

Retinotopic scaffolding of high-level vision

Nicholas M. Blauch^{1,2,*}, Marlene Behrmann^{1,3,4}, and David C. Plaut^{1,3}

¹*Neuroscience Institute, Carnegie Mellon University*

²*Department of Psychology, Harvard University*

³*Department of Psychology, Carnegie Mellon University*

⁴*Department of Ophthalmology, University of Pittsburgh*

*Corresponding author: nblauch@gmail.com

June 16, 2025

Abstract

Functional specialization within high-level vision is reflected in the topographic organization of the ventral temporal cortex (VTC). The presence and consistent locations of small areas responding selectively to particular visual categories – such as faces and scenes – has led to proposals of innate domain-specific modules. However, such proposals do not easily explain other aspects of an apparently multi-scale topographic organization of high-level visual features in VTC. Computational models have recently accounted for the presence of domain-selective areas and other facets of topographic organization from a basic optimization process with local topographic pressures, such as locally constrained connectivity, but fail to account for the consistent location of category-selectivity. In the current work, we extend a recent computational model to demonstrate how this consistency may emerge from wiring constraints external to VTC, focusing on the role of retinotopically organized early visual representations. After training several random initializations of the model, we find consistent global topographic organization, with face- and scene-selectivity emerging on opposite ends of a medial-lateral gradient corresponding to eccentricity bias, similar to human VTC. As in human VTC, the eccentricity-biased topography persists across viewing sizes under sufficiently broad viewing bias distributions, suggesting that it is a learned bias for efficiently organizing representations proximal to the most useful inputs, rather than merely an explicit retinotopic code. Abolishing the retinotopic constraint abolishes topographic consistency, but not topographic organization. Our work suggests that the organization of high-level visual cortex may emerge from domain-relevant interactions between viewing biases and task demands with an innate retinotopic scaffold. More generally, we suggest that both local and global connectivity constraints interact with representation learning to produce mature cortical organization: local constraints pressure the system towards smooth organization, whereas long-range constraints encourage a consistent global layout.

1 Introduction

Neuroscientists have long sought to parcellate the human brain into a discrete set of canonical brain areas involved in particular mental functions or neural computations. Early work parcellated the brain anatomically, based on differences in cytoarchitectural properties (Brodmann, 1909). More recent approaches aim to localize regions *functionally* within individual subjects, based on patterns of responses to stimuli and/or tasks, allowing for individual differences in precise anatomical localization (Kanwisher, 2010). Even still, some approaches opt to combine structural and functional features into multimodal parcellations (Glasser et al., 2016). This discrete areal view has provided a useful lens on the large-scale organization of the brain. However, in addition to this view, many cases of smooth gradations of functional organization have been found, including maps of retinotopic space (Daniel and Whitteridge, 1961) and local orientation pinwheel topography (Hubel and Wiesel, 1969) in primary visual cortex, large-scale feature tuning in higher-level visual areas (Konkle and Caramazza, 2013; Bao et al., 2020; Yao et al., 2023), and large-scale semantic maps that cover much of the human cerebral cortex (Huth et al., 2016). How can we understand both the discrete and graded forms of organization in terms of basic principles?

1.1 Optimization and wiring constraints

One hypothesis is that cortical organization emerges under some optimization procedure in combination with powerful biological constraints. Biological brains have evolved to increase the fitness of organisms, and this has led to powerful learning algorithms that allow for behavioral optimization within the lifespan of an organism. This optimization procedure can be seen as encompassing both evolutionary and ontogenetic factors. The optimization framework has provided strong explanatory power in computational cognitive neuroscience (Richards et al., 2019), owing to advances in deep learning systems, such that artificial neural networks can be optimized to exhibit human-level behavioral performance on several key cognitive and perceptual tasks.

A fundamental constraint on the cortical organization of brain functions is the spatial cost of wiring connections between neurons. To illustrate this, consider a simplified human brain containing 100 billion neurons placed on a spherical cortical

surface, with axons of $0.1 \mu\text{m}$ cross-sectional radius. Nelson and Bower (1990) found that this would require a sphere radius of 10km to have sufficient volume to fit all of the axons – $\approx 140,000$ x the linear size of our brains (70mm radius). Even under a more realistic scenario of neurons embedded *within* the sphere¹ (Behrmann and Plaut, 2020), the required radius is ≈ 8.6 km for full connectivity; assuming a realistic proportion of connectivity, with each neuron connecting to 10^4 others, still requires a radius of ≈ 5.55 m. Thus, sparse connectivity is not sufficient to fully alleviate the cost; connections must also be distributed primarily locally in order to minimize the wiring cost to biologically reasonable levels. This perspective has been termed the wiring minimization or wiring economy principle (Ramón y Cajal et al., 1899; Koulakov and Chklovskii, 2001; Chen et al., 2006). Phenomena as diverse as orientation pinwheel topography in V1 (Koulakov and Chklovskii, 2001) to the layout of neurons in *C. elegans* (Chen et al., 2006) can all be explained in terms of arranging neurons so as to minimize wiring costs given particular connectivity patterns.

1.2 The organization of high-level visual cortex

Here, we apply the wiring minimization principle to better understand the organization of high-level visual cortex, expanding upon our earlier computational modeling work (Blauch et al., 2022; Plaut and Behrmann, 2011). The ventral temporal cortex – the primary site of high-level visual representations in humans – contains spatial clusters of neurons that exhibit preferential activation in response to the viewing of certain visual categories, across a large degree of variation in the precise visual inputs. Notably, these clusters emerge in consistent locations across subjects (Kanwisher et al., 1997; Epstein and Kanwisher, 1998; Chao et al., 1999; Grill-Spector and Weiner, 2014). Considering the ecological importance of the tasks thought to be implemented in these clusters – facial recognition, spatial navigation, place and landmark recognition – it is not unreasonable to surmise that these structures evolved as domain-specific modules for particularly relevant tasks (Kanwisher, 2010; Fodor, 1983; Mahon and Caramazza, 2011).

However, the modular notion of VTC has been repeatedly challenged by alternative accounts of more generic forms of organization. Ishai et al. (1999) and Haxby et al. (2001) noted the presence of broadly distributed information patterns across VTC, where weak responses to a category in non-selective parts of VTC were sufficient to decode whether a subject was viewing that or another non-selective category; using invasive techniques in non-human primates, distributed information was later confirmed in precisely localized face-selective neuronal population activity (Meyers et al., 2015). Although the information for non-preferred categories was found to be weaker than that of the more selective category (Spiridon and Kanwisher, 2002; Meyers et al., 2015), it nevertheless suggested a less modular and more graded organization of function that could be understood in terms of basic principles of cooperation, competition, and wiring efficiency in a flexible learning system (Plaut and Behrmann, 2011; Behrmann and Plaut, 2015). Key to this idea are two broad sets of findings detailing topographic organization in VTC beyond the category-level: 1) spatially-organized tuning to low-level retinotopic information, and 2) spatially-organized tuning to high-level visual features.

First, Levy et al. (2001); Hasson et al. (2002) demonstrated that localized category selectivity could be systematically related to broad retinotopic biases in VTC, with activation to faces and words emerging in an area preferring foveal inputs, and activation to scenes located in an area preferring peripheral inputs. This led to the *eccentricity bias* theory, which accounted for category-selective organization in terms of early retinotopic organization and category-specific viewing biases. Critical to this theory is the large-scale eccentricity organization of early visual areas (Daniel and Whitteridge, 1961; Sereno et al., 1995), in which the fovea is overrepresented compared to the periphery, and eccentricity and polar angle form the two organizing dimensions. An implication of this eccentricity-based organization of early visual areas is that stimuli that extend into the periphery will be represented more medially along the ventral surface compared to stimuli restricted to the center of gaze. As a result, any consistent viewing biases for particular semantic categories will impose a retinotopic *eccentricity bias* on the organization of that category, under pressure to minimize wiring costs in the progression from early visual cortex to high-level visual cortex. Indeed, real-world viewing of scenes stimulates the full visual field, whereas faces, objects and words are often viewed at small sizes and require the high acuity of central gaze for accurate recognition in real world tasks. The eccentricity biased organization of VTC has been replicated using more sophisticated population receptive field (pRF) mapping methods (Silson et al., 2016, 2021; Finzi et al., 2021). Critically, retinotopic organization appears to be present at birth (Arcaro and Livingstone, 2017), developing *in utero* through gradients of axon guidance molecules connecting retinal ganglion cells (RGCs) to V1 via the lateral geniculate nucleus (LGN), as well as through activity-dependent mechanisms operating over spontaneous retinal waves (McLaughlin and O’Leary, 2005; Huberman et al., 2008). This has led to the idea that retinotopy serves as a “proto-architecture” on which high-level visual cortical organization is scaffolded (Hasson et al., 2002; Arcaro and Livingstone, 2017; Groen et al., 2022).

Second, category selectivity has been shown to be systematically related to broader feature tuning (Konkle and Caramazza, 2013), with large-scale topographic preferences for animacy and real-world object size that encompass smaller category-selective regions; whereas face and body selectivity lie in a broad zone preferentially responsive to animate information, scene selectivity is found within a broad zone preferentially responsive to large inanimate objects. Replacing visual stimuli with *texforms* – stimuli designed to preserve mid-level visual features of images while eliminating the semantic content – Long et al. (2018) demonstrated preserved selectivity, suggesting that mid-level visual features, rather than explicitly semantic information, may be at the core of VTC organization. Thus, rather than representing animacy and real-world size explicitly, these large-scale dimensions may represent their visual correlates in mid-level features. In line with this idea, Bao et al. (2020) and Yao et al. (2023) have demonstrated that an organization of abstract “object space” dimensions is present in VTC, in which category selective clusters emerge naturally. Whether this organization is fully

¹Given that the average distance between two points within a sphere of radius r is given as $\bar{d} = \frac{36}{35}r$, we solve for the radius of the sphere whose volume is filled with $m = n^2p$ axons of cross-sectional radius r_c connecting n neurons with probability p , ignoring spatial overlap: $V = \frac{4}{3}\pi r^3 = n^2p\pi\bar{d}r_c^2$. This yields $r = r_c n \sqrt{\frac{27}{35}p}$

experience-dependent is not fully determined, but it appears to arise from the natural statistics of images, as the object space dimensions are directly derived from the representations in deep neural networks trained for object recognition. Relatedly, in humans, the use of naturalistic movie stimuli and sophisticated representational alignment techniques (e.g., hyperalignment) has demonstrated that category-selective topography can be predicted from high-dimensional cortical responses to movies, in line with a more generic organization of content encompassing but not solely defined by category-selectivity (Haxby et al., 2011, 2020). The high-dimensional organization of human visual cortex has been shown to follow precise mathematical trends, with both the variance explained and the spatial scale decaying with a power law of the component number (Gauthaman et al., 2024, 2025).

1.3 The current approach

While the retinotopic proto-architecture and generic-object-space perspectives have been explored in isolation as theories regarding the organization of high-level visual cortex, here, we suggest that they can be generically cast under the ideas of wiring minimization and task optimization. In previous work, we developed a computational model to account for the topographic organization of high-level visual cortex. By training a deep learning model to perform to an expert level in recognizing different faces, objects, and scenes, while reducing its wiring cost within high-level visual processing areas, we found that highly domain-specialized organization emerged in the context of more generic self-organized topographic organization (Blauch et al., 2022). Several additional computational models have recently been developed that account for the topographic organization of ventral temporal cortex in generic terms, via optimization in a deep learning model (Keller and Welling, 2021; Doshi and Konkle, 2023; Zhang et al., 2024; Margalit et al., 2024; Qian et al., 2024; Deb et al., 2025; Lu et al., 2025). Notably, Doshi and Konkle (2023) and Keller and Welling (2021) studied the relationship between animacy and object size with category-selectivity, finding a systematic relationship as predicted by human fMRI data. Although each of the topographic models is subtly different, they can all be considered as "local constraint" models, per the terminology of Mahon (2022), where the key driver of functional organization is local to high-level vision, rather than dependent on the long-range interactions of high-level vision with broader brain networks.

These local-constraint models account for the systematic relationship between category selectivity and generic high-level feature tuning. However, they have left unexplained the global organization of VTC, which is the key property of the eccentricity bias theory. In this paper, we incorporate the eccentricity bias theory within the broad Interactive Topographic Network (ITN) framework that combines wiring and task optimization to explain replicable and mature topographic visual organization (Blauch et al., 2022). We demonstrate that a simple implementation of connectivity constraints with retinotopically organized mid-level inputs constrains the global layout of topographic model layers, yielding an organization that is remarkably similar to human VTC in its global category selectivity for objects, faces, and scenes. This allows the topography of models to be analyzed at the group level via anatomical alignment, yielding strong group-level consistency in topographic organization, as in the human brain. Critically, we demonstrate that the eccentricity-biased topography is *learned* from domain-specific viewing biases, yielding a strong degree of invariance to size after learning, similar to the findings of human empirical data (Park et al., 2024), but dependent on the degree of spatial invariance exposed to the model through its viewing biases. Our results provide strong computational support for the idea that retinotopic connectivity constraints shape the global layout of topographic organization in VTC. More broadly, we suggest that eccentricity bias is merely one example of a long-range connectivity constraint that shapes the optimization of topographic organization in VTC.

2 Methods

We extend the Interactive Topographic Network framework (ITN; Blauch et al., 2022) to test whether additionally constraining connectivity between posterior VTC with early retinotopic areas, in combination with domain-biased viewing conditions, produces the global layout of a single hemisphere of human VTC. As the focus of this computational model is the human rather than macaque brain, we depart from the nomenclature of Blauch et al. (2022), referring to model areas as "VTC" areas rather than "IT" areas; however, they are functionally the same. The architecture of the model is diagrammed in Figure 1A.

2.1 Task and wiring optimization

The key idea of our topographic network approach is that biological neural networks must minimize wiring costs while maximizing behavioral performance. Wiring cost can be quantified with a loss function \mathcal{L}_w that decreases as wiring cost is reduced, and behavioral performance can be quantified with a loss function \mathcal{L}_t that decreases as errors are reduced. We then compute a global loss function \mathcal{L} , where λ is a hyper-parameter controlling the relative weighting on the wiring cost:

$$\mathcal{L} = \mathcal{L}_t + \lambda \mathcal{L}_w \quad (1)$$

We set $\lambda_w = 0.1$, unless indicated otherwise. Both the task and wiring costs are flexible. Here, we focus on visual recognition in a supervised learning setting, and thus use the cross-entropy loss. Given one-hot labels y_{ic} – where each trial i is associated with a 1 at the index corresponding to category c , and 0 elsewhere – and predictions \hat{y}_{ic} , the cross entropy loss is given as:

$$\mathcal{L}_t = - \sum_{i=1}^N \sum_{c=1}^C y_{ic} \log(\hat{y}_{ic}) \quad (2)$$

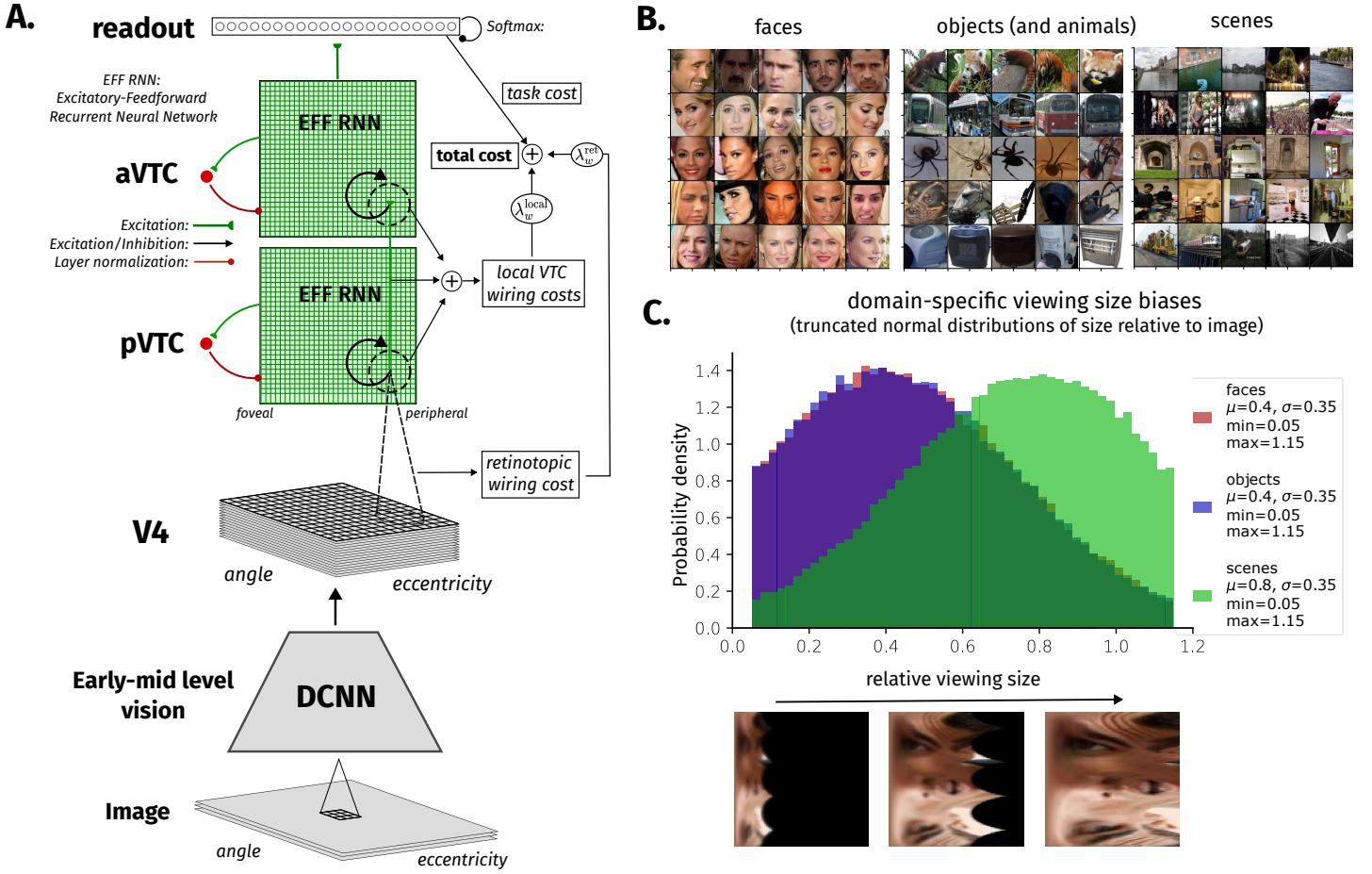


Figure 1: Modeling eccentricity bias with connectivity constraints, task optimization, and viewing biases. **A.** Architecture of the computational model and wiring cost computations. **B.** Example images from face, object/animal, and scene datasets. Each row contains 5 example images of a given category, with 5 categories shown per domain. **C.** Illustration of domain-specific viewing size biases. Faces and objects are presented at medium sizes with identical distributions, while scenes are presented at a larger size. Below the plot, we show an example face image at three relative viewing sizes (0.2, 0.6, 1.0) to demonstrate the polar coordinate images and relative viewing sizes.

157 We use the same wiring cost explored in our previous work that introduced the ITN (Blauch et al., 2022), modified from
 158 the earlier work of Jacobs and Jordan (1992). Specifically, given the Euclidean distance matrix $D^{(a,b)}$ between all pairs of
 159 units in areas a and b , and the corresponding weight matrix $W^{(a,b)}$ we used the following cost function:

$$\mathcal{L}_w^{(a,b)} = \sum_{i,j} \frac{\left(D_{ij}^{(a,b)}\right)^2 \left(W_{ij}^{(a,b)}\right)^2}{1 + \left(W_{ij}^{(a,b)}\right)^2} \quad (3)$$

160 We refer to this wiring cost as $\mathcal{L}_w^{\text{local}}$, to contrast with the retinotopic input wiring cost discussed in 2.3. The entire
 161 network is trained end-to-end from initial random weights, allowing the viewing biases to influence the representations
 162 of the entire network. We begin with a learning rate of 0.01 and reduce it 5 times by a factor of 10 upon plateau of the
 163 validation error; after the fifth learning rate decay, the next validation error plateau determined the end of training, up to a
 164 maximum of 100 epochs. Stochastic gradient descent with momentum ($\rho = 0.9$) and $L2$ weight decay ($\lambda = 0.0001$) was
 165 used, with batch size of 256 on a single GPU.

2.2 Model architecture

167 As in the ITN, we adopt a three-stage model made up of an encoder, high-level topographic areas, and a readout layer. A
 168 modified ResNet-18 encoder using reduced strides and pooling is used in order to allow for an output feature map size of
 169 9x9 for a 112x112 image, serving as the "V4" inputs to pVTC, with V4 as the final stage of the encoder. We implement two
 170 topographic layers of the excitatory-feedforward recurrent neural network (EFF-RNN) ITN variant (Blauch et al., 2022),
 171 which was shown to exhibit the major hallmarks of topographic organization in ITN models, while demanding substantially
 172 less memory and being simpler to analyze than fully separated E/I networks. This model circuit is similar to a standard

173 or "vanilla" RNN, but restricts the feedforward connectivity sign to be positive, reflecting the dominance of excitatory
174 connectivity in communication between different cortical areas (Crick et al., 1986; Laszlo and Plaut, 2012).

175 2.2.1 RNN equations

176 Where $\mathbf{x}^{(a)}$ is the vector of pre-activation activities in area a of IT, $\mathbf{r}^{(a)}$ is the vector of post-activation activities in area
177 a , $\mathbf{b}^{(a)}$ is the vector of baseline activities in area a , α is the discrete neuronal time constant, and $W^{(a,b)}$ is the matrix of
178 weights from area a to area b , we have the discrete time update as:

$$179 \mathbf{x}_t^{(a)} = (1 - \alpha)\mathbf{x}_{t-1}^{(a)} + \alpha \left(W^{(a,a)}\mathbf{r}_{t-1}^{(a)} + [W^{(a-1,a)}]_+ \mathbf{r}_{t-1}^{(a-1)} + \mathbf{b}^{(a)} \right) \quad (4)$$

180 The activation function $\mathbf{r}_t^{(a)} = [\mathbf{x}_t^{(a)}]_+$ is positive rectification, also called a Rectified Linear Unit (ReLU). Additionally,
note that the feedforward weights are restricted to be positive $[W^{(a-1,a)}]_+$.

181 2.2.2 Layer normalization

182 As in standard ITN models, layer normalization (Ba et al., 2016) is used to stabilize training (Blauch et al., 2022), without
183 the trainable scaling parameter that is sometimes used (see Ba et al., 2016, for more details). Where $\mu(\mathbf{x})$ is the mean of \mathbf{x} ,
184 and $\sigma(\mathbf{x})$ is the standard deviation of \mathbf{x} , and \mathbf{b} is the learned bias term (moved outside of the layer normalization), the
185 layer-normalized activities are given as:

$$\begin{aligned} z_t &= \frac{\mathbf{x}_t - \mu(\mathbf{x}_t)}{\sigma(\mathbf{x}_t)} + \mathbf{b} \\ \mathbf{r}'_t &= [z_t]_+ \end{aligned}$$

186 Incorporating layer normalization into our update equation yields the update equation:

$$187 \mathbf{x}_t^{(a)} = (1 - \alpha)\mathbf{z}_{t-1}^{(a)} + \alpha \left(W^{(a,a)}\mathbf{r}'_{t-1}^{(a)} + [W^{(a-1,a)}]_+ \mathbf{r}'_{t-1}^{(a-1)} \right) \quad (5)$$

188 2.2.3 Time constant

189 For computational convenience in training models end-to-end with greater retinotopic resolution – both of which increase
190 the computational demands on training – we train only on 5 time steps in our models here (in contrast to 20 time steps
191 used in previous work; Blauch et al., 2022). Thus, we opt to set $\alpha = 1$, meaning that the networks' activity would decay
192 exactly to zero in the absence of feedforward and recurrent connectivity, the standard choice in machine learning RNN
applications. This yields the simplified RNN update:

$$193 \mathbf{x}_t^{(a)} = W^{(a,a)}\mathbf{r}'_{t-1}^{(a)} + [W^{(a-1,a)}]_+ \mathbf{r}'_{t-1}^{(a-1)} \quad (6)$$

194 2.2.4 Connection noise

195 As in previous work (Blauch et al., 2022), to approximate axon-specific variability in instantaneous firing rate (Cipollini
196 and Cottrell, 2013), we apply multiplicative noise on the individual connections between neurons. This helps to regularize
197 the activations in the network, encouraging a more distributed representation that aids the formation of topography across
198 a range of models. Noise is sampled independently from a Gaussian distribution \mathcal{N} centered at 0 with variance σ^2 at each
199 time step of each trial, and is squashed by a sigmoidal function $S(x) = \frac{2}{1+e^{-x}}$, ensuring that the sign of each weight is not
200 changed and each magnitude does not change by more than 100%. Thus, the noisy weight matrix $W_n^{(a,b)}$ from area a to
area b on a given trial and time step is:

$$201 W_n^{(a,b)} = S(\mathcal{N}(0, \sigma)) * W^{(a,b)} \quad (7)$$

202 2.2.5 Readout

203 An excitatory-only linear readout is applied to the final (anterior) VTC layer. Given its weights W_r and the rectified final
layer (aVTC) activations \mathbf{r}'_t , this layer yields the vector of output predictions $\hat{\mathbf{y}}_t$ at time t as:

$$204 \hat{\mathbf{y}}_t = [W_r]_+ \mathbf{r}'_{t-1} \quad (8)$$

205 There are 300 output units, corresponding to 100 face identities, and 100 object categories, and 100 scene categories
(described below).

206 2.3 Retinotopic constraints

207 2.3.1 Wiring costs

208 Retinotopic constraints are modeled by adding a spatial cost on feedforward connections from V4 into posterior VTC
209 (pVTC) of the model. Whereas VTC layers are non-convolutional, and their units are thus simply placed on a 2D grid,
210 layers before VTC (such as V4) are convolutional and can be understood as 3D blocks of feature maps. Focused on the
211 influence of retinotopy on high-level vision, and lacking a principled manner to organize feature channels and retinotopic
212 position jointly, we opt to model only the spatial component of V4 organization as a wiring constraint in feedforward
213 connectivity to pVTC. The spatial position of each V4 unit is thus determined solely by its feature map location – not its
214 feature index; given these spatial positions, spatial connectivity costs are then computed as for the other layers. Specifically,
215 we can consider the wiring cost as a sum over local VTC wiring costs $\mathcal{L}_w^{\text{local}}$ and retinotopic costs $\mathcal{L}_w^{\text{ret}}$, weighted separately,
216 and then added to the task cost to compute the overall optimization objective:

$$217 \mathcal{L} = L_t + \lambda_w^{\text{local}} \mathcal{L}_w^{\text{local}} + \lambda_w^{\text{ret}} \mathcal{L}_w^{\text{ret}} \quad (9)$$

218 Unless otherwise stated, we use $\lambda_w^{\text{local}} = 0.1$ and $\lambda_w^{\text{ret}} = 0.01$.

219 2.3.2 Polar coordinates

220 Critical to the eccentricity bias theory of cortical organization is that eccentricity is a dominant spatial dimension of early
221 retinotopic maps. While cortical organization is better described by a logarithmic sampling of eccentricity (Schwartz,
222 1980), such a mapping is largely unsuitable for low-resolution images, due to the excessive magnification of the fovea.
223 Thus, following our previous work (Plaut and Behrmann, 2011), we approximate the cortical mapping function as a
224 mapping from Cartesian coordinates (x, y) to polar coordinates (r, θ) , where $r = \sqrt{x^2 + y^2}$ is referred to as eccentricity and
225 $\theta = \arctan(y, x)$ is the polar angle. Polar images are sampled equally in r and θ , with the top of the image corresponding
226 to the upper vertical meridian ($\theta = \pi/2$). For a given square image with side-length resolution s , θ is then incremented in
227 steps of $\Delta_\theta = 2\pi/s$ until the maximum value of $\theta = 3\pi/2 - \Delta_\theta$ is reached. This allows for the top half of polar images
228 to correspond to the left visual field, and the bottom half to the right visual field. Similarly, for a square image of r is
229 sampled linearly from 0 to $\sqrt{2}s$ in increments of $\sqrt{2}$. The sampled coordinates are then converted to Cartesian coordinates
230 for sampling the image. In practice, we use the `warp_polar` function implemented in `scikit-image`. To enable learning of
231 visual representation with polar coordinates, we use wrap-around distances in the angular dimension throughout all layers
232 of the network. That is, in convolutional layers, wrap-around padding is performed in order to complete receptive fields
233 along the vertical meridian. Similarly, in VTC layers, distances are computed with wrap-around boundary conditions along
the polar angle dimension, but not the eccentricity dimension.

234 2.3.3 Viewing biases

235 A critical component of retinotopic constraints is that particular categories or domains of stimuli have distinct eccentricity
236 biases in natural viewing experience (Levy et al., 2001; Hasson et al., 2002). In natural viewing, scenes typically take
237 up the full visual field, and thus place strong demands on peripheral vision. Computational work has demonstrated the
238 preferential informativeness of the periphery, relative to the fovea, in scene recognition (Wang and Cottrell, 2017). In
239 contrast, faces and objects are frequently viewed at smaller sizes and fixated. We model these viewing biases across domains
240 in a simple, yet sufficient way to explore how they might constrain the topographic layout of VTC. Specifically, we use
241 truncated Gaussian distributions of relative image size, where the relative size of each domain d is specified by a mean μ_d
242 and standard deviation σ_d , along with bounds $[s_{d,\min}, s_{d,\max}]$. For each training image i , given the domain d_i , a candidate
243 size s_i is drawn from the Gaussian distribution $N(\mu_d, \sigma_d)$; if the sample drawn is outside the bounds $[s_{d,\min}, s_{d,\max}]$, further
244 samples are drawn until $s_i \in [s_{d,\min}, s_{d,\max}]$. Specific viewing biases used in the main model are given in Figure 1B, with an
example image at two characteristic sizes.

246 2.3.4 Stimuli

247 Three domains of stimuli are used: objects (ImageNet; Deng et al., 2009; Russakovsky et al., 2015), faces (VGGFace2; Cao
248 et al., 2018), and scenes (Places365; Zhou et al., 2018), as in Blauch et al. (2022). In contrast to this previous work, we
249 cropped objects within ImageNet images in order to more precisely control viewing biases as can be done for the cropped
250 face images. To do so, we used ImageNet images for which bounding boxes were available, and selectively cropped the main
251 category of interest. Each dataset was generated to contain 100 categories, with 50,000 training images (500 per category)
252 and 10,000 validation images (100 per category).

253 2.4 Analysis techniques

254 2.4.1 Selectivity

255 We follow the approach from functional neuroimaging (Kanwisher, 2010) in understanding topographic visual organization
256 through selectivity analyses, contrasting particular sets of stimuli. The functional relevance of selectivity in DNNs has
257 been shown by us and others in recent work (Blauch et al., 2022; Prince et al., 2024). Here, we adopt the use of an effect

size metric, namely Cohen's d , applied to a vector of responses for a target domain \mathbf{x}_1 over n_1 examples, and a vector of responses to off-target domains \mathbf{x}_2 over n_2 examples:

$$d = \frac{\mu_1 - \mu_2}{\sigma} \quad (10)$$

Where μ_1 and μ_2 are the means of \mathbf{x}_1 and \mathbf{x}_2 , and s is the pooled standard deviation, computed from the individual standard deviations σ_1 and σ_2 :

$$\sigma = \sqrt{\frac{(n_1 - 1)\sigma_1^2 + (n_2 - 1)\sigma_2^2}{n_1 + n_2 + 2}} \quad (11)$$

Selectivity is computed independently for each unit in the network. To summarize selectivity for faces, objects, and scenes simultaneously, we plot Cohen's d for each domain as an RGB channel normalized to a maximum of $d = 1$.

2.4.2 Eccentricity-based lesions

To analyze the functional consequences of selectivity, we analyze performance after lesions to the model, following our previous work (Blauch et al., 2022). Here, since we are interested in eccentricity-biased organization, we perform eccentricity-biased lesions, corresponding to the deletion of a vertical strip of units with a similar eccentricity bias. We use a sliding window of lesions to map out the impact of lesioning parts of the network with different eccentricity biases. Each retinotopically-biased lesion is applied to 25% of the units in the layer, corresponding to 8 columns of units. Eight lesions equally spaced across the eccentricity-biased (i.e. column) dimension of the map, corresponding to column-wise jumps of 4 units, are performed.

2.4.3 Model consistency across random seeds

To examine the consistency of topographic organization across models, we test multiple random seeds of the model with otherwise identical architectural and training hyper-parameters. The seed controls both the random initialization of weights, as well as the random order of training stimuli. In previous work without retinotopic constraints, changing this random seed was sufficient to rearrange the global layout of category selectivity, while leaving other aspects of representation and topographic organization highly similar (Blauch et al., 2022). To demonstrate inter-subject consistency, we perform group analyses over models with identical hyperparameters except for a different random seed controlling the stimulus presentation order and random initialization of network weights.

2.5 Neuroimaging data

We make use of data from two neuroimaging datasets for comparisons with our models. We briefly describe the datasets below, and refer the reader to the original sources for more information.

2.5.1 Natural Scenes Dataset

For analyses of category-selective topographic organization that do not rely on comparisons with retinotopy data, we made use of the Natural Scenes Dataset (Allen et al., 2022) (NSD), which provides excellent functional localizer data at 7T resolution. We focused on the *fLoc* category-selective mapping data, using pre-computed selectivity contrasts. We loaded these contrasts for each of eight NSD subjects, and the domains of *faces*, *places*, and *objects*, using the `nsdaccess` toolbox (10.5281/zenodo.14165749). We projected these maps into *fsaverage* group cortical surface space using the `nsdcode` toolbox (<https://github.com/cvnlab/nsdcode>).

Following our previous work (Blauch et al., 2025b), we constructed a large ventral temporal cortex (VTC) anatomical parcel from the Human Connectome Project Multi-Modal Parcellation (HCPMMP) group-level atlas in *fsaverage* surface space. We constructed the VTC ROI as the union of several HCPMMP parcels: V8, pIT, FFC, VVC, PHA1, PHA2, PHA3, TE2p, TF, PH, VMV1, VMV2, VMV3. This constrained all further analyses of NSD. For each voxel, given a vector of selectivity t -values across subjects with mean μ and standard deviation σ , we constructed group selectivity maps using the one-sample version of Cohen's d compared to a baseline of 0: $d = \frac{\mu}{\sigma}$

2.5.2 Human Connectome Project

To analyze the relationship between early visual cortex eccentricity bias and high-level visual cortex category-selectivity, we make use of the Human Connectome Project (HCP) (Glasser et al., 2013), focusing on the visual working memory task which can serve as a localizer for visual category-selectivity (Glasser et al., 2016), along with the HCP 7T retinotopy dataset.

Visual working memory task. The visual working memory (WM) task consisted of a block design of visual categories – faces, bodies, tools, and places – with a 0-back or 2-back working-memory task, which was used in order to index category selectivity in VTC and neural mechanisms for visual working memory. Our focus is on the former usage. We focus on face and place (i.e. scene) selectivity, using the standard contrasts "FACE-AVG" and "PLACE-AVG", which contrasts the responses of faces and places to those of all other categories. We downloaded pre-computed category selective maps available on Amazon Web Services S3 storage, under the HCP_1200 folder. Specifically, we extracted surface-based contrast maps in the *32k_fs_LR* group CIFTI space of the HCP preprocessing pipeline (Glasser et al., 2016), using standard sulcal-based

307 alignment rather than the multimodal-surface-matching procedure (Robinson et al., 2014). We acquired the maps smoothed
308 with a Gaussian kernel of 2mm FWHM, the smallest amount of smoothing available for download. We used the same VTC
309 ROI as described above. We used the `neuromaps` and `hcp_utils` packages to transform maps from the `32k_fs_LR` space to
310 the `fsaverage` template, for comparison with HCP 7T retinotopy data and visualization in `PyCortex` (Gao et al., 2015).

311 **7T retinotopy.** The HCP 7T retinotopy dataset (Benson et al., 2018) provides retinotopic localizers for 181 of the
312 original HCP subjects at 7T, allowing for precise retinotopic mapping. Stimuli subtended up to 8 deg in the periphery, in
313 contrast to the 4.2 deg subtended by NSD retinotopic stimuli, making it more suitable for analyzing peripherally-responsive
314 cortex. We downloaded from a publicly available repository on OSF <https://osf.io/bw9ec> the pre-processed population
315 receptive field (pRF) mapping results, described further in (Benson et al., 2018). Specifically, we focused on the estimated
316 pRF eccentricity maps in `fsaverage` group space.

317 **Computing geodesic distance between V4 voxels and category-selective regions.** To assess the retinotopic
318 wiring constraints on category-selective regions in VTC, we computed geodesic (surface) distances between these category-
319 selective regions and each voxel in V4. We computed category-selective regions by constraining a functional mask of
320 selectivity with an anatomical mask spanning VTC. We used the same VTC parcel as described in 2.5.1, and adopted a
321 relatively high threshold of functional selectivity ($t > 6$) to avoid the presence of spurious voxels which would artificially
322 deflate the minimum geodesic distance values to V4 vertices. For a V4 parcel, we noted that the HCPMP1 atlas (Glasser
323 et al., 2016) and Wang atlas (Wang et al., 2015) gave very different anatomical definitions of the early visual regions
324 V2-V4. Whereas the HCPMP1 atlas provides ring-like regions for V2, V3, and V4, spanning both ventral and dorsal
325 cortex, the Wang atlas divides V2 and V3 into separate dorsal and ventral regions, and defines a single ventral "hV4"
326 region. This hV4 regions spans only a very small degree of the eccentricity gradient visible in the group-level map. We
327 thus opted to construct a ventral V4 parcel by constraining the HCPMP V4 parcel within the ventral cortex. To do
328 so, we intersected the V4 parcel with a ventral ROI composed of the "midventral" and "ventral" ROIs of the NSD (Allen
329 et al., 2022) *nsdgeneral* atlas. We found that the mean pRF eccentricity within this region yielded the predicted smooth
330 lateral-to-medial progression from foveal to peripheral eccentricity bias, as shown in Figure 3A. To compute geodesics, we
331 used an implementation of the efficient heat-based geodesic method (Crane et al., 2013) implemented within the PyCortex
332 software package (Gao et al., 2015). Given an ROI, this method efficiently computes the minimum geodesic distance to
333 every vertex on the cortical surface, separately for each hemisphere. Thus, we computed the minimum geodesic distance
334 between face- and place-selective regions within VTC to each V4 vertex.

3 Results

3.1 Globally consistent topographic organization

337 Following training, we test the selectivity of units in the VTC model layers to faces, objects, and scenes, using held-out
338 images not seen during training, presented at random relative sizes that are uniformly drawn from 0.3 to 1 in the distribution
339 of size. We compare the selectivity of 8 instances of the model with the 8 subjects of the Natural Scenes Dataset (NSD).
340 Human VTC and model VTC category selectivity are shown in Figure 2. In the first column of each row, we plot the
341 group effect size map using Cohen's d , demonstrating strong topographic selectivity in both human and model VTC. In
342 the remaining columns, we show selectivity in each individual human or model. As can be seen in Figure 2A, while there
343 are some differences across individual human VTCs, the broad theme of vertical strips of selectivity running along the
344 posterior-anterior dimension, with domains organized across the medial-lateral dimension, is apparent in every subject. In
345 Figure 2B, we show this same pattern in each model instance, with a striking degree of consistency.

3.2 Abolishing retinotopic constraints abolishes the consistency, but not presence, of topographic organization

346 What happens if we remove the retinotopic constraint? To answer this question, we set the retinotopic wiring penalty
347 $\mathcal{L}_w^{\text{ret}} = 0$, increasing $\mathcal{L}_w^{\text{local}}$ from 0.1 to 0.5 to account for the lack of retinotopic wiring cost. We trained 8 models like
348 this, plotting the results in Figure 2C; we additionally plot models with $\mathcal{L}_w^{\text{local}} = 0.1$ in Figure S1. On the left, we show a
349 group Cohen's d selectivity map, computed identically to the one shown in Figure 2B. On the right, we show selectivity
350 maps for individual models. Here, we see strong topographic clustering, but an inconsistent global layout. This leads to
351 weak group-level selectivity, in contrast to the main model. Thus, the retinotopic constraint here is necessary to sculpt
352 the global organization of topography, but not to produce topography in the first place. Notably, the organization here is
353 not completely random (in contrast to the model of Blauch et al., 2022), with a preference for vertical strips of selectivity
354 similar to that seen in the retinotopic model. This is due to the use of wrap-around distances along the angular dimension,
355 which makes communication between the top and bottom of VTC more efficient than communication between the left and
356 right. This thus yields a bias towards organizing specialized representations vertically, since this allows for the use of shorter
357 connections. Nevertheless, these results demonstrate that the consistent global layout is determined specifically by the
358 retinotopic connectivity constraint.

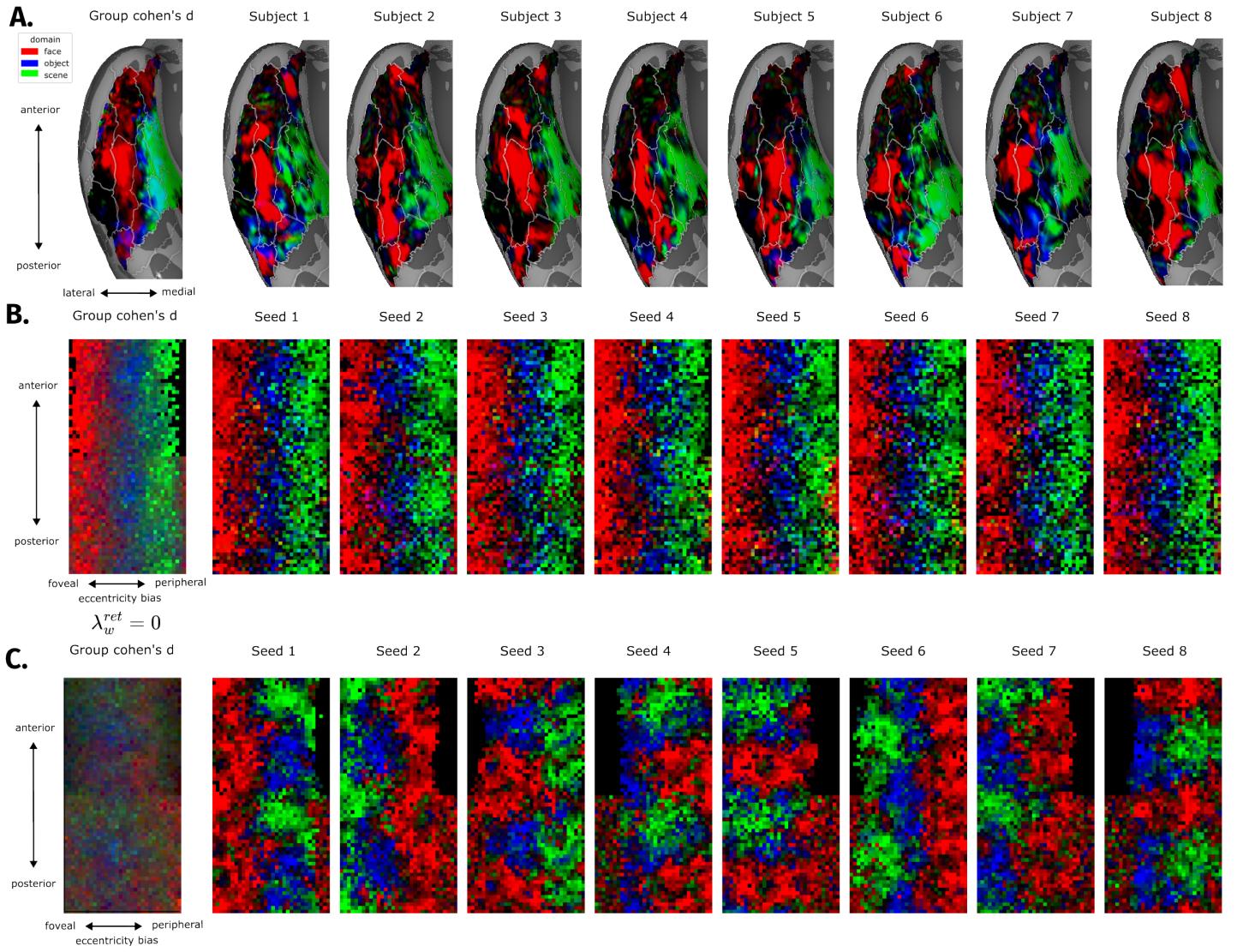


Figure 2: Consistent topographic organization from retinotopic connectivity constraints. **A.** Individual model domain-level selectivity plots. **B.** Group-level selectivity probability maps. For each domain, we plot the probability of selectivity at a given VTC location, computed over the 8 model runs shown in panel A., at a threshold of Cohen's $d > 0.3$. **C.** Group-level selectivity probability maps in models trained without a wiring cost ($\mathcal{L}_w^{\text{ret}} = 0$, increasing $\mathcal{L}_w^{\text{local}}$ to 0.5 (from 0.1) to account for the lack of retinotopic wiring cost.

3.3 Eccentricity-biased organization in human and model high-level visual cortex

Next, we examine the effects of eccentricity bias in greater detail. In Figure 3A, we plot category selectivity in VTC alongside eccentricity preference in V4, in both the model and human VTC. (Note that, due to the use of HCP, we do not plot object selectivity.) We see that the category selectivity in VTC is neatly organized along the eccentricity gradient in V4, in both the model and human VTC. We quantify the proximity between V4 voxels and pVTC category selective regions in individual subjects/models and aggregate results at the group level (see 2.5.2 for details on geodesic calculations in human fMRI data). As shown in Figure 3B, in both cases, we find that the scene-selective area is significantly closer to more peripherally preferring V4 voxels/units, whereas the face-selective area is significantly closer to foveally-prefering V4 voxels/units. In Figure 3C, we plot the eccentricity bias of pVTC. For each pVTC unit, we compute a weighted sum of feedforward V4 weights, weighting the input eccentricity according to the learned weight magnitude (note that feedforward weights are strictly positive), and summing across V4 channels; we then divide by the number of input units to yield the weighted eccentricity bias preference map. A smooth gradient is seen, allowing for the minimization of the retinotopic wiring cost $\mathcal{L}_w^{\text{ret}}$. Note that this eccentricity bias is not hard-coded into the connectivity between V4 and pVTC, but rather emerges to minimize the overall objective function.

We next demonstrate the consequences of this eccentricity bias for domain-level representations. In Figure 3D, we progressively lesion vertical strips corresponding to 25% of pVTC, and measure the accuracy of within-domain recognition for faces, objects, and scenes. We plot the result as a function of the difference from unlesioned performance, ignoring absolute differences between domains. The results reveal a strong eccentricity bias in the locus of damage to each domain:

foveally-biased lesions primarily affect face recognition, whereas peripherally-biased lesions primarily affect scene recognition; intermediate lesions primarily affect object recognition. These results allow us to conclude that the observed functional organization is functionally significant.

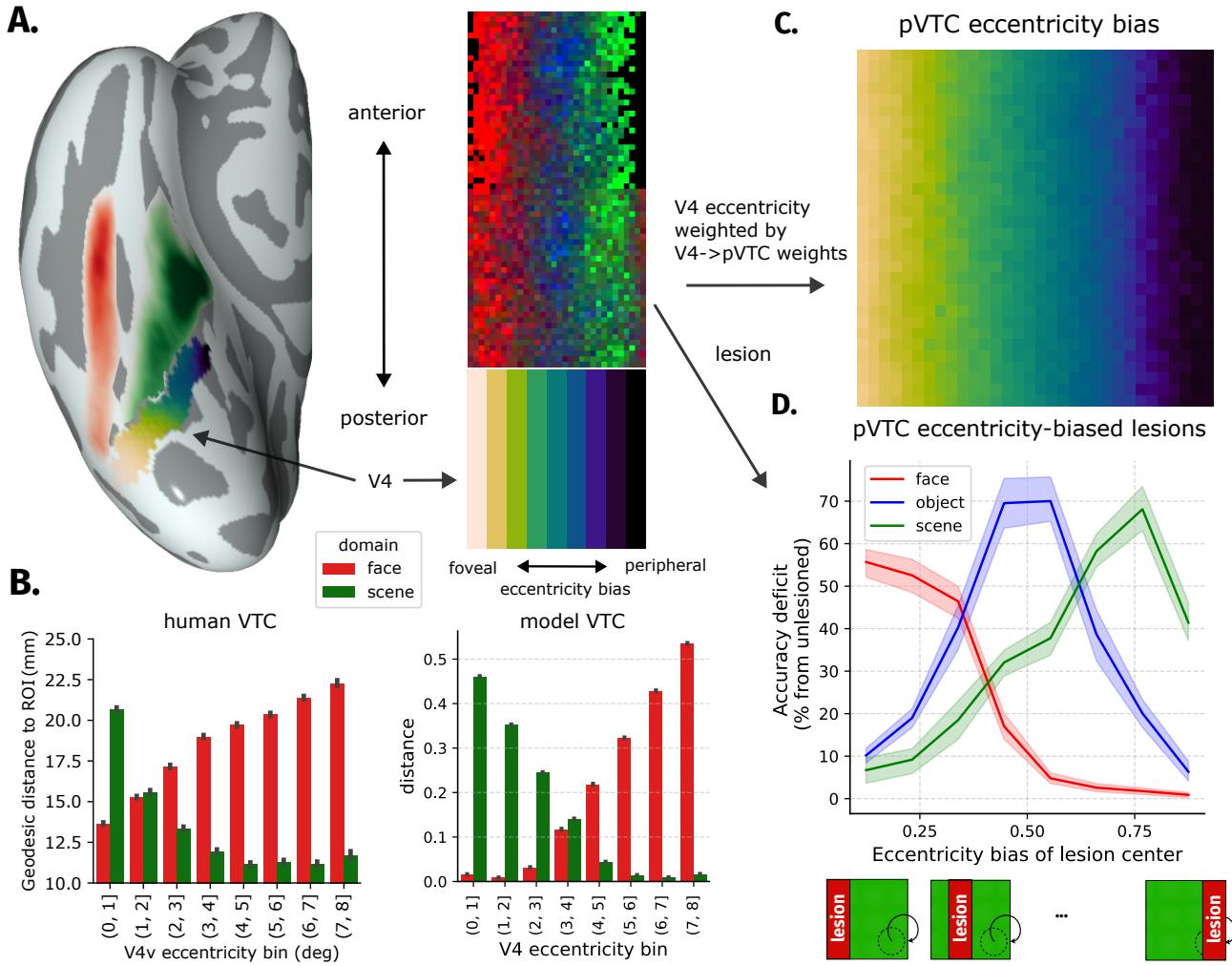


Figure 3: Retinotopic biases on category-selective organization. **A.** Eccentricity bias and category-selective organization in the human ventral temporal cortex (left) and computational model (right). In both cases, we show the eccentricity representation of the "V4" region, and category selectivity in the "VTC" region. In the human brain, we plot the average eccentricity over all subjects in the HCP 7T retinotopy dataset, masked within the ventral portion of the HCPMMP V4 region. In the model, we plot eccentricity as the column index of the polar feature map of the final convolutional layer, which we call "V4". Both human and model VTC show a layout of face and scene selectivity that corresponds to the V4 eccentricity gradient, with place selectivity located more medially, proximal to peripheral V4, and face representations located more laterally, proximal to foveal V4. **B.** Quantification of proximity between category-selective regions and V4 eccentricity in individual subjects. Left: Human VTC. Right: Model VTC. Y-axis: Geodesic distance to ROI (mm) or distance. X-axis: V4 eccentricity bin (deg).

3.4 Eccentricity bias persists across retinotopic variation

One question that arises from these results is whether the retinotopically-biased topographic organization is indeed a learned, efficient organization optimized for the viewing biases and retinotopic connectivity constraints, versus merely an explicit retinotopic code inherited from the V4 region (owing to that region's convolutional nature). The results thus far have collapsed across responses to images presented at a range of sizes, generally suggesting that the organization generalizes across size and is, therefore, a learned organization. However, the extent of generalization is yet unclear. To assess this, we analyze the topographic organization explicitly at a range of relative viewing sizes, rather than collapsing across them. The results, shown in Figure 4A, demonstrate a striking consistency of the topographic organization across a large range of viewing sizes, indicating a relative invariance to changes in the retinotopic inputs. We quantify this by localizing the selective regions at the intermediate size of 0.6, and analyzing domain-level responses across the full range of relative sizes. In both pVTC and aVTC, we see that face, object, and scene selectivity is nearly invariant to relative size, remaining strong

393 across the full range of sizes. These results are in line with the idea that the retinotopically-biased topographic organization
394 of VTC is an efficient organization of category-representations based on experienced viewing biases, rather than on an
395 explicit retinotopic code.

396 3.5 Topography becomes more explicitly retinotopic with less variation in viewing 397 conditions

398 How does the learned efficient organization optimized for viewing biases and retinotopic connectivity constraints arise?
399 One hypothesis is that it arises due to sufficient variation in viewing conditions, allowing for the extraction of an invariant
400 representation whose location is optimized to the mean of viewing conditions. To test this hypothesis, we trained another
401 model with less within-domain variation in viewing size biases. As shown in Figure 4B, we kept the means of relative size
402 distributions the same, but reduced the variance. After training, this model showed a similar domain-level topography,
403 however, it was less invariant to relative size. As inputs were presented at smaller sizes, the intermediate and far periphery
404 reduced their consistency in selectivity for scenes and objects, with more intermixing (Figure 4B). We quantified this by
405 analyzing the mean responses in category selective VTC regions computed at the intermediate size of 0.6, as in the previous
406 set of analyses. This confirms what is apparent in the selectivity maps, that at smaller relative sizes, the scene selective area
407 becomes less selective for scenes and the object selective area becomes less selective for objects; this is driven by a reduction
408 in responses to scenes in the scene-selective area, and a large increase in responses to scenes in the object-selective area. In
409 contrast, selectivity for faces is strongly invariant to relative size. These results demonstrate that the extent to which the
410 high-level topography retains a retinotopic *representational* bias, vs. purely a retinotopic *location* bias, is dependent on the
411 learned viewing conditions, and how sufficient they are for developing an invariant representation. Further viewing bias
412 experiments are presented in the Appendix, Figure S2.

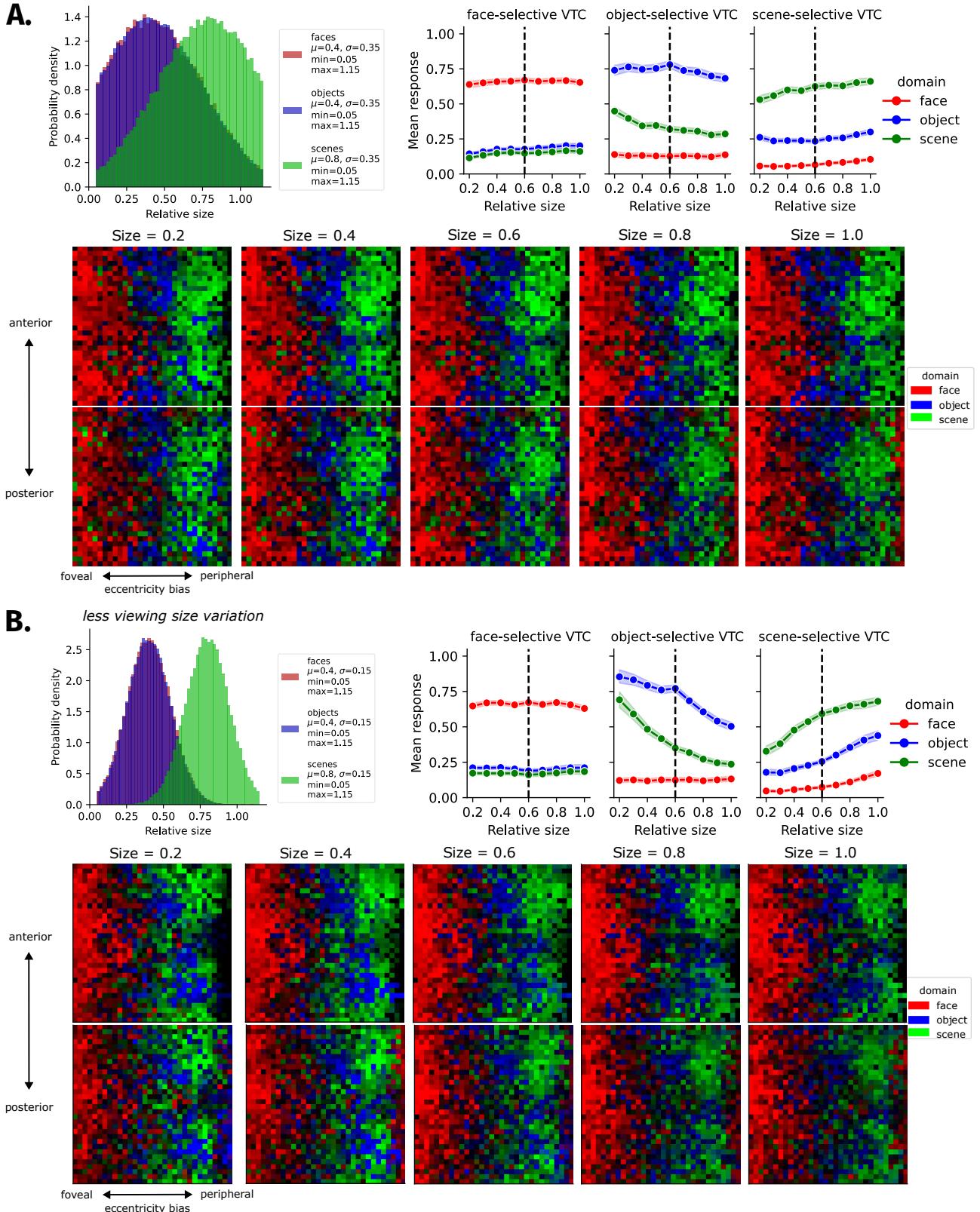


Figure 4: Eccentricity-biased topographic organization generalizes across large changes in retinotopic input size, given sufficient viewing bias variation. **A.** Domain-level topographic organization in the main model VTC across variation in the relative viewing size. Top-left: viewing size distributions during training. Top-right: mean responses to each domain in each domain-selective area (computed at the intermediate size of 0.6) across the full range of sizes. Bottom: VTC domain-selectivity across 5 representative relative sizes. **B.** Domain-level topographic organization in a model trained with less variation in domain-specific viewing size biases. The layout is the same as in panel **A**.

4 Discussion

The wiring minimization principle has been successful in explaining aspects of both local smooth functional organization (Mitchison, 1991; Koulakov and Chklovskii, 2001; Chklovskii and Koulakov, 2004) as well as the organization of entire connectomes, including *C. elegans* (Chen et al., 2006), and, with less precision, human connectomes measured with functional and structural imaging (Bassett et al., 2010; Raj and Chen, 2011). Recent topographic models of high-level vision have shown how wiring constraints from local to high-level visual processing layers can give rise to many key characteristics of VTC topographic organization (Blauch et al., 2022; Keller and Welling, 2021). Similar results emerge in broader *local-constraint* accounts that do not model local connectivity directly, but arise through other smoothness constraints on high-level visual representations (Doshi and Konkle, 2023; Margalit et al., 2024). These local-constraint models have explained the *presence* and form of certain topographic visual features of the high-level visual cortex, suggesting that long-range connectivity is not necessary for such organization to emerge, and that an efficient mapping of features that describe the visual world may give rise to a large degree of the organization of VTC. However, purely local models provide no explanation for the *consistency* of the global layout of cortical organization across individuals, suggesting that additional, long-range connectivity constraints are needed.

We argue that a specific *connectivity-based* view on local constraints allows for local constraints to naturally be subsumed within a more generic optimization-based view on VTC cortical organization, where both local and long-range connectivity constraints play an important role. To do so, we consider such optimization to occur over both phylogenetic and ontogenetic timescales. Over evolution, nature has selected for brains that do not connect indiscriminately, but exhibit a bias towards local connectivity, mitigating the extreme spatial costs of full random connectivity. However, to support efficient network integration, specific long-range fiber bundles are required (Bullmore and Sporns, 2012) and have been selected for, as indicated by comparative studies and heritability (Ardesch et al., 2019; Miranda-Dominguez et al., 2018). Last, evolution has selected for powerful learning algorithms that allow for ontogenetic development to further optimize brain representations for behavioral and energetic demands (Richards et al., 2019). Such learning algorithms allow for representations to be sculpted into the brain according to the experience of the organism (Dehaene and Cohen, 2007; Arcaro et al., 2019). Critically, such learning is constrained by the particular long-range pathways that exist, the particular developmental processes that constrain local branching patterns, and the particular mechanisms of plasticity that allow neurons to modify the weights from and onto other neurons.

Our work provides a computational demonstration of how the innate retinotopic organization of early visual representations may serve as a long-range constraint, or scaffold, on the organization of high-level vision, in line with earlier theories (Hasson et al., 2002; Arcaro et al., 2019). We demonstrated that this retinotopic bias was not necessary for the emergence of high-level topographic organization (which arose due to local connectivity constraints) but was critical for the emergence of a *spatially consistent* layout. This spatially consistent layout exhibited a remarkable retinotopic bias, with topographic clusters responsive to faces and scenes laying in areas proximal to foveal and peripheral inputs, respectively, as in the human brain. We found that this did not necessitate strong retinotopic coding; under sufficiently broad distributions of experienced viewing biases, the emergent category-selective areas yielded strong invariance to the particular retinotopic properties of stimuli, despite their locations being optimized for the learned retinotopic biases. However, under narrower distributions of learned viewing biases, a stronger degree of retinotopic sensitivity was seen. This thus provides a general principle for the retinotopic influences on high-level vision, in which retinotopic biases shape the organization of high-level vision but can be abstracted away under sufficient conditions for invariant representation learning.

One limitation of our model is its relatively simplified implementation of retinotopic input constraints, using a polar-coordinate transformation of an input image. This implementation does not capture the precise details of foveated retinal sampling (Watson, 2014) and the exact magnification of the fovea in human visual cortical maps (Daniel and Whitteridge, 1961), in favor of a simplified model that captures a key critical detail – mediolateral organization of eccentricity. A more realistic approach would use a logarithmic sampling of the input corresponding to both spatially-variant retinal sampling and cortical magnification (Schwartz, 1980). However, accurately characterizing foveated sensing requires the use of a curved manifold rather than a rectangular image, and an associated specialized architecture (Blauch et al., 2025a). Additionally, properly modeling foveal sampling necessitates the use of multiple fixations to aggregate information over time, and ideally, an intelligent mechanism for choosing where to fixate. Thus, in this work we model the minimal important detail of early retinotopic organization – mediolateral eccentricity organization – and delegate the modeling of foveated sensing, and the study of its impact on high-level visual representations, to future work.

The other main limitation of our model is its limited scope on the possible long-range influences on high-level vision. For example, the impact of language on the emergence of lateralized responses to words is firmly established, for example, with demonstrations that individual variability in VTC word laterality is directly related to language laterality (Gerrits et al., 2019) and word responses in traditional language-selective areas (Blauch et al., 2025b). Moreover, many others have argued that long-range connectivity constraints shape the organization of high-level visual cortex beyond the lateralization of word-selective responses (Saygin et al., 2012, 2016; Powell et al., 2018; Ratan Murty et al., 2020). The large-scale impact of language processing on broad cortical networks appears to shape the hemispheric organization of domains beyond word recognition, for example, through competitive influences social processing (Rajimehr et al., 2022), that may impose long-range coupling constraints on high-level visual cortex for domains such as faces (Blauch et al., 2025b). Modeling all of the constraints on high-level visual cortex within task-optimized deep learning is a challenging endeavor, and it is our view that systematically incorporating additional constraints is a good route for tempering this complexity and precisely understanding each constraint in turn. Thus, the current work should not be taken to claim that the retinotopic organization of EVC is the only external constraint on VTC organization. Rather, our argument is that it is a very powerful constraint shaping VTC, and we provide computational simulations as a proof of concept.

We find that simple biases in viewing, combined with a constraint to minimize connectivity between retinotopically organized early-mid-level visual areas with high-level visual areas, results in a consistently biased global layout that places high-level scene representations nearer to peripheral visual input, and face representations nearer to foveal visual input. Owing to a requirement to represent these high-level visual categories across a range of viewing sizes, this organization is not *retinotopic per se*, but *retinotopically biased*, with stable category preferences across a range of retinotopic variation, as has been found in the human brain (Hasson et al., 2002; Silson et al., 2021; Park et al., 2024). Critically, topographic organization in the model emerges regardless of whether the retinotopic connectivity constraint is included; the retinotopic constraint ensures that the global organization is consistent in order to minimize the connectivity with retinotopic input areas.

If other long-range connectivity constraints do shape the organization of high-level visual cortex, we propose that the retinotopic constraint may have served as an important evolutionary scaffold on which these other long-range connections have evolved. While high-level vision in higher primates interfaces with complex cognitive behaviors such as social perception (Isik et al., 2017), tool use (Mahon et al., 2007), and, specifically in humans, reading (Dehaene and Cohen, 2007), spatial vision pre-dates all of these behaviors phylogenetically, and retinotopic maps exist across the animal kingdom (Cisek, 2019). To the extent that long-range connectivity has been genetically modified to support particular human cognitive tasks, it had to interface with an already retinotopically-biased visual system – this retinotopic bias would have provided an inherently systematic organization on which long-range connectivity could operate, leading to natural selection of particularly effective fiber tracts. Thus, while there is some evidence that category-selective organization persists in the absence of visual input (Mahon et al., 2009; van den Hurk et al., 2017; Ratan Murty et al., 2020), such evidence should not be taken to imply that retinotopic influences and visual experience do not shape VTC organization. Indeed, in sighted monkeys deprived of experience with faces, topographic face selectivity does not emerge (Arcaro et al., 2017, 2019; Arcaro and Livingstone, 2021). We leave for future work a systematic exploration of long-range connectivity constraints. With the emergence of agentic modeling frameworks and ever more powerful and available computational resources, systematic investigation of increasingly plausible interactions between vision and other cognitive functions will be possible. Such agentic modeling will allow for viewing biases to emerge entirely naturally due to the agent's actions in the world, including active foveated eye movements, which will depend on higher cognitive demands. Such advances will provide important refinements to the wiring minimization account of cortical organization in high-level vision and beyond.

Acknowledgments

We thank Michael Arcaro and Talia Konkle for useful discussions. This research was supported by a grant from the National Science Foundation to M.B. and D.C.P. (BCS 2123069), and by a grant from the National Institute of Health (R01EY027018) to M.B. and D.C.P. M.B. also acknowledges support from grant R01EY026701 and a P30 CORE award EY08098 from the NEI, as well as unrestricted supporting funds from The Research to Prevent Blindness Inc, NY and the Eye & Ear Foundation of Pittsburgh.

Author contributions

Nicholas M. Blauch: Conceptualization, Methodology, Software, Formal analysis, Investigation, Funding acquisition, Writing - Original Draft, Writing - Review & Editing, Visualization. Marlene Behrmann: Conceptualization, Supervision, Funding acquisition, Writing - Review & Editing. David C. Plaut: Conceptualization, Supervision, Funding acquisition, Writing - Review & Editing.

Data availability

NSD and HCP datasets are publicly available as described in Methods. All code necessary to reproduce the results of the paper will be made available upon publication.

References

- Allen, E. J., St-Yves, G., Wu, Y., Breedlove, J. L., Prince, J. S., Dowdle, L. T., Nau, M., Caron, B., Pestilli, F., Charest, I., et al. (2022). A massive 7T fMRI dataset to bridge cognitive neuroscience and artificial intelligence. *Nature neuroscience*, 25(1):116–126.
- Arcaro, M. J. and Livingstone, M. S. (2017). A hierarchical , retinotopic proto- organization of the primate visual system at birth. *eLife*, pages 1–24.
- Arcaro, M. J. and Livingstone, M. S. (2021). On the relationship between maps and domains in inferotemporal cortex. *Nature Reviews Neuroscience*.
- Arcaro, M. J., Schade, P. F., and Livingstone, M. S. (2019). Universal Mechanisms and the Development of the Face Network: What You See Is What You Get. *Annual Review of Vision Science*, 5(1):341–372.
- Arcaro, M. J., Schade, P. F., Vincent, J. L., Ponce, C. R., and Livingstone, M. S. (2017). Seeing faces is necessary for face-domain formation. *Nature Neuroscience*, (September).
- Ardesch, D. J., Scholtens, L. H., and Van Den Heuvel, M. P. (2019). The human connectome from an evolutionary perspective. In *Progress in Brain Research*, volume 250, pages 129–151. Elsevier.
- Ba, J. L., Kiros, J. R., and Hinton, G. E. (2016). Layer Normalization. *arXiv:1607.06450 [cs, stat]*.
- Bao, P., She, L., Mcgill, M., and Tsao, D. Y. (2020). A map of object space in primate inferotemporal cortex. *Nature*, (January 2019).
- Bassett, D. S., Greenfield, D. L., Meyer-Lindenberg, A., Weinberger, D. R., Moore, S. W., and Bullmore, E. T. (2010). Efficient Physical Embedding of Topologically Complex Information Processing Networks in Brains and Computer Circuits. *PLoS Computational Biology*, 6(4):e1000748.
- Behrmann, M. and Plaut, D. C. (2015). A vision of graded hemispheric specialization. *Annals of the New York Academy of Sciences*, 1359(1):30–46.
- Behrmann, M. and Plaut, D. C. (2020). Hemispheric Organization for Visual Object Recognition: A Theoretical Account and Empirical Evidence. *Perception*, 49(4):373–404.
- Benson, N. C., Jamison, K. W., Arcaro, M. J., Vu, A. T., Glasser, M. F., Coalson, T. S., Van Essen, D. C., Yacoub, E., Ugurbil, K., Winawer, J., and Kay, K. (2018). The Human Connectome Project 7 Tesla retinotopy dataset: Description and population receptive field analysis. *Journal of Vision*, 18(13):23.
- Blauch, N. M., Alvarez, G. A., and Konkle, T. (2025a). Foveated sensing with KNN-convolutional neural networks based on isotropic cortical magnification.
- Blauch, N. M., Behrmann, M., and Plaut, D. C. (2022). A connectivity-constrained computational account of topographic organization in primate high-level visual cortex. *Proceedings of the National Academy of Sciences of the United States of America*, 119(3).
- Blauch, N. M., Plaut, D. C., Vin, R., and Behrmann, M. (2025b). Individual variation in the functional lateralization of human ventral temporal cortex: Local competition and long-range coupling. *Imaging Neuroscience*, 3:imag_a_00488.
- Brodmann, K. (1909). *Vergleichende Lokalisationslehre Der Grosshirnrinde in Ihren Prinzipien Dargestellt Auf Grund Des Zellenbaues*. Barth.
- Bullmore, E. and Sporns, O. (2012). The economy of brain network organization. *Nature Reviews Neuroscience*, 13(5):336–349.
- Cao, Q., Shen, L., Xie, W., Parkhi, O. M., and Zisserman, A. (2018). VGGFace2: A dataset for recognising faces across pose and age. In *International Conference on Automatic Face and Gesture Recognition*.
- Chao, L. L., Haxby, J. V., and Martin, A. (1999). Attribute-based neural substrates in temporal cortex for perceiving and knowing about objects. *Nature Neuroscience*, 2(10):913–919.
- Chen, B. L., Hall, D. H., and Chklovskii, D. B. (2006). Wiring optimization can relate neuronal structure and function. *Proceedings of the National Academy of Sciences*, 103(12):4723–4728.
- Chklovskii, D. B. and Koulakov, A. A. (2004). MAPS IN THE BRAIN: What Can We Learn from Them? *Annual Review of Neuroscience*, 27(1):369–392.
- Cipollini, B. and Cottrell, G. (2013). Uniquely human developmental timing may drive cerebral lateralization and interhemispheric collaboration. In *Proceedings of the Cognitive Science Society*, 35(35), pages 334–339.

- 566 Cisek, P. (2019). Resynthesizing behavior through phylogenetic refinement. *Attention, Perception, & Psychophysics*,
567 81(7):2265–2287.
- 568 Crane, K., Weischedel, C., and Wardetzky, M. (2013). Geodesics in Heat. *ACM Transactions on Graphics*, 32(5):1–11.
- 569 Crick, F., Asanuma, C., et al. (1986). Certain aspects of the anatomy and physiology of the cerebral cortex. *Parallel
570 distributed processing*, 2:333–371.
- 571 Daniel, P. M. and Whitteridge, D. (1961). The representation of the visual field on the cerebral cortex in monkeys. *The
572 Journal of Physiology*, 159(2):203–221.
- 573 Deb, M., Deb, M., and Murty, N. A. R. (2025). TopoNets: High Performing Vision and Language Models with Brain-Like
574 Topography.
- 575 Dehaene, S. and Cohen, L. (2007). Cultural Recycling of Cortical Maps. *Neuron*, 56(2):384–398.
- 576 Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L. (2009). ImageNet: A Large-Scale Hierarchical Image
577 Database. In *CVPR*, pages 248–255. IEEE.
- 578 Doshi, F. R. and Konkle, T. (2023). Cortical topographic motifs emerge in a self-organized map of object space. *Science
579 Advances*, 9(25):eade8187.
- 580 Epstein, R. and Kanwisher, N. (1998). A cortical representation of the local visual environment. *Nature*, 392(6676):598–601.
- 581 Finzi, D., Gomez, J., Nordt, M., Rezai, A. A., Poltoratski, S., and Grill-Spector, K. (2021). Differential spatial computations
582 in ventral and lateral face-selective regions are scaffolded by structural connections. *Nature Communications*, 12(1):2278.
- 583 Fodor, J. (1983). *The Modularity of Mind: An Essay on Faculty Psychology*. The MIT Press, Cambridge, MA.
- 584 Gao, J. S., Huth, A. G., Lescroart, M. D., and Gallant, J. L. (2015). Pycortex: An interactive surface visualizer for fMRI.
585 *Frontiers in Neuroinformatics*, 9(September):1–12.
- 586 Gauthaman, R. M., Ménard, B., and Bonner, M. F. (2024). Universal scale-free representations in human visual cortex.
- 587 Gauthaman, R. M., Ménard, B., and Bonner, M. F. (2025). Spatial-scale invariant properties of primary visual cortex in
588 humans and mice.
- 589 Gerrits, R., Van der Haegen, L., Brysbaert, M., and Vingerhoets, G. (2019). Laterality for recognizing written words and
590 faces in the fusiform gyrus covaries with language dominance. *Cortex*, 117:196–204.
- 591 Glasser, M. F., Coalson, T. S., Robinson, E. C., Hacker, C. D., Harwell, J., Yacoub, E., Ugurbil, K., Andersson, J.,
592 Beckmann, C. F., Jenkinson, M., Smith, S. M., and Van Essen, D. C. (2016). A multi-modal parcellation of human
593 cerebral cortex. *Nature*, 536:171–178.
- 594 Glasser, M. F., Sotiroopoulos, S. N., Wilson, J. A., Coalson, T. S., Fischl, B., Andersson, J. L., Xu, J., Jbabdi, S., Webster,
595 M., Polimeni, J. R., Van Essen, D. C., and Jenkinson, M. (2013). The minimal preprocessing pipelines for the Human
596 Connectome Project. *NeuroImage*, 80:105–124.
- 597 Grill-Spector, K. and Weiner, K. S. (2014). The functional architecture of the ventral temporal cortex and its role in
598 categorization. *Nature reviews Neuroscience*, 15(8):536–548.
- 599 Groen, I. I., Dekker, T. M., Knapen, T., and Silson, E. H. (2022). Visuospatial coding as ubiquitous scaffolding for human
600 cognition. *Trends in Cognitive Sciences*, 26(1):81–96.
- 601 Hasson, U., Levy, I., Behrmann, M., Hendler, T., and Malach, R. (2002). Eccentricity bias as an organizing principle for
602 human high-order object areas. *Neuron*, 34(3):479–490.
- 603 Haxby, J. V., Gobbini, M. I., Furey, M. L., Ishai, A., Schouten, J. L., and Pietrini, P. (2001). Distributed and overlapping
604 representations of faces and objects in ventral temporal cortex. *Science*, 293(September):2425–2430.
- 605 Haxby, J. V., Guntupalli, J. S., Connolly, A. C., Halchenko, Y. O., Conroy, B. R., Gobbini, M. I., Hanke, M., and Ramadge,
606 P. J. (2011). A common, high-dimensional model of the representational space in human ventral temporal cortex. *Neuron*,
607 72(2):404–416.
- 608 Haxby, J. V., Guntupalli, J. S., Nastase, S. A., and Feilong, M. (2020). Hyperalignment: Modeling shared information
609 encoded in idiosyncratic cortical topographies. *eLife*, 9:e56601.
- 610 Hubel, D. H. and Wiesel, T. N. (1969). Anatomical Demonstration of Columns in the Monkey Striate Cortex. *Nature*,
611 221(5182):747–750.
- 612 Huberman, A. D., Feller, M. B., and Chapman, B. (2008). Mechanisms Underlying Development of Visual Maps and
613 Receptive Fields. *Annual Review of Neuroscience*, 31(1):479–509.

- 614 Huth, A. G., de Heer, W. A., Griffiths, T. L., Theunissen, F. E., and Gallant, J. L. (2016). Natural speech reveals the
615 semantic maps that tile human cerebral cortex. *Nature*, 532(7600):453–458.
- 616 Ishai, A., Ungerleider, L. G., Martin, A., Schouten, J. L., and Haxby, J. V. (1999). Distributed representation of objects in
617 the human ventral visual pathway. *Proceedings of the National Academy of Sciences of the United States of America*,
618 96(16):9379–9384.
- 619 Isik, L., Koldewyn, K., Beeler, D., and Kanwisher, N. (2017). Perceiving social interactions in the posterior superior
620 temporal sulcus. *Proceedings of the National Academy of Sciences*, 114(43).
- 621 Jacobs, R. A. and Jordan, M. I. (1992). Computational consequences of a bias toward short connections. *Journal of
622 Cognitive Neuroscience*, 4(4):323–336.
- 623 Kanwisher, N. (2010). Functional specificity in the human brain: A window into the functional architecture of the mind.
624 *Proceedings of the National Academy of Sciences of the United States of America*, 107(25):11163–11170.
- 625 Kanwisher, N., McDermott, J., and Chun, M. M. (1997). The Fusiform Face Area: A Module in Human Extrastriate Cortex
626 Specialized for Face Perception. *Journal of Neuroscience*, 17(11):4302–4311.
- 627 Keller, T. A. and Welling, M. (2021). Topographic VAEs learn Equivariant Capsules. *arXiv*.
- 628 Konkle, T. and Caramazza, A. (2013). Tripartite Organization of the Ventral Stream by Animacy and Object Size. *Journal
629 of Neuroscience*, 33(25):10235–10242.
- 630 Koulakov, A. A. and Chklovskii, D. B. (2001). Orientation Preference Patterns in Mammalian Visual Cortex: A Wire
631 Length Minimization Approach. *Neuron*, 29:519–527.
- 632 Laszlo, S. and Plaut, D. C. (2012). A neurally plausible Parallel Distributed Processing model of Event-Related Potential
633 word reading data. *Brain and Language*, 120(3):271–281.
- 634 Levy, I., Hasson, U., Avidan, G., Handler, T., and Malach, R. (2001). Center-periphery organization of human object areas.
635 *Nature Neuroscience*, 4(5):533–539.
- 636 Long, B., Yu, C.-P., and Konkle, T. (2018). Mid-level visual features underlie the high-level categorical organization of the
637 ventral stream. *Proceedings of the National Academy of Sciences*, 115(38):E9015–E9024.
- 638 Lu, Z., Doerig, A., Bosch, V., Krahmer, B., Kaiser, D., Cichy, R. M., and Kietzmann, T. C. (2025). End-to-end topographic
639 networks as models of cortical map formation and human visual behaviour. *Nature Human Behaviour*.
- 640 Mahon, B. Z. (2022). Domain-specific connectivity drives the organization of object knowledge in the brain. *Handbook of
641 Clinical Neurology*, 187:221–244.
- 642 Mahon, B. Z., Anzellotti, S., Schwarzbach, J., Zampini, M., and Caramazza, A. (2009). Category-specific organization in
643 the human brain does not require visual experience. *Neuron*, 63(3):397–405.
- 644 Mahon, B. Z. and Caramazza, A. (2011). What drives the organization of object knowledge in the brain? *Trends in
645 Cognitive Sciences*, 15(3):97–103.
- 646 Mahon, B. Z., Milleville, S. C., Negri, G. A. L., Rumia, R. I., Caramazza, A., and Martin, A. (2007). Action-related
647 properties shape object representations in the ventral stream. *Neuron*, 55(3):507–520.
- 648 Margalit, E., Lee, H., Finzi, D., DiCarlo, J. J., Grill-Spector, K., and Yamins, D. L. K. (2024). A unifying framework for
649 functional organization in early and higher ventral visual cortex. *Neuron*, 112(14):2435–2451.e7.
- 650 McLaughlin, T. and O’Leary, D. D. (2005). MOLECULAR GRADIENTS AND DEVELOPMENT OF RETINOTOPIC
651 MAPS. *Annual Review of Neuroscience*, 28(1):327–355.
- 652 Meyers, E. M., Borzello, M., Freiwald, W. A., and Tsao, D. (2015). Intelligent Information Loss: The Coding of Facial Identity,
653 Head Pose, and Non-Face Information in the Macaque Face Patch System. *Journal of Neuroscience*, 35(18):7069–7081.
- 654 Miranda-Dominguez, O., Feczkó, E., Grayson, D. S., Walum, H., Nigg, J. T., and Fair, D. A. (2018). Heritability of the
655 human connectome: A connectotyping study. *Network Neuroscience*, 2(2):175–199.
- 656 Mitchison, G. (1991). Neuronal branching patterns and the economy of cortical wiring. *Proceedings of the Royal Society of
657 London. Series B: Biological Sciences*, 245(1313):151–158.
- 658 Nelson, M. E. and Bower, J. M. (1990). Brain maps and parallel computers. *Trends in Neurosciences*, 13(10):403–408.
- 659 Park, J., Soucy, E., Segawa, J., Mair, R., and Konkle, T. (2024). Immersive scene representation in human visual cortex
660 with ultra-wide-angle neuroimaging. *Nature Communications*, 15(1):5477.
- 661 Plaut, D. C. and Behrmann, M. (2011). Complementary neural representations for faces and words: A computational
662 exploration. *Cognitive Neuropsychology*, 28(3&4):251–275.

- 663 Powell, L. J., Kosakowski, H. L., and Saxe, R. (2018). Social Origins of Cortical Face Areas. *Trends in cognitive sciences*,
664 22(9):752–763.
- 665 Prince, J. S., Alvarez, G. A., and Konkle, T. (2024). Contrastive learning explains the emergence and function of visual
666 category-selective regions. *Science Advances*, 10(39):eadl1776.
- 667 Qian, X., Dehghani, A. O., Farahani, A. B., and Bashivan, P. (2024). Local lateral connectivity is sufficient for replicating
668 cortex-like topographical organization in deep neural networks.
- 669 Raj, A. and Chen, Y.-h. (2011). The Wiring Economy Principle: Connectivity Determines Anatomy in the Human Brain.
670 *PLoS ONE*, 6(9):e14832.
- 671 Rajimehr, R., Firooz, A., Rafipoor, H., Abbasi, N., and Duncan, J. (2022). Complementary hemispheric lateralization of
672 language and social processing in the human brain. *Cell reports*, 41(6):111617.
- 673 Ramón y Cajal, S. et al. (1899). *Textura del sistema nervioso del hombre y de los vertebrados: estudios sobre el plan
674 estructural y composición histológica de los centros nerviosos adicionados de consideraciones fisiológicas fundadas en los
675 nuevos descubrimientos. Volumen I.* Madrid: Nicolás Moya.
- 676 Ratan Murty, N. A., Teng, S., Beeler, D., Mynick, A., Oliva, A., and Kanwisher, N. (2020). Visual experience is not
677 necessary for the development of face-selectivity in the lateral fusiform gyrus. *Proceedings of the National Academy of
678 Sciences*, 117(37):23011–23020.
- 679 Richards, B. A., Lillicrap, T. P., Beaudoin, P., Bengio, Y., Bogacz, R., Christensen, A., Clopath, C., Costa, R. P., de
680 Berker, A., Ganguli, S., Gillon, C. J., Hafner, D., Kepcs, A., Kriegeskorte, N., Latham, P., Lindsay, G. W., Miller,
681 K. D., Naud, R., Pack, C. C., Poirazi, P., Roelfsema, P., Sacramento, J., Saxe, A., Scellier, B., Schapiro, A. C., Senn, W.,
682 Wayne, G., Yamins, D., Zenke, F., Zylberberg, J., Therien, D., and Kording, K. P. (2019). A deep learning framework for
683 neuroscience. *Nature Neuroscience*, 22(11):1761–1770.
- 684 Robinson, E. C., Jbabdi, S., Glasser, M. F., Andersson, J., Burgess, G. C., Harms, M. P., Smith, S. M., Van Essen, D. C.,
685 and Jenkinson, M. (2014). MSM: A new flexible framework for multimodal surface matching. *Neuroimage*, 100:414–426.
- 686 Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M.,
687 Berg, A. C., and Fei-Fei, L. (2015). ImageNet Large Scale Visual Recognition Challenge. *International Journal of
688 Computer Vision*, 115(3):211–252.
- 689 Saygin, Z. M., Osher, D. E., Koldewyn, K., Reynolds, G., Gabrieli, J. D., and Saxe, R. R. (2012). Anatomical connectivity
690 patterns predict face selectivity in the fusiform gyrus. *Nature Neuroscience*, 15(2):321–327.
- 691 Saygin, Z. M., Osher, D. E., Norton, E. S., Youssoufian, D. A., Beach, S. D., Feather, J., Gaab, N., Gabrieli, J. D.,
692 and Kanwisher, N. (2016). Connectivity precedes function in the development of the visual word form area. *Nature
693 Neuroscience*, 19(9):1250–1255.
- 694 Schwartz, E. L. (1980). Computational anatomy and functional architecture of striate cortex: A spatial mapping approach
695 to perceptual coding. *Vision Research*, 20(8):645–669.
- 696 Sereno, M. I., Dale, A. M., Reppas, J. B., Kwong, K. K., Belliveau, J. W., Brady, T. J., Rosen, B. R., and Tootell, R. B. H.
697 (1995). Borders of Multiple Visual Areas in Humans Revealed by Functional Magnetic Resonance Imaging. *Science*,
698 268(5212):889–893.
- 699 Silson, E. H., Groen, I. I. A., and Baker, C. I. (2021). Direct comparison of contralateral bias and face/scene selectivity in
700 human occipitotemporal cortex. *BioRxiv*.
- 701 Silson, E. H., Groen, I. I. A., Kravitz, D. J., and Baker, C. I. (2016). Evaluating the correspondence between face-, scene-,
702 and object-selectivity and retinotopic organization within lateral occipitotemporal cortex. *Journal of Vision*, 16(6):14.
- 703 Spiridon, M. and Kanwisher, N. (2002). How distributed is visual category information in human occipito-temporal cortex?
704 An fMRI study. *Neuron*, 35(6):1157–1165.
- 705 van den Hurk, J., Van Baelen, M., and Op de Beeck, H. P. (2017). Development of visual category selectivity in ventral
706 visual cortex does not require visual experience. *Proceedings of the National Academy of Sciences*, 114(22):E4501–E4510.
- 707 Wang, L., Mruczek, R. E., Arcaro, M. J., and Kastner, S. (2015). Probabilistic Maps of Visual Topography in Human
708 Cortex. *Cerebral Cortex*, 25(10):3911–3931.
- 709 Wang, P. and Cottrell, G. W. (2017). Central and peripheral vision for scene recognition: A neurocomputational modeling
710 exploration. *Journal of Vision*, 17(4):9.
- 711 Watson, A. B. (2014). A formula for human retinal ganglion cell receptive field density as a function of visual field location.
712 *Journal of Vision*, 14(7):15.

- 713 Yao, M., Wen, B., Yang, M., Guo, J., Jiang, H., Feng, C., Cao, Y., He, H., and Chang, L. (2023). High-dimensional
714 topographic organization of visual features in the primate temporal lobe. *Nature Communications*, 14(1):5931.
- 715 Zhang, Y., Zhou, K., Bao, P., and Liu, J. (2024). A biologically inspired computational model of human ventral temporal
716 cortex. *Neural Networks*, 178:106437.
- 717 Zhou, B., Lapedriza, A., Khosla, A., Oliva, A., and Torralba, A. (2018). Places: A 10 Million Image Database for Scene
718 Recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 40(6):1452–1464.

5 Appendix

Non-retinotopically constrained models

In our main analyses, we presented results with a non-retinotopically constrained model that used a higher λ_w^{local} value to accommodate the reduced overall wiring cost resulting from $\lambda^{\text{ret}} = 0$. Here, we compare this model with a model trained with $\lambda_w^{\text{local}} = 0.1$, as in the retinotopically constrained model (where $\lambda_w^{\text{ret}} = 0.01$), shown in Figure S1. We see that the additional local wiring cost was necessary to encourage a smoother topographic organization, but in both cases, no global consistency is seen due to the lack of retinotopic wiring constraints.

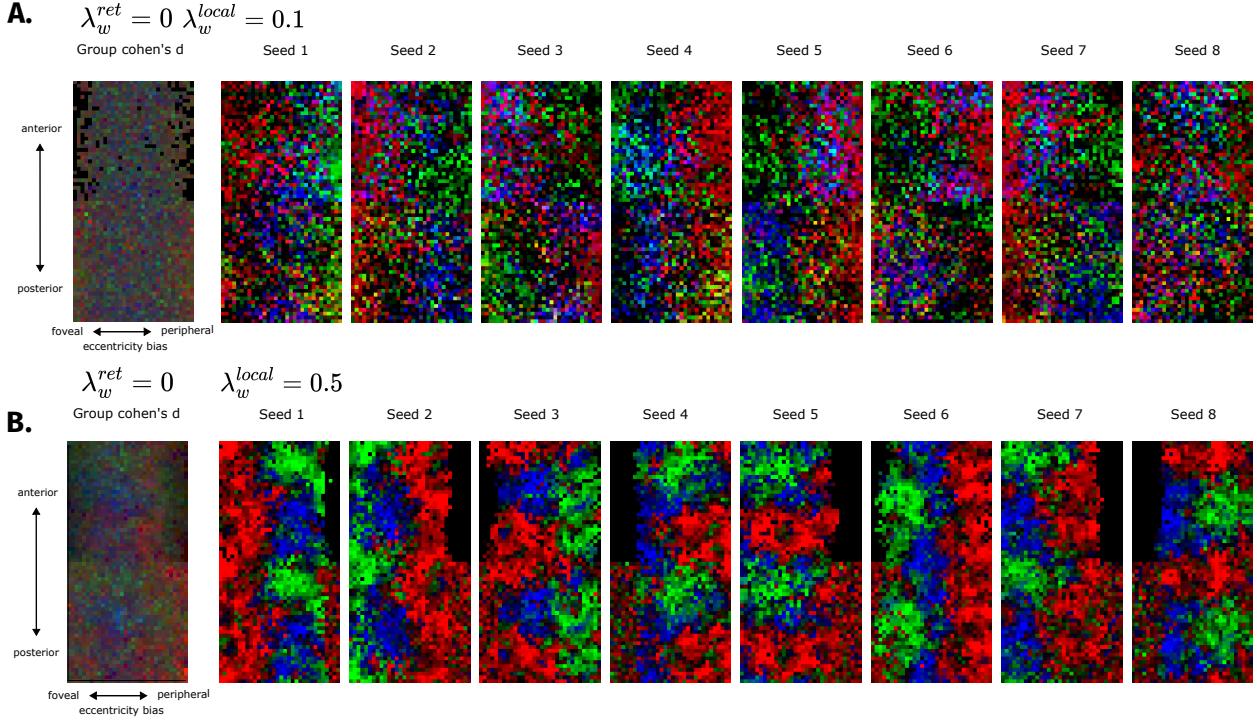


Figure S1: Removing the retinotopic constraint by setting $\lambda^{\text{ret}} = 0$. **A.** Using the same value of $\lambda_w^{\text{local}} = 0.1$ as in the retinotopically constrained model (where $\lambda_w^{\text{ret}} = 0.01$) Left: group effect size (d) map. Right: individual model selectivity effect size (d) maps for 5 random model seeds. **B.** Same organization as **A.**, but for a set of models trained with a stronger local wiring cost ($\lambda_w^{\text{local}} = 0.5$) to make up for the lack of retinotopic wiring cost.

Differential viewing biases differentially bias topographic organization

Thus far we have explored viewing biases in which objects and faces were presented at distributionally smaller sizes than scenes. Nevertheless, faces have taken the foveally biased area consistently when retinotopic connectivity constraints were enabled, both with large and small amounts of within-domain viewing variation. To explore how the emergent global topographic layout depends on the viewing biases, we trained 3 additional models with different sets of viewing biases constructed from the distributions used in the model trained with less within-domain variation (Figure 4B), shown in Figure S2. We first plot the model from Figure 4B as a reference, in Figure S2A. This model shows a similar pattern of lesion deficits in pVTC as in the main model with greater within-domain viewing bias variation.

Next, we plot the first model variant in Figure S2B. This model was trained without any differential viewing biases: faces, objects, and scenes were all presented at a shared large size. Interestingly, this produced the same rough global layout as in the main model with differential viewing biases; however, the effects of eccentricity-biased lesions were less sharply centered, suggesting a greater overlap of representations. Assuming that this result did not emerge by chance, it suggests that the image features of faces may lend themselves towards affinity towards the foveal representations, perhaps due to their canonical structure or a greater importance of foveal features at the same retinotopic size. In the next model variant, we trained objects at a small size, and faces at the same large size as scenes, to determine whether this would lead objects to take the fovea. Interestingly, domain-selective topography and lesions revealed a more intermixed foveal territory – with representations for faces and objects intermixed. We hypothesized that this was due to a pressure for objects and scenes to partially share resources. To test this hypothesis, we trained a final model variant with scenes presented at the same smaller size as objects. Despite no change in the sizes of faces or objects, in support of our hypothesis, this led objects to more convincingly take over the foveal territory, with scenes emerging intermediate and faces spread across intermediate and peripheral territory. Collectively, these results suggest that viewing biases, combined with retinotopic constraints, indeed influence but do not uniquely determine the global topographic layout; the interrelationships between domains, and the informativeness of features at different eccentricities for fixed viewing biases, also appear to influence the layout.

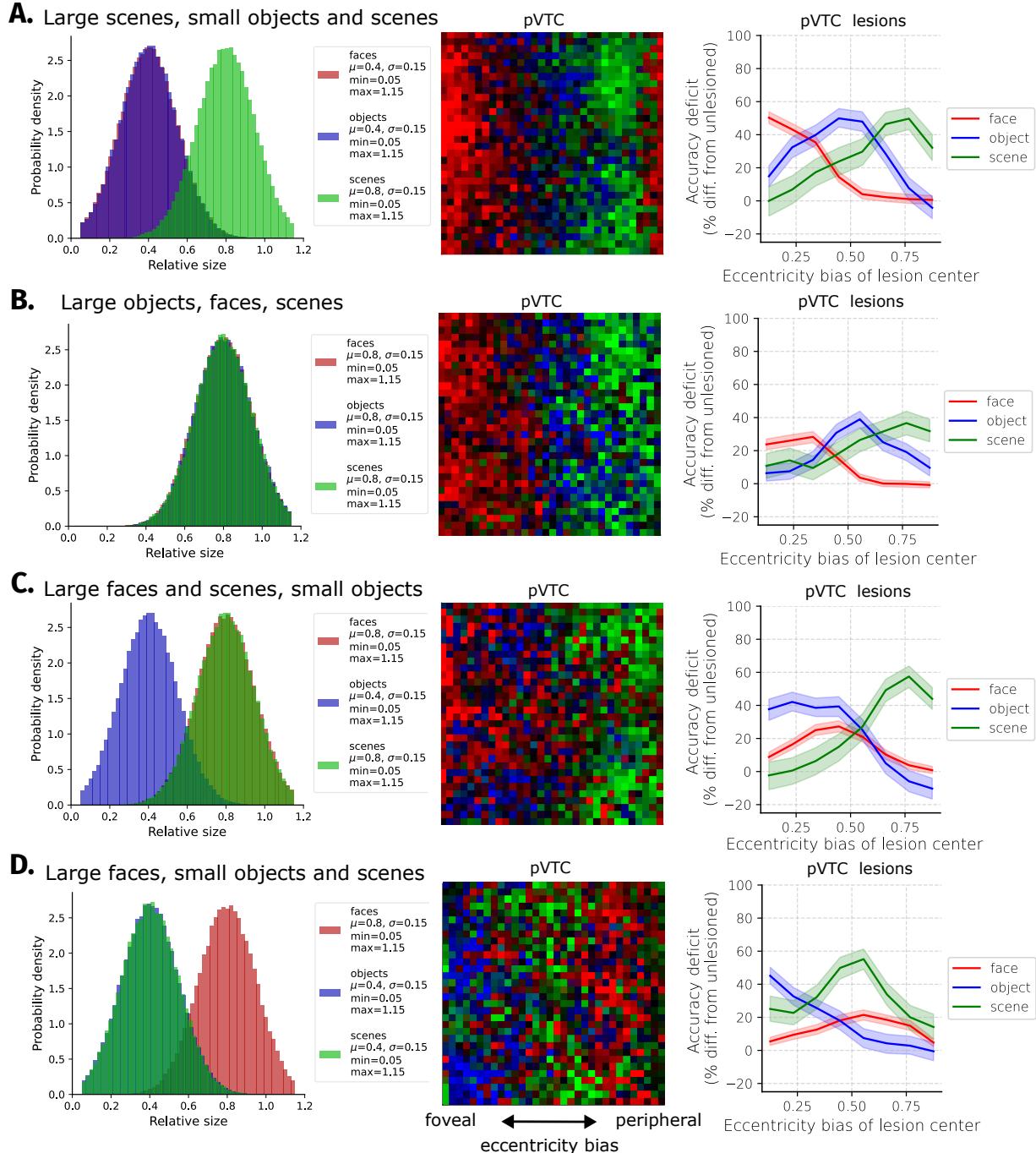


Figure S2: Analyzing emergent topography under different viewing biases. In each panel, we plot the viewing size distributions (left), pVTC domain-level topography (middle), and damage resulting from eccentricity-biased pVTC lesions (right).