{i'} p(f_{i'} | f_j \forall j \in N(i))} \quad (8)$$

$$= \sum_i \sum_{j \in N(i)} \log \psi_{ij}(f_i, f_j) - \sum_i \log \left( \sum_{i'} \exp \left[ \sum_{j \in N(i)} \log \psi_{ij}(f_{i'}, f_j) \right] \right), \quad (9)$$

116 where  $N(i)$  is the set of neighbors of  $i$ .

117 This loss is consistent with the prediction made by the whole Markov random field, but is relatively  
 118 inefficient, because the predicted distribution  $p(f_i | f_j \forall j \in N(i))$  and the normalization constants for  
 119 these conditional distributions are different for every location  $i$ . Thus, the second term in equation  
 120 (9) cannot be reused across the locations  $i$ . Instead, we need to compute the second term for each  
 121 location separately, which requires a similar amount of memory as the whole feature representation  
 122 for each negative sample  $i'$  and each neighbor.

123 To enable a sufficiently large set of negative points  $i'$  with the available memory, we compute this loss  
 124 multiple times with few negative samples and sum the gradients. This trick saves memory, because  
 125 we can free the memory for the loss computation after each repetition. As the initial computation  
 126 of the feature maps is the same for all negative samples, we can save some computation for this  
 127 procedure by computing the feature maps only once. To propagate the gradients through this single  
 128 computation, we add up the gradients of the loss repetitions with regard to the feature maps and then  
 129 propagate this summed gradient through the feature map computation. This procedure does not save  
 130 computation time compared to the loss with many negative samples, as we still need to calculate the  
 131 evaluation for each position and each sample in the normalization set.

<sup>2</sup>Instead,  $p(w_{ij} = 1)$  will be higher, because other factors increase the precision for the feature vectors, which makes the normalization constants more similar.

132 **2.1.2 Factor loss**

133 The *factor loss* instead maximizes each individual factor for the correct feature vectors relative to  
 134 random pairs of feature vectors sampled from different locations and images:

$$l_{\text{fact}} = \sum_{i,j} \log \frac{\psi_{ij}(f_i, f_j)}{\sum_{i',j'} \psi_{ij}(f_{i'}, f_{j'})} \quad (10)$$

$$= \sum_{i,j} \log \psi_{ij}(f_i, f_j) - \sum_{i,j} \log \sum_{i',j'} \psi_{ij}(f_{i'}, f_{j'}), \quad (11)$$

135 where  $i, j$  index the correct locations and  $i', j'$  index randomly drawn locations, in our implementation  
 136 generated by shuffling the feature maps and taking all pairs that occur in these shuffled maps.

137 This loss does not lead to a consistent estimation of the MRF model, because the prediction  $p(f_i|f_j)$   
 138 should not be based only on the factor  $\psi_{ij}$ , but should include indirect effects as  $f_j$  also constrains  
 139 the other neighbors of  $i$ . Optimizing each factor separately will thus overaccount for information  
 140 that could be implemented in two factors. However, this loss has the distinct advantage that the same  
 141 noise evaluations can be used for all positions and images in a minibatch, which enables a much  
 142 larger number of noise samples and thus much faster convergence.

143 **2.1.3 Optimization**

144 We optimize all weights of the neural network used for feature extraction and the parameters of the  
 145 random field, i.e. the connectivity matrices  $C$  and the  $p_{ij}$  for the different relative spatial locations  
 146 simultaneously. As an optimization algorithm we use stochastic gradient descent with momentum.  
 147 Further details of the optimization can be found in the supplementary materials.

148 **2.2 Segmentation inference**

149 Computing the probability for any individual pair of locations  $(i, j)$  to be connected, i.e. computing  
 150  $p(w_{ij} = 1|\mathbf{f})$ , depends only on the two connected feature vectors  $f_i$  and  $f_j$ :

$$\frac{p(w_{ij} = 1|\mathbf{f})}{p(w_{ij} = 0|\mathbf{f})} = \frac{p_{ij}}{(1 - p_{ij})} \frac{Z(w_{ij} = 0, C_{ij})}{Z(w_{ij} = 1, C_{ij})} \exp(-(f_i - f_j)^T C_{ij} (f_i - f_j)) \quad (12)$$

151 This inference effectively yields a connectivity measure for each pair of neighboring locations, i.e. a  
 152 sparse connectivity matrix. Given that we did not apply any prior information enforcing continuous  
 153 objects or contours, the inferred  $w_{ij}$  do not necessarily correspond to a valid segmentation or set of  
 154 contours. Finding the best fitting contours or segmentation for given probabilities for the  $w$ s is an  
 155 additional process, which in humans appears to be an attention-dependent serial process [35, 63].

156 To evaluate the detected boundaries in computer vision benchmarks, we nonetheless need to convert  
 157 the connectivity matrix we extracted into a contour image. To do so, we use the spectral-clustering-  
 158 based globalization method developed by [2]. This method requires that all connection weights  
 159 between nodes are positive. To achieve this, we transform the log-probability ratios for the  $w_{ij}$  as  
 160 follows: For each image, we find the 30% quantile of the values, subtract it from all log-probability  
 161 ratios, and set all values below 0.01 to 0.01. We then compute the smallest eigenvectors of the graph  
 162 Laplacian as in graph spectral clustering. These eigenvectors are then transformed back into image  
 163 space and are filtered with simple edge detectors to find the final contours.

164 **3 Evaluation**

165 We implement 3 model types implementing feature extractions of increasing complexity in PyTorch  
 166 [56]:

167 **Pixel value model.** For illustrative purposes, we first apply our ideas to the rgb pixel values of an  
 168 image as features. This provides us with an example, where we can easily show the feature values  
 169 and connections. Additionally, this model provides an easy benchmark for all evaluations.

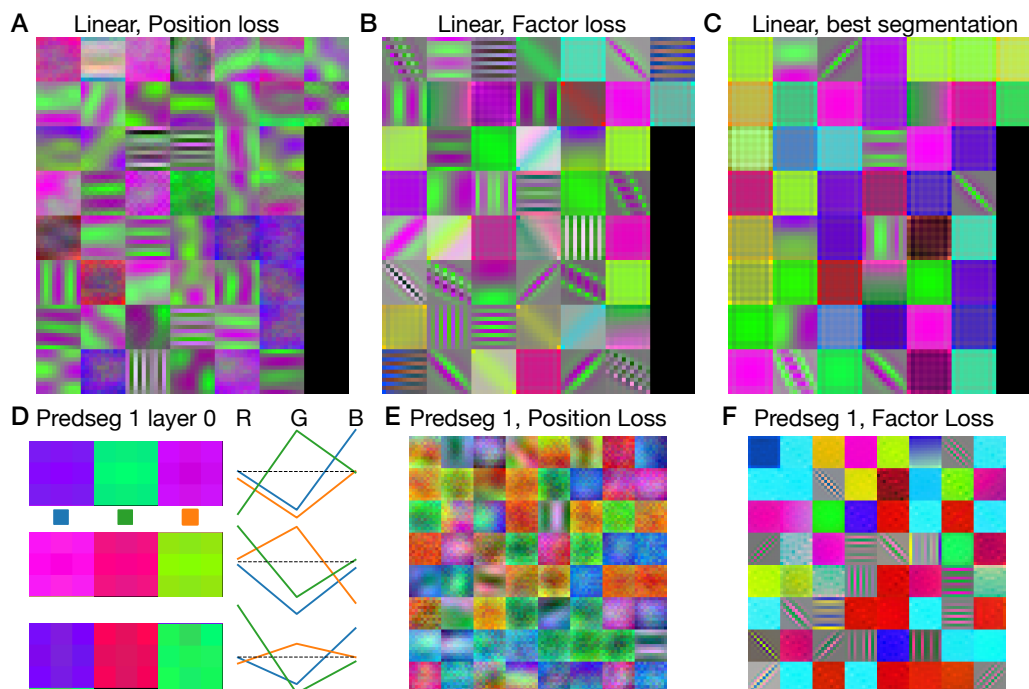


Figure 2: Example linear filter weights learned by our models. Each individual filter is normalized to minimum 0 and maximum 1. As weights can be negative even a zero weight can lead to a pixel having some brightness. For example, a number of channels load similarly on red and green across positions. Where these weights are positive the filter appears yellow and where the weights are negative filter appears blue, even if the blue channel has a zero weight. **A-C**: Feature weights learned by the linear model. **A**: Using the position loss. **B**: Using the factor loss. **C**: The weights of the model that leads to the best segmentation performance, i.e. the one shown in Figure 3. **D**: Weights of the first convolution in predseg1. Next to the filter shapes, which are nearly constant, we plot the average weight of each channel onto the three color channels of the image. **E** Predseg1 filters in the second convolution for a network trained with the position based loss. **F**: Predseg1 filters in the second convolution for a network trained with the factor based loss.

170 **Linear model.** As the simplest kind of model that allows learning features, we use a single convolutional deep neural network layer as our feature model. Here, we use 50  $11 \times 11$  linear features.

172 **Predseg1:** To show that our methods work for more complex architecture with non-linearities, we use a relatively small deep neural network with 4 layers (2 convolutional layers and 2 residual blocks with subsampling layers between them, see supplement for details).

175 For each of these architectures, we train 24 different networks with all combinations of the following settings: 4 different sizes of neighborhoods (4, 8, 12, or 20 neighbors, see Fig. 1D); 3 different noise levels (0, 0.1, 0.2) and the two learning objectives. As a training set, we used the unlabeled image set from MS COCO [48], which contains 123,404 color images with varying resolution. To enable batch processing, we randomly crop these images to  $256 \times 256$  pixel resolution, but use no other data augmentation (See supplementary information for further training details).

181 We want to evaluate whether our models learn meaningful features and segmentations. To do so, we first analyze the features in the first layers of our networks where we can judge whether features are representative of biological visual systems. In particular, we extract segmentations from our activations and evaluate those on the Berkeley Segmentation Dataset [2, BSDS500]

### 185 3.1 Learned features

186 **Linear Model** We first analyze the weights in our linear models (Fig 2 A-C). All instances learn local averages and Gabor-like striped features, i.e. spatial frequency and orientation tuned features



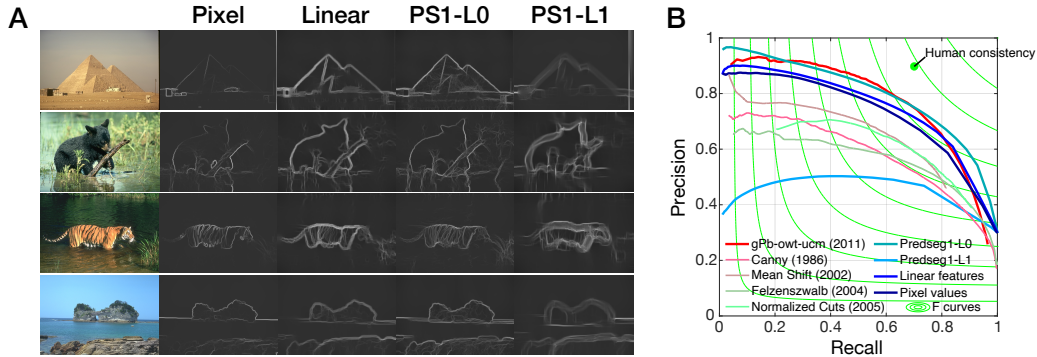


Figure 3: Contour detection results. **A**: Example segmentations from our models. **B**: Precision-recall curves for our models on the Berkeley segmentation dataset, with some other models for comparison as evaluated by [2]: gPb-owt-ucm, the final algorithm combining all improvements [2], Canny’s classical edge detector [7], the mean shift algorithm [12], Felzenszwalb’s algorithm [16] and segmentation based on normalized cuts [13]. For all comparison algorithms evaluations on BSDS were extracted from the figure by [2]

188 with limited spatial extent. These features clearly resemble receptive fields of neurons in primary  
 189 visual cortex. Additionally, there appears to be some preference for features that weight the red and  
 190 green color channels much stronger than the blue channel, similar to the human luminance channel,  
 191 which leads to the yellow-blue contrasts in the plots. There is some difference between the two  
 192 learning objectives though. The position based loss generally leads to lower frequency and somewhat  
 193 noisier features. This could either be due to the higher learning efficiency of the factor based loss, i.e.  
 194 the factor based loss is closer to convergence, or due to a genuinely different optimization goal.

195 **Predseg1** In Predseg1, we first analyze the layer 0 convolution (Fig. 2D), which has only 3 channels  
 196 with  $3 \times 3$  receptive fields, which we originally introduced as a learnable downsampling. This layer  
 197 consistently converges to applying near constant weights over space. Additionally, exactly one of  
 198 the channels has a non-zero mean (the 3rd, 1st and 3rd in Fig. 2D) and the other two take balanced  
 199 differences between two of the channels (red vs green and green vs. blue in the examples). This  
 200 parallels the luminance and opponent color channels of human visual perception.

201 In the second convolution, we observe a similar pattern of oriented filters and local averages as in  
 202 the linear model albeit in false color as the input channels are rotated by the weighting of the layer 0  
 203 convolution (Fig. 2E & F).

### 204 3.2 Contour detection

205 To evaluate whether the connectivity information extracted by our model corresponds to human  
 206 perceived segmentation, we extract contours from our models and compare them to contours reported  
 207 by humans for the Berkeley Segmentation database [2, 53]. This database contains human drawn  
 208 object boundaries for 500 natural images and is accompanied by methods for evaluating segmentation  
 209 models. Using the methods provided with the database, we compute precision-recall curves for each  
 210 model and use the best F-value (geometric mean of precision and recall) as the final evaluation metric.

211 As we had multiple models to choose from, we choose the models from each class that perform  
 212 best on the *training data* for our reports. For all models this was one of the models with the largest  
 213 neighborhood, i.e. using 20 neighbors, and the factor loss. It seems the factor loss performed  
 214 better simply due to its technical efficiency advantage as discussed above. Performance increases  
 215 monotonically with neighborhood size and Markov random field based approaches to semantic  
 216 segmentation also increased their performance with larger neighborhoods up to fully connected  
 217 Markov random fields [43, 8, 9]. We thus expect that larger neighborhoods could work even better.

218 Qualitatively, we observe that all our models yield sensible contour maps (see Fig. 3A). Even the  
 219 contours extracted from the pixel model yield sensible contours. Additionally, we note that the linear  
 220 model and Layer 1 of the predseg model tend to produce double contours, i.e. they tend to produce

Table 1: Numerical evaluation for various algorithms on the BSDS500 dataset. Precision and recall are only given for ODS, i.e. with a the threshold fixed across the whole dataset.

model	Recall	Precision	F(ODS)	F(OIS)	Area_PR
Deep Contour** [65]	–	–	0.76	0.78	0.80
HED** [83]	–	–	0.79	0.81	0.84
RCF** [50]	–	–	0.81	0.83	–
Deep Boundary** [39]	–	–	0.813	0.831	0.866
BDCN** [27]	–	–	0.83	0.84	0.89
Canny* [7]	–	–	0.60	0.63	0.58
Mean Shift* [12]	–	–	0.64	0.68	0.56
Felzenszwalb* [16]	–	–	0.61	0.64	0.56
Normalized Cuts* [13]	–	–	0.64	0.68	0.45
gPb-owt-ucm [2]	0.73	0.73	0.73	0.76	0.73
Pixel	0.73	0.66	0.69	0.69	0.73
linear	0.78	0.66	0.72	0.73	0.75
Predseg1-Layer 0	0.79	0.69	0.74	0.73	0.80
Predseg1-Layer 1	0.74	0.47	0.57	0.59	0.45

\*: Evaluation of these algorithms taken from [2]. \*\*: Supervised DNNs, evaluation taken from [27].

221 two contours on either side of the contour reported by human subjects with some area between them  
 222 connected to neither side of the contour.

223 Quantitatively, our models also perform well except for the deeper layers of Predseg 1 (Fig. 3B and  
 224 Table 1). The other models beat most hand-crafted contour detection algorithms that were tested  
 225 on this benchmark [7, 12, 13, 16] and perform close to the gPb-owt-ucm contour detection and  
 226 segmentation algorithm [2] that was the state of the art at the time. Layer-0 of Predseg 1 performs  
 227 best followed by the linear feature model and finally the pixel value model. Interestingly, the best  
 228 performing models seem to be mostly the local averaging models (cf. Fig. 2C). In particular, the  
 229 high performance of the first layer of Predseg 1 is surprising, because it uses only  $3 \times 3$  pixel local  
 230 color averages as features.

231 Since the advent of deep neural network models, networks trained to optimize performance on  
 232 image segmentation have reached much higher performance on the BSDS500 benchmark, essentially  
 233 reaching perfect performance up to human inconsistency [e.g. 27, 39, 49, 50, 65, 71, 83, see Table 1].  
 234 However, these models all require direct training on human reported contours and often use features  
 235 learned for other tasks. There are also a few deep neural network models that attempt unsupervised  
 236 segmentation [e.g. 10, 47, 82], but we were unable to find any that were evaluated on the contour  
 237 task of BSD500. The closest is perhaps the W-net [82], which used an autoencoder structure with  
 238 additional constraints and was evaluated on the segmentation task on BSDS500 performing slightly  
 239 better than gPb-owt-ucm.

## 240 4 Discussion

241 We present a model that can learn features and local segmentation information from images without  
 242 further supervision signals. This model integrates the prediction task used for feature learning and the  
 243 segmentation task into the same coherent probabilistic framework. This framework and the dual use  
 244 for the connectivity information make it seem sensible to represent this information. Furthermore,  
 245 the features learned by our models resemble receptive fields in the retina and primary visual cortex  
 246 and the contours we extract from connectivity information match contours drawn by human subject  
 247 fairly well, both without any training towards making them more human-like.

248 To improve biological plausibility, all computations in our model are local and all units are connected  
 249 to the same small, local set of other units throughout learning and inference, which matches early  
 250 visual cortex, in which the lateral connections that follow natural image statistics are implemented  
 251 anatomically [6, 31, 59, 70]. This in contrast to other ideas that require flexible pointers to arbitrary  
 252 locations and features [as discussed by 64] or capsules that flexibly encode different parts of the input  
 253 [14, 42, 61, 62]. Nonetheless, we employ contrastive learning objectives and backpropagation here,



254 for which we do not provide a biologically plausible implementations. However, there is currently  
255 active research towards biologically plausible alternatives to these algorithms [e.g. 32, 84].

256 Selecting the neurons that react to a specific object appears to rely on some central resource [72, 73]  
257 and to spread gradually through the feature maps [34, 35, 63]. We used a computer vision algorithm  
258 for this step, which centrally computes the eigenvectors of the connectivity graph Laplacian [2],  
259 which does not immediately look biologically plausible. However, a recent theory for hippocampal  
260 place and grid cells suggests that these cells compute the same eigenvectors of a graph Laplacian  
261 of a prediction network, albeit of a successor representation, i.e. of predictions of the animals state  
262 transitions [68, 69]. Thus, this might be an abstract description of an operation brains are capable of.  
263 In particular, earlier accounts that model the selection as a marker that spreads to related locations  
264 [e.g. 17, 58, 67] have some similarities with iterative algorithms to compute eigenvectors. Originally,  
265 phase coherence between the neurons encoding the same object was proposed [17, 57, 67], but a gain  
266 increase with object based attention [58] or a known random modulation is also sufficient to select a  
267 task relevant set of neurons [25, 26]. Regardless of the mechanistic implementation of the marker,  
268 connectivity information of the type our model extracts would be extremely helpful to explain the  
269 gradual spread of object selection.

270 Our implementation of the model is not fully optimized, as it is meant as a proof of concept. In  
271 particular, we did not optimize the architectures or training parameters of our networks for the  
272 task, like initialization, optimization algorithm, learning rate, or regularization. Presumably, better  
273 performance in all benchmarks could be reached by adjusting any or all of these parameters.

274 One possible next step for our model would be to train deeper architectures, such that the features  
275 could be used for complex tasks like object detection and classification. Contrastive losses like the  
276 one we use here are successfully applied for such pretraining purposes even for large scale tasks such  
277 as ImageNet [60] or MS Coco [48]. These large scale applications often require modifications for  
278 better learning though [11, 15, 22, 28, 29, 75]. For example: Image augmentations to explicitly train  
279 networks to be invariant to some image changes, prediction heads that allow more complex shapes  
280 for the predictions, and memory banks or other methods to decrease the reliance on many negative  
281 samples. Similar modifications might be necessary to apply our formulation to deeper architectures  
282 for pretraining purposes. For understanding human vision, this line of reasoning opens the exciting  
283 possibility that higher visual cortex could be explained based on similar principles, as representations  
284 from contrastive learning also yield high predictive power for these cortices [86].

285 The model we propose here is a probabilistic model of the feature maps. One implication of this  
286 is that we could also infer the feature values if they were not fixed based on the input. Thus, our  
287 model implies a pattern how neurons should combine their bottom-up inputs with predictions from  
288 nearby other neurons, once we include some uncertainty for the bottom-up inputs. In particular, the  
289 combination ought to take into account which nearby neurons react to the same object and which  
290 ones do not. Investigating this pooling could provide insights and predictions for phenomena that  
291 are related to local averaging like crowding for example [3, 18, 30, 77-79], where summary statistic  
292 models currently capture perceptual limitations best [3, 18, 78], but deviations from these predictions  
293 suggest that object boundaries change processing [30, 77, 79].

294 Another promising extension of our model would be processing over time, because predictions over  
295 time were found to be a potent signal for contrastive learning [15] and because coherent object motion  
296 is among the strongest grouping signals for human observers [38] and computer vision systems [85].  
297 Beside the substantial increases in processing capacity necessary to move to video processing instead  
298 of image processing, this step would require some extension of our framework to include object  
299 motion into the prediction. Nonetheless, including processing over time seems to be an interesting  
300 avenue for future research, especially because segmentation annotations for video are extremely  
301 expensive to collect such that unsupervised learning is particularly advantageous and popular in  
302 recent approaches [1, 33, 44].

303 This work aims to move us closer to understanding how human visual perception can take object  
304 structure into account in retinotopic feature map processing and may help us to build systems with  
305 similar capabilities in the future. We acknowledge that such technological progress can have unknown  
306 societal consequences, but we do not foresee specific negative consequences of this work.

307 **References**

- 308 [1] Nikita Araslanov, Simone Schaub-Meyer, and Stefan Roth. Dense Unsupervised Learning for  
309 Video Segmentation. *Advances in Neural Information Processing Systems*, 35:12, 2021.
- 310 [2] P Arbeláez, M Maire, C Fowlkes, and J Malik. Contour Detection and Hierarchical Image  
311 Segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 33(5):898–916,  
312 2011.
- 313 [3] B. Balas, L. Nakano, and R. Rosenholtz. A summary-statistic representation in peripheral vision  
314 explains visual crowding. *Journal of Vision*, 9(12):13–13, 2009.
- 315 [4] Daniel Baldauf and Robert Desimone. Neural mechanisms of object-based attention. *Science*,  
316 344(6182):424–427, 2014.
- 317 [5] Jacob Beck, Anne Sutter, and Richard Ivry. Spatial frequency channels and perceptual grouping  
318 in texture segregation. *Computer Vision, Graphics, and Image Processing*, 37(2):299–325,  
319 1987.
- 320 [6] Péter Buzás, Krisztina Kovács, Alex S. Ferecskó, Julian M.L. Budd, Ulf T. Eysel, and Zoltán F.  
321 Kisvárdy. Model-based analysis of excitatory lateral connections in the visual cortex. *The*  
322 *Journal of Comparative Neurology*, 499(6):861–881, 2006.
- 323 [7] John Canny. A computational approach to edge detection. *IEEE Transactions on pattern*  
324 *analysis and machine intelligence*, (6):679–698, 1986.
- 325 [8] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L. Yuille.  
326 Semantic image segmentation with deep convolutional nets and fully connected crfs. *arXiv*  
327 *preprint arXiv:1412.7062*, 2014.
- 328 [9] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L. Yuille.  
329 Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution,  
330 and fully connected crfs. *IEEE transactions on pattern analysis and machine intelligence*,  
331 40(4):834–848, 2017.
- 332 [10] Mickaël Chen, Thierry Artières, and Ludovic Denoyer. Unsupervised Object Segmentation by  
333 Redrawing. *Advances in Neural Information Processing Systems*, 33:12, 2019.
- 334 [11] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework  
335 for contrastive learning of visual representations. In *International conference on machine*  
336 *learning*, pages 1597–1607. PMLR, 2020.
- 337 [12] D. Comaniciu and P. Meer. Mean shift: A robust approach toward feature space analysis. *IEEE*  
338 *Transactions on Pattern Analysis and Machine Intelligence*, 24(5):603–619, 2002.
- 339 [13] T. Cour, F. Benezit, and Jianbo Shi. Spectral Segmentation with Multiscale Graph Decomposition.  
340 In *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*  
341 *(CVPR’05)*, volume 2, pages 1124–1131, San Diego, CA, USA, 2005. IEEE.
- 342 [14] Adrien Doerig, Lynn Schmittwilken, Bilge Sayim, Mauro Manassi, and Michael H. Herzog.  
343 Capsule networks as recurrent models of grouping and segmentation. *PLOS Computational*  
344 *Biology*, 16(7):e1008017, 2020.
- 345 [15] Christoph Feichtenhofer, Haoqi Fan, Bo Xiong, Ross Girshick, and Kaiming He. A large-scale  
346 study on unsupervised spatiotemporal representation learning. In *Proceedings of the IEEE/CVF*  
347 *Conference on Computer Vision and Pattern Recognition*, pages 3299–3309, 2021.
- 348 [16] Pedro F. Felzenszwalb and Daniel P. Huttenlocher. Efficient Graph-Based Image Segmentation.  
349 *International Journal of Computer Vision*, 59(2):167–181, 2004.
- 350 [17] Holger Finger and Peter König. Phase synchrony facilitates binding and segmentation of natural  
351 images in a coupled neural oscillator network. *Frontiers in Computational Neuroscience*, 7,  
352 2014.

- 353 [18] Jeremy Freeman and Eero P Simoncelli. Metamers of the ventral stream. *Nature Neuroscience*,  
354 14(9):1195–1201, 2011.
- 355 [19] Wilson S. Geisler and Jeffrey S. Perry. Contour statistics in natural images: Grouping across  
356 occlusions. *Visual Neuroscience*, 26(1):109–121, 2009.
- 357 [20] W.S. Geisler, J.S. Perry, B.J. Super, and D.P. Gallogly. Edge co-occurrence in natural images  
358 predicts contour grouping performance. *Vision Research*, 41(6):711–724, 2001.
- 359 [21] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for  
360 accurate object detection and semantic segmentation. In *Proceedings of the IEEE conference on  
361 computer vision and pattern recognition*, pages 580–587, 2014.
- 362 [22] Jean-Bastien Grill, Florian Strub, Florent Alché, Corentin Tallec, Pierre Richemond, Elena  
363 Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar,  
364 et al. Bootstrap your own latent—a new approach to self-supervised learning. *Advances in neural  
365 information processing systems*, 33:21271–21284, 2020.
- 366 [23] David H. Grosz, Robert M. Shapley, and Michael J. Hawken. Macaque VI neurons can signal  
367 ‘illusory’ contours. *Nature*, 365(6446):550–552, 1993.
- 368 [24] Michael Gutmann and Aapo Hyvarinen. Noise-contrastive estimation: A new estimation  
369 principle for unnormalized statistical models. *International Conference on Artificial Intelligence  
370 and Statistics (AISTATS)*, pages 297–304, 2010.
- 371 [25] Caroline Haimerl, Douglas A Ruff, Marlene R Cohen, Cristina Savin, and Eero P Simoncelli.  
372 Targeted comodulation supports flexible and accurate decoding in V1. *bioRxiv : the preprint  
373 server for biology*, 2021.
- 374 [26] Caroline Haimerl, Cristina Savin, and Eero Simoncelli. Flexible information routing in neural  
375 populations through stochastic comodulation. *Advances in Neural Information Processing  
376 Systems*, 32, 2019.
- 377 [27] Jianzhong He, Shiliang Zhang, Ming Yang, Yanhu Shan, and Tiejun Huang. Bi-Directional  
378 Cascade Network for Perceptual Edge Detection. In *2019 IEEE/CVF Conference on Computer  
379 Vision and Pattern Recognition (CVPR)*, pages 3823–3832, Long Beach, CA, USA, 2019. IEEE.
- 380 [28] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum Contrast for  
381 Unsupervised Visual Representation Learning. In *2020 IEEE/CVF Conference on Computer  
382 Vision and Pattern Recognition (CVPR)*, pages 9726–9735, Seattle, WA, USA, 2020. IEEE.
- 383 [29] Olivier J Hénaff, Aravind Srinivas, Jeffrey De Fauw, Ali Razavi, Carl Doersch, and S M Ali  
384 Eslami. Data-Efficient Image Recognition with Contrastive Predictive Coding. *Proceedings of  
385 the 37th International Conference on Machine Learning (PMLR)*, page 119, 2020.
- 386 [30] Michael H. Herzog, Bilge Sayim, Vitaly Chicherov, and Mauro Manassi. Crowding, grouping,  
387 and object recognition: A matter of appearance. *Journal of Vision*, 15(6):5, 2015.
- 388 [31] Jonathan J Hunt, William H Bosking, and Geoffrey J Goodhill. Statistical structure of lateral  
389 connections in the primary visual cortex. *Neural Systems & Circuits*, 1(1):3, 2011.
- 390 [32] Bernd Illing, Jean Ventura, Guillaume Bellec, and Wulfram Gerstner. Local plasticity rules can  
391 learn deep representations using self-supervised contrastive predictions. *Advances in Neural  
392 Information Processing Systems*, 35:15, 2021.
- 393 [33] Allan A Jabri, Andrew Owens, and Alexei A Efros. Space-Time Correspondence as a Contrastive  
394 Random Walk. *Advances in Neural Information Processing Systems*, 34:16, 2020.
- 395 [34] Danique Jeurissen, Matthew W. Self, and Pieter R. Roelfsema. Surface reconstruction, figure-  
396 ground modulation, and border-ownership. *Cognitive neuroscience*, 4(1):50–52, 2013.
- 397 [35] Danique Jeurissen, Matthew W. Self, and Pieter R. Roelfsema. Serial grouping of 2D-image  
398 regions with object-based attention in humans. *Elife*, 5:e14320, 2016.

- 399 [36] Diederik P. Kingma and Max Welling. Auto-Encoding Variational Bayes. *arXiv:1312.6114 [cs,*  
400 *stat]*, 2014.
- 401 [37] K Koffka. Principles of gestalt psychology. 1935.
- 402 [38] Wolfgang Köhler. Gestalt psychology. *Psychologische Forschung*, 31(1):XVIII–XXX, 1967.
- 403 [39] Iasonas Kokkinos. Pushing the Boundaries of Boundary Detection using Deep Learning.  
404 *arXiv:1511.07386 [cs]*, 2016.
- 405 [40] Hidehiko Komatsu. The neural mechanisms of perceptual filling-in. *Nature Reviews Neuro-*  
406 *science*, 7(3):220–231, 2006.
- 407 [41] Iliia Korjoukov, Danique Jeurissen, Niels A. Kloosterman, Josine E. Verhoeven, H. Steven  
408 Scholte, and Pieter R. Roelfsema. The Time Course of Perceptual Grouping in Natural Scenes.  
409 *Psychological Science*, 23(12):1482–1489, 2012.
- 410 [42] Adam R. Kosiorek, Sara Sabour, Yee Whye Teh, and Geoffrey E. Hinton. Stacked Capsule  
411 Autoencoders. *arXiv:1906.06818 [cs, stat]*, 2019.
- 412 [43] Philipp Krähenbühl and Vladlen Koltun. Efficient Inference in Fully Connected CRFs with  
413 Gaussian Edge Potentials. *arXiv:1210.5644 [cs]*, 2012.
- 414 [44] Zihang Lai, Erika Lu, and Weidi Xie. MAST: A Memory-Augmented Self-Supervised Tracker.  
415 *CVPR*, pages 6479–6488, 2020.
- 416 [45] Michael S Landy and James R Bergen. Texture segregation and orientation gradient. *Vision*  
417 *research*, 31(4):679–691, 1991.
- 418 [46] SZ Li. *Markov Random Field Modeling in Computer Vision*. Springer Science & Business  
419 Media, 2012.
- 420 [47] Qinghong Lin, Weichan Zhong, and Jianglin Lu. Deep Superpixel Cut for Unsupervised Image  
421 Segmentation. In *2020 25th International Conference on Pattern Recognition (ICPR)*, pages  
422 8870–8876, Milan, Italy, 2021. IEEE.
- 423 [48] Tsung-Yi Lin, Michael Maire, Serge Belongie, Lubomir Bourdev, Ross Girshick, James Hays,  
424 Pietro Perona, Deva Ramanan, C. Lawrence Zitnick, and Piotr Dollár. Microsoft COCO:  
425 Common Objects in Context. *arXiv:1405.0312 [cs]*, 2015.
- 426 [49] Drew Linsley, Junkyung Kim, Alekh Ashok, and Thomas Serre. Recurrent neural circuits for  
427 contour detection. In *International Conference on Learning Representations*, 2020.
- 428 [50] Yun Liu, Ming-Ming Cheng, Xiaowei Hu, Kai Wang, and Xiang Bai. Richer Convolutional  
429 Features for Edge Detection. In *Proceedings of the IEEE Conference on Computer Vision and*  
430 *Pattern Recognition*, pages 3000–3009, 2017.
- 431 [51] William Lotter, Gabriel Kreiman, and David Cox. Deep Predictive Coding Networks for Video  
432 Prediction and Unsupervised Learning. *arXiv:1605.08104 [cs, q-bio]*, 2017.
- 433 [52] William Lotter, Gabriel Kreiman, and David Cox. A neural network trained to predict fu-  
434 ture video frames mimics critical properties of biological neuronal responses and perception.  
435 *arXiv:1805.10734 [cs, q-bio]*, 2018.
- 436 [53] D. Martin, C. Fowlkes, D. Tal, and J. Malik. A database of human segmented natural images  
437 and its application to evaluating segmentation algorithms and measuring ecological statistics.  
438 In *Proc. 8th Int’l Conf. Computer Vision*, volume 2, pages 416–423, July 2001.
- 439 [54] Shervin Minaee, Yuri Y. Boykov, Fatih Porikli, Antonio J Plaza, Nasser Kehtarnavaz, and  
440 Demetri Terzopoulos. Image Segmentation Using Deep Learning: A Survey. *IEEE Transactions*  
441 *on Pattern Analysis and Machine Intelligence*, pages 1–1, 2021.
- 442 [55] Bruno A Olshausen and David J Field. Emergence of simple-cell receptive field properties by  
443 learning a sparse code for natural images. *Nature*, 381(6583):607–609, 1996.

- 444 [56] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan,  
445 Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas  
446 Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy,  
447 Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style,  
448 high-performance deep learning library. In *Advances in Neural Information Processing Systems*  
449 32, pages 8024–8035. 2019.
- 450 [57] Alina Peter, Cem Uran, Johanna Klön-Lipok, Rasmus Roese, Sylvia van Stijn, William Barnes,  
451 Jarrod R Dowdall, Wolf Singer, Pascal Fries, and Martin Vinck. Surface color and predictability  
452 determine contextual modulation of V1 firing and gamma oscillations. *eLife*, 8:e42101, 2019.
- 453 [58] Pieter R Roelfsema. Cortical algorithms for perceptual grouping. *Annu. Rev. Neurosci.*, 29:203–  
454 227, 2006.
- 455 [59] Pieter R Roelfsema, Victor AF Lamme, and Henk Spekreijse. Object-based attention in the  
456 primary visual cortex of the macaque monkey. *Nature*, 395(6700):376–381, 1998.
- 457 [60] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng  
458 Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei.  
459 ImageNet Large Scale Visual Recognition Challenge. *arXiv:1409.0575 [cs]*, 2015.
- 460 [61] Sara Sabour, Nicholas Frosst, and Geoffrey E. Hinton. Dynamic Routing Between Capsules.  
461 *arXiv:1710.09829 [cs]*, 2017.
- 462 [62] Sara Sabour, Andrea Tagliasacchi, Soroosh Yazdani, Geoffrey E. Hinton, and David J. Fleet.  
463 Unsupervised part representation by Flow Capsules. *arXiv:2011.13920 [cs]*, 2021.
- 464 [63] Matthew W. Self, Danique Jeurissen, Anne F. van Ham, Bram van Vugt, Jasper Poort, and  
465 Pieter R. Roelfsema. The Segmentation of Proto-Objects in the Monkey Primary Visual Cortex.  
466 *Current Biology*, 29(6):1019–1029, 2019.
- 467 [64] Michael N Shadlen and J Anthony Movshon. Synchrony unbound: A critical evaluation of the  
468 temporal binding hypothesis. *Neuron*, 24(1):67–77, 1999.
- 469 [65] Wei Shen, Xinggang Wang, Yan Wang, Xiang Bai, and Zhijiang Zhang. DeepContour: A deep  
470 convolutional feature learned by positive-sharing loss for contour detection. In *2015 IEEE*  
471 *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3982–3991, Boston,  
472 MA, USA, 2015. IEEE.
- 473 [66] Jianbo Shi and Jitendra Malik. Normalized Cuts and Image Segmentation. *IEEE Transactions*  
474 *On Pattern Analysis and Machine Intelligence*, 22(8):18, 2000.
- 475 [67] Wolf Singer and Charles M Gray. Visual feature integration and the temporal correlation  
476 hypothesis. *Annual review of neuroscience*, 18(1):555–586, 1995.
- 477 [68] Kimberly L Stachenfeld, Matthew Botvinick, and Samuel J Gershman. Design Principles of  
478 the Hippocampal Cognitive Map. *Advances in Neural Information Processing Systems*, page 9,  
479 2014.
- 480 [69] Kimberly L Stachenfeld, Matthew M Botvinick, and Samuel J Gershman. The hippocampus as  
481 a predictive map. *Nature Neuroscience*, 20(11):1643–1653, 2017.
- 482 [70] Dan D Stettler, Aniruddha Das, Jean Bennett, and Charles D Gilbert. Lateral Connectivity and  
483 Contextual Interactions in Macaque Primary Visual Cortex. *Neuron*, 36(4):739–750, 2002.
- 484 [71] Zhuo Su, Wenzhe Liu, Zitong Yu, Dewen Hu, Qing Liao, Qi Tian, Matti Pietikäinen, and Li Liu.  
485 Pixel Difference Networks for Efficient Edge Detection. In *Proceedings of the IEEE/CVF*  
486 *International Conference on Computer Vision (ICCV)*, pages 5117–5127, 2021.
- 487 [72] Anne Treisman. The binding problem. *Current opinion in neurobiology*, 6(2):171–178, 1996.
- 488 [73] Anne M. Treisman and Garry Gelade. A Feature-Integration Theory of Attention. *Cognitive*  
489 *Psychology*, 12:97–136, 1980.

- 490 [74] Jonathan Vacher, Claire Launay, and Ruben Coen-Cagli. Flexibly regularized mixture models  
491 and application to image segmentation. *Neural Networks*, 149:107–123, 2022.
- 492 [75] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation Learning with Contrastive  
493 Predictive Coding. *arXiv:1807.03748 [cs, stat]*, 2019.
- 494 [76] R. von der Heydt, E. Peterhans, and G. Baumgartner. Illusory Contours and Cortical Neuron  
495 Responses. *Science*, 224(4654):1260–1262, 1984.
- 496 [77] Thomas S. A. Wallis, Matthias Bethge, and Felix A. Wichmann. Testing models of peripheral  
497 encoding using metamerism in an oddity paradigm. *Journal of Vision*, 16(2):4, 2016.
- 498 [78] Thomas S. A. Wallis, Christina M. Funke, Alexander S. Ecker, Leon A. Gatys, Felix A.  
499 Wichmann, and Matthias Bethge. A parametric texture model based on deep convolutional  
500 features closely matches texture appearance for humans. *Journal of Vision*, 17(12):5–5, 2017.
- 501 [79] Thomas SA Wallis, Christina M Funke, Alexander S Ecker, Leon A Gatys, Felix A Wichmann,  
502 and Matthias Bethge. Image content is more important than Bouma’s Law for scene metamers.  
503 *eLife*, 8:e42512, 2019.
- 504 [80] Laurenz Wiskott and Terrence J. Sejnowski. Slow Feature Analysis: Unsupervised Learning of  
505 Invariances. *Neural Computation*, 14(4):715–770, 2002.
- 506 [81] S Sabina Wolfson and Michael S Landy. Discrimination of orientation-defined texture edges.  
507 *Vision research*, 35(20):2863–2877, 1995.
- 508 [82] Xide Xia and Brian Kulis. W-net: A deep model for fully unsupervised image segmentation.  
509 *arXiv preprint arXiv:1711.08506*, 2017.
- 510 [83] Saining Xie and Zhuowen Tu. Holistically-Nested Edge Detection. In *Proceedings of the IEEE  
511 International Conference on Computer Vision*, pages 1395–1403, 2015.
- 512 [84] Yuwen Xiong, Mengye Ren, and Raquel Urtasun. LoCo: Local Contrastive Representation  
513 Learning. *Advances in Neural Information Processing Systems*, 34:12, 2020.
- 514 [85] Charig Yang, Hala Lamdouar, Erika Lu, Andrew Zisserman, and Weidi Xie. Self-supervised  
515 Video Object Segmentation by Motion Grouping. *arXiv:2104.07658 [cs]*, 2021.
- 516 [86] Chengxu Zhuang, Siming Yan, Aran Nayebi, Martin Schrimpf, Michael C. Frank, James J.  
517 DiCarlo, and Daniel L. K. Yamins. Unsupervised neural network models of the ventral visual  
518 stream. *Proceedings of the National Academy of Sciences*, 118(3):e2014196118, 2021.



519 **Checklist**

- 520 1. For all authors...
- 521 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
522 contributions and scope? [Yes]
- 523 (b) Did you describe the limitations of your work? [Yes]
- 524 (c) Did you discuss any potential negative societal impacts of your work? [Yes]
- 525 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
526 them? [Yes]
- 527 2. If you are including theoretical results...
- 528 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 529 (b) Did you include complete proofs of all theoretical results? [N/A]
- 530 3. If you ran experiments...
- 531 (a) Did you include the code, data, and instructions needed to reproduce the main ex-  
532 perimental results (either in the supplemental material or as a URL)? [Yes] In the  
533 supplementary material.
- 534 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
535 were chosen)? [Yes] In the supplementary material.
- 536 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
537 ments multiple times)? [No] This would require undue amounts of computation for  
538 results, which we interpret only qualitatively anyway.
- 539 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
540 of GPUs, internal cluster, or cloud provider)? [Yes] In the supplementary material.
- 541 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 542 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 543 (b) Did you mention the license of the assets? [N/A]
- 544 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 545
- 546 (d) Did you discuss whether and how consent was obtained from people whose data you’re  
547 using/curating? [N/A]
- 548 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
549 information or offensive content? [N/A]
- 550 5. If you used crowdsourcing or conducted research with human subjects...
- 551 (a) Did you include the full text of instructions given to participants and screenshots, if  
552 applicable? [N/A]
- 553 (b) Did you describe any potential participant risks, with links to Institutional Review  
554 Board (IRB) approvals, if applicable? [N/A]
- 555 (c) Did you include the estimated hourly wage paid to participants and the total amount  
556 spent on participant compensation? [N/A]