# A Unifying Principle for the Functional Organization of Visual Cortex

Eshed Margalit <sup>1</sup>,<sup>[]</sup>, Hyodong Lee<sup>2</sup>, Dawn Finzi<sup>3,4</sup>, James J. DiCarlo <sup>2,5,6</sup>, Kalanit Grill-Spector <sup>3,7,\*</sup>, and Daniel L. K.
 Yamins<sup>3,4,7,\*</sup>

<sup>5</sup> <sup>1</sup>Neurosciences Graduate Program, Stanford University, Stanford, CA 94305

- <sup>6</sup> <sup>2</sup>Department of Brain and Cognitive Sciences, Massachusetts Institute of Technology, Cambridge, MA 02139
- <sup>7</sup> <sup>3</sup>Department of Psychology, Stanford University, Stanford, CA 94305
- <sup>8</sup> <sup>4</sup>Department of Computer Science, Stanford University, Stanford, CA 94305
- <sup>9</sup> <sup>5</sup>McGovern Institute for Brain Research, Massachusetts Institute of Technology, Cambridge, MA 02139
- <sup>10</sup> <sup>6</sup>Center for Brains Minds and Machines, Massachusetts Institute of Technology, Cambridge, MA 02139
- <sup>11</sup> <sup>7</sup>Wu Tsai Neurosciences Institute, Stanford University, Stanford, CA 94305

<sup>12</sup> co-senior author

A key feature of many cortical systems is functional organization: the arrangement of neurons with specific 13 functional properties in characteristic spatial patterns across the cortical surface. However, the principles 14 underlying the emergence and utility of functional organization are poorly understood. Here we develop 15 the Topographic Deep Artificial Neural Network (TDANN), the first unified model to accurately predict the 16 functional organization of multiple cortical areas in the primate visual system. We analyze the key factors 17 responsible for the TDANN's success and find that it strikes a balance between two specific objectives: 18 achieving a task-general sensory representation that is self-supervised, and maximizing the smoothness of 19 responses across the cortical sheet according to a metric that scales relative to cortical surface area. In 20 turn, the representations learned by the TDANN are lower dimensional and more brain-like than those in 21 models that lack a spatial smoothness constraint. Finally, we provide evidence that the TDANN's functional 22 organization balances performance with inter-area connection length, and use the resulting models for 23 a proof-of-principle optimization of cortical prosthetic design. Our results thus offer a unified principle 24 for understanding functional organization and a novel view of the functional role of the visual system in 25 particular. 26

27 Correspondence: eshed.margalit@gmail.com

## <sup>28</sup> Introduction

Neurons in sensory cortical systems support two kinds of measurements: their response patterns as a function 29 of stimulus input and their spatial arrangement across the cortical surface. The confluence of these observations 30 is referred to as functional organization, the reproducible spatial arrangement of neurons within a cortical area 31 according to their response properties. Functional organization is among the most ubiquitous of neuroscience 32 findings, appearing in the topographic maps of the visual system [1], and in auditory [2], parietal [3], sensorimotor [4], 33 and entorhinal areas [5, 6]. These organized structures anchor our understanding of cortical development, function, 34 and dysfunction, yet it remains a mystery what processes govern their emergence, and what computational function 35 they serve. 36

Any theory of functional organization must explain both neuronal response properties and the physical arrangement 37 of neurons within a cortical area. Furthermore, a unified theory should account for the observed functional 38 organization in multiple cortical areas. Prior computational models of the organization within single cortical areas 39 have been developed [7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22], but these approaches do not 40 generalize to multiple cortical areas. Moreover, many of these models operate from a hand-crafted set of stimulus 41 features, and thus cannot explain how neuronal response properties are learned from realistic sensory inputs. 42 On the other hand, deep artificial neural networks (DANNs) trained with large quantities of naturalistic data are 43 increasingly being used to model neuronal responses in regions responsible for vision, audition, and language 44 processing [23, 24, 25, 26, 27, 28, 29, 30, 31]. However, standard DANNs impose no spatial arrangement among 45 model units that differ in their stimulus tuning, and thus cannot explain the observed organization of neurons across 46 the cortical surface. 47

Here, we introduce the Topographic Deep Artificial Neural Network (TDANN), a unified framework for predicting 48 functional organization in sensory systems. The TDANN implements the hypothesis that neural systems are 49 optimized to address two key goals: they must support ecologically-relevant behaviors by producing useful neural 50 representations [32], and they must do so in a biophysically efficient manner, using as few resources as possible. A 51 critical component of biophysical efficiency is the minimization of neuronal wiring length, which is theorized to result 52 in the smooth topographic organization observed in many cortical areas [33, 19, 18]. The TDANN begins with a 53 standard DANN and spatially augments it by embedding each layer's units in a two-dimensional simulated cortical 54 sheet. The TDANN then optimizes a composite objective function with two components: a functional objective 55 that drives the learning of useful representations, and a spatial constraint that encourages efficiency with smooth 56 response patterns across the simulated cortical sheet. We test this framework in the primate ventral visual stream, 57 a cortical system in which functional organization has been extensively documented. 58 The ventral stream is a hierarchical series of cortical areas that support visual recognition, beginning with primary 59

visual cortex (V1) and ascending through intermediate areas (e.g., V4) to high-level regions: inferotemporal (IT) 60 cortex in macaques and ventral temporal cortex (VTC) in humans. Well-known neuronal response properties in V1 61 include tuning to edge orientation [1, 34, 35], spatial frequency [36], and color [37, 38]. These response properties 62 are coupled with topographic signatures: orientation preferences form a smooth cortical map with pinwheel-like 63 discontinuities [39, 40, 41, 42, 43]; spatial frequency tuning is organized in a guasi-periodic map with isolated 64 low-frequency domains [42, 43, 44]; and color-preferring neurons cluster in punctate blobs [38] across the V1 surface. 65 Higher-level regions such as primate IT [45, 46, 47, 48] and the analogous human VTC contain neurons with stronger 66 responses for items of specific categories vs. others (e.g., faces vs non-faces), a property known as category 67 selectivity. A core characteristic of functional organization in IT [48, 49] and VTC [50, 51, 52, 53, 54, 55, 56] is that 68 neurons selective for certain ecologically-relevant categories - including faces, places, limbs, and visual wordforms 69 – cluster into spatial patches, with characteristic patch sizes, counts, and relative inter-patch distances. 70

We find that the TDANN reproduces the functional organization of the ventral stream, including smooth orientation 71 maps with pinwheels in an earlier model layer, and category-selective patches in a later layer that match the number, 72 size, and relative geometry of patches in human VTC. To understand the principles underlying the emergence of the 73 ventral stream's functional organization, we then test which specific functional and spatial constraints of the TDANN 74 are critical to the TDANN's success by insantiating alternative models and measuring their capacity to predict neural 75 data. We find that the specific combination of task and spatial objectives that best matches the functional organization 76 of the ventral stream also makes learned representations more brain-like by constraining their intrinsic dimensionality. 77 The TDANN learns these representations while minimizing the network's inter-layer wiring length, suggesting that 78 brain-like functional organization effectively balances performance with metabolic costs. 79

Finally, because the the TDANN accurately predicts the functional organization of the ventral stream, it provides an exciting new platform for simulating experiments that are challenging to implement empirically. As a proof of principle, we perform *in silico* experiments simulating the effect of cortical microstimulation devices that vary in their

- spatial precision and cortical coverage. Taken together, our results show that the TDANN serves both as a unified
- explanation for the functional organization of the visual system and as a platform to fuel discovery in neuroscience.

## **Results**

90

91

92

93

94

95

96

97

98

#### <sup>86</sup> Instantiating models that balance task performance with spatial smoothness

<sup>87</sup> Building on optimization-based approaches in computational neuroscience [57, 58], we seek a model architecture

and objective function that generate a neural network which matches the neuronal responses and topography of the

<sup>89</sup> primate ventral visual stream.

Because standard DANNs have no within-area spatial structure beyond retinotopy, we must augment their architecture to model spatial topography. Specifically, we take the ResNet-18 architecture [59], a DANN that achieves strong object recognition performance and accurate prediction of neuronal responses throughout the ventral visual stream [30], and augment it by embedding the units of each convolutional layer into a two-dimensional simulated cortical sheet (Figure 1a). Given that neurons in visual cortex are organized retinotopically at birth [60], we assign model unit positions retinotopically, such that units responding to similar regions of the input images are nearby in the simulated cortical sheet. Then, prior to training, unit positions are locally shuffled to circumvent limitations of weight-shared convolution (see Methods). The size of the simulated cortical sheet in each layer is anchored by estimates of cortical surface area in the human ventral visual stream (Figure 1a). We refer to the resulting model as

<sup>99</sup> the Topographic DANN (TDANN).



**Figure 1.** Constructing a unified model of the functional and spatial constraints of ventral visual cortex. (a) TDANNs are a family of deep artificial neural networks whose units are assigned positions in a two-dimensional simulated cortical sheet in each layer. Position assignments are retinotopic, such that location in the cortical sheet corresponds to position in the visual field. Each individual dot is a single model unit. The degree of overlap between a unit's spatial receptive field (RF) and the purple square marked on the input image is indicated by the shade of purple; RFs from gray units do not overlap the marked region at all. The TDANN is trained to minimize the sum of a task loss and a spatial loss (SL).  $\alpha$  is a free parameter controlling the relative weight of the SL. (b) The SL encourages nearby units to develop strong response correlations. Plotted: pairwise similarity of unit responses as a function of pairwise cortical distance in the final layer of a TDANN model; each dot represents one pair of units. (c) The TDANN is evaluated on a battery of quantitative benchmarks that measure its correspondence to topographic features throughout the ventral visual stream. Left: orientation preference map in the V1-like TDANN layer (see Figure 2 for details). Right: category selectivity map in the VTC-like layer (see Figure 3 for details).

Having selected the architecture, our goal is to discover the objective whose optimization yields an accurate model of both response properties and their topographic arrangement. The core of the TDANN approach is a composite objective that is a weighted sum of two components: a task objective encouraging the learning of behaviorally-useful functional representations, and a spatial objective driving the emergence of topographic properties. Following recent progress in training neural networks without explicit category labels [61, 62], we use an unsupervised algorithm that performs *contrastive self-supervision*, SimCLR [63], as the task objective. For the spatial loss (SL), we introduce an objective that encourages nearby pairs of units to have more correlated responses than distant pairs of units (Figure 1b, see Methods). The SL is computed separately in each convolutional layer, then summed across layers for each

<sup>108</sup> batch of training data:

100

10

102

103

104

105

106

107

$$TDANN Loss = L_{task} + \sum_{k \in layers} \alpha_k SL_k$$
(1)

where  $\alpha_k$  is the weight of the spatial loss in the *k*th layer, set to  $\alpha_k = 0.25$  for all layers. The TDANN architecture is trained to optimize this objective using conventional back-propagation with stochastic gradient descent.

Training the TDANN on ImageNet [64] resulted in successful minimization of both task and spatial losses 111 (Supplementary Figure S1). We tested if adding the spatial loss interferes with visual representation learning by 112 measuring the model's object categorization performance with a linear readout. Categorization accuracy was slightly 113 but significantly lower for the TDANN (median across random initialization seeds = 43.9%) than "Task Only" models 114 with no spatial loss ( $\alpha = 0$ , median = 48.5%; Mann-Whitney U = 25, p = .008). Despite the modest decrease in 115 categorization performance, adding the spatial loss term had the intended effect: in each layer, the correlation 116 between units' responses increased with spatial proximity (Supplementary Figure S1c,d). To determine if this 117 learned correlation structure corresponds to brain-like topographic maps, we constructed a battery of quantitative 118 benchmarks comparing model predictions with neural data in primary visual cortex (V1) and ventral temporal cortex 119 (VTC), (Figure 1c). To compare against these benchmarks, we needed to identify the TDANN layers that would be 120 our models of V1 and VTC. As in prior work [28, 25], we find that earlier model layers best predict V1 responses and 121 later layers best predict responses in higher visual cortex (Supplementary Figure S2). Accordingly, we designate the 122 fourth and ninth convolutional layers as the "V1-like" and "VTC-like" layers, respectively. 123

#### 124 The TDANN predicts the functional organization of primary visual cortex

Neurons in primate V1 are organized into maps of preferred stimulus orientation, spatial frequency, and color 125 [38, 43, 65]. Because high-resolution data at the scale necessary to visualize these maps is not available for 126 human V1, we compare the TDANN to macaque V1 data using scale-invariant metrics. We tested if the V1-like 127 TDANN layer captures the functional organization of macague V1 with three kinds of guantitative benchmarks. First, 128 we evaluate functional correspondence by asking if model units in the TDANN V1-like layer have similar preferred 129 orientations and orientation tuning strengths as neurons in macaque V1. Second, we assay the structure of cortical 130 maps by measuring pairwise similarity of tuning for orientations, spatial frequencies, and colors as a function of 131 cortical distance. Third, we measure the density of pinwheel-like discontinuities in the orientation preference map, 132 a hallmark of V1 functional organization in many species [41, 66]. In addition to the TDANN, we also evaluate four 133 control models on these benchmarks: the Unoptimized TDANN, in which model weights and unit positions are left 134 randomly initialized, the Task Only variant in which  $\alpha = 0$ , and two kinds of self-organizing maps (SOMs), which have 135 been proposed as models of V1 functional organization [11, 10]. We refer to the traditional SOM in which feature 136 dimensions are manually predetermined (as in Swindale and Bauer [11]), as the Hand-Crafted SOM, and a novel 137 SOM that organizes the output of an AlexNet V1-like layer (inspired by Doshi and Konkle [13], Zhang et al. [12]) as 138 the DNN-SOM. 139

The TDANN matches orientation tuning in V1 We measured orientation tuning strength by presenting a set of 140 oriented sine grating images to the model (Figure 2a), computing a tuning curve for each unit, and calculating 141 the circular variance (CV; lower values for sharper tuning) of each tuning curve. Setting a selectivity threshold of CV 142 < 0.6, we find that the TDANN V1-like layer has a significantly greater proportion of selective units (range across 143 model seeds: [20%, 31%]) than Unoptimized models ([1%, 3%]; Mann-Whitney U = 25; p = .008, Figure 2b), but 144 fewer than Task Only models ([35%, 50%]; U = 25; p = .008) or macaque V1 (45%; Supplementary Figure S3c). In 145 contrast, neither the Hand-Crafted SOM nor the DNN-SOM exhibited any units with sharp orientation tuning. We 146 also find that TDANN and Task Only models (but not SOMs or Unoptimized models) show an over-representation of 147 cardinal orientations (0 and 90 degrees) as in macaque V1 [35] (Supplementary Figure S3b, see also Henderson 148 and Serences [67]). 149

The TDANN predicts the arrangement of orientation-selective V1 neurons To evaluate whether the TDANN V1-like layer captures the topographic properties of macaque V1, we consider the spatial distribution of orientation-selective units – the orientation preference map (OPM) – and find a smooth progression of preferred orientations that

resembles macaque V1 (Figure 2c, d). Following prior work [68, 69, 70], we quantify this structure by measuring the 153 absolute pairwise difference in preferred orientation as a function of cortical distance. In both the TDANN and 154 macaque V1 (data from Nauhaus et al. [43]), we find that nearby units have smaller differences in orientation 155 preference than distant pairs (Figure 2e). In contrast, orientation preference similarity does not vary with cortical 156 distance in Task Only or Unoptimized models, and both the Hand-Crafted and DNN-SOMs exhibit OPMs with 157 abnormally high orientation tuning similarity (Figure 2e, Supplementary Figure S3). We summarize these profiles 158 by computing a *smoothness score* that measures the increase in tuning similarity for nearby unit pairs compared to 159 distant unit pairs. Smoothness of TDANN OPMs ([min, max] across random initialization: [.64, .83]) was consistent 160 with macaque V1 (.68); however, OPMs in the Hand-Crafted SOM ([.92, .92]) and DNN-SOMs ([.81, .86]) were 161 smoother than in macaque V1. In turn, macaque V1 OPMs were smoother than Unoptimized ([.03, .04]) and 162 Task Only ([.28, .39]) models. Jointly comparing each model to macaque V1 orientation tuning strength and OPM 163 smoothness highlights that the TDANN is the only model class that satisfies both criteria (Figure 2j). 164



Figure 2. The TDANN reproduces V1-like topography. (a) Example sine grating stimuli used to assess tuning for orientation, spatial frequency, and color. (b) Orientation tuning curves (top) and spatial frequency tuning curves (bottom) for four example units in the V1-like layer. (c) Smoothed orientation preference map (OPM) in the V1-like layer of the TDANN. Box corresponds to inset at right, where individual model units are labeled by their preferred orientation. Results for additional model seeds shown in Supplementary Figure S10. (d) OPMs for Macaque V1 (data from Nauhaus et al. [43]), TDANN, and an Unoptimized control model. (e) Left: Pairwise difference in preferred orientations as a function of pairwise cortical distance, normalized to the chance level expected by random sampling of pairs. Right: Map smoothness for OPMs in macague V1 (dashed green line, data from Nauhaus et al. [43]) and four candidate models: the TDANN (purple), the Hand-Crafted self-organizing map (SOM, squares), deep neural network SOM (DNN-SOM, plus signs), and Task Only (diamonds) trained without the spatial term of the loss function. Error bar: 95% CI across random model seeds and sampling of cortical neighborhoods. (f) Spatial frequency preference, shown for the same region of the TDANN V1-like layer and macaque V1 as in panel (d). (g) Change in preferred spatial frequency as a function of cortical distance, normalized to chance, for macaque V1 and each model type. (h) Preference for chromatic stimuli for the same region of the TDANN V1-like layer. Dark-colored dots: stronger responses to chromatic than achromatic gratings. Macague data: reconstruction of cytochrome oxidase staining data from Livingstone and Hubel [38]. (i) Fraction of units differing in their chromatic preference as a function of cortical distance, normalized to chance. (i) Similarity of models to the distribution of orientation tuning strengths in macaque V1 (data from Ringach et al. [34]) on the x-axis, and similarity to the smoothness of macaque OPMs (data from Nauhaus et al. [43]) on the y-axis. Multiple markers of the same type indicate different random initial seeds for each model. A value of 1.0 (dashed green) indicates perfect correspondence. (k) Density of pinwheels detected in TDANNs, Hand-Crafted SOMs, Task Only models, and Unoptimized models. Error bars: CI across random model seeds. Green: putative macaque V1 pinwheel density.

As a more stringent test of OPM structure, we counted the number of periodic pinwheel-like discontinuities in the 165 OPM [41] and compared to the expected value of  $\sim$ 3.1 pinwheels /  $mm^2$  in macaque V1 [66]. Multiple pinwheels are 166 apparent in both the TDANN and the Hand-Crafted SOM (Figure 2k). To facilitate quantitative comparison across 167 models, we compute pinwheel density - the number of pinwheels normalized by the average spacing between 168 "columns", i.e. clusters of units preferring the same orientation. We find that the TDANN has lower pinwheel density 169 (range across seeds = [2.0, 2.3] pinwheels / column spacing<sup>2</sup>) than macaque V1, but significantly higher than either 170 the Task Only ([0.2, 0.8]; Mann-Whitney U = 25, p = .008) or Unoptimized models (0 pinwheels; Figure 2k). The 171 Hand-Crafted SOM has higher pinwheel density ([3.7, 4.5]) than the TDANN, but the DNN-SOM has no detectable 172 pinwheels. Although the TDANN has pinwheel density approaching that of macaque V1, we note that the orientation 173 174 column spacing in the TDANN ( $\sim 3.5$ mm width) does not match macaque V1 ( $\sim 1$ mm). This mismatch, caused in part by our commitment of the TDANN as a model of human visual cortex and not macaque visual cortex, can 175 also be overcome by increasing the number of units in the network at the expense of increased computational cost 176 (Supplementary Figure S5). 177

The TDANN predicts maps of spatial frequency and color preference in V1 While OPMs are the best-studied feature 178 of V1 functional organization, the cortical sheet simultaneously accommodates organized maps of spatial frequency 179 [43] and chromatic tuning [71, 38]. An accurate model of V1 should also predict these aspects of V1 functional 180 organization. We compared spatial frequency preference maps in macaque V1 (data from [43]) and in the TDANN 181 V1-like layer and found a smooth progression of preferred spatial frequency in both (Figure 2f). Quantifying the 182 difference in spatial frequency tuning as a function of cortical distance indicates that the TDANN map ([min, max] 183 of smoothness across random initializations = [.38, .54]) is as smooth as the map in macague V1 (0.53; Figure 2g), 184 whereas maps from Task Only ([.23, .36]) and Unoptimized models ([.02, .03]) are far less smooth than macague V1, 185 and both the Hand-Crafted SOM ([.79, .81]) and the DNN-SOM ([.83, .86]) are again far smoother than the neural 186 data. We observe similar results for maps of chromatic preference (Figure 2h, i), where comparisons are made 187 to imaging of cytochrome oxidase (CO) uptake that is prevalent in color-tuned neurons (data from Livingstone and 188 Hubel [38]). In the TDANN chromatic map, the fraction of units with opposite color-tuning increases with cortical 189 distance, again exhibiting comparable smoothness to macaque V1 (TDANN smoothness: [.38, .54], macaque: .53). 190 Together, our analyses demonstrate that the TDANN predicts the multifaceted functional organization of macaque 191

<sup>192</sup> V1, providing a stronger match to neural data than existing models such as the standard Hand-Crafted SOM.

### 193 The TDANN reproduces the functional organization of higher visual cortex

Because benchmarks measuring the topographic similarity between models and higher visual cortex, i.e. primate 194 inferior temporal (IT) and human ventral temporal cortex (VTC), are still underdeveloped, we introduce five 195 quantitative benchmarks that compare both responses and topography. Response properties are compared by 196 measuring the similarity of population category selectivity patterns with representational similarity analysis (RSA; 197 Kriegeskorte et al. [72]), as in Margalit et al. [73], Haxby et al. [74]). Topographic properties are then compared 198 against four complementary benchmarks: 1) the smoothness of category selectivity maps, 2) the number of category 199 selective patches, 3) the area occupied by those patches, and 4) the spatial overlap of units selective for different 200 categories. We compute these metrics for the TDANN's VTC-like layer and for VTC data from eight human subjects 201 in the Natural Scenes Dataset (NSD) [75] (Supplementary Figure S6). We also evaluate two alternative models of 202 VTC topography: an SOM trained on the outputs of a categorization-pretrained AlexNet (DNN-SOM, cf Doshi and 203 Konkle [13], Zhang et al. [12]) and a variant of the Interactive Topographic Network (ITN) that is trained on the same 204 dataset (ImageNet) we used (Blauch et al. [20]: Supplementary Figure S19C). Human subjects and models were all 205 presented a common set of 1,440 object category images [76] composed of five categories: faces, bodies, written 206 characters, places, and objects (cars and instruments). Selectivity was computed as the t-value for each category, 207 for each human voxel and model unit. 208

The TDANN predicts patterns of category selectivity We characterize neuronal responses in VTC by computing a 209 representational similarity matrix (RSM): the similarity between pairs of distributed selectivity patterns to each of the 210 five object categories. The average RSM from human VTC indicates high similarity between patterns of selectivity 211 for faces and bodies, and low similarity between selectivity for faces and places (Figure 3a). The alignment between 212 any two RSMs is computed as Kendall's  $\tau$ . RSMs from different subjects and hemispheres were very similar, with 213 the 95% CI of Kendall's  $\tau = [.72, .75]$ . We then compute RSMs for each model and compare against the human data, 214 finding that some models provide a closer match to human VTC than others (ANOVA  $F(4, 331) = 630; p < 10^{-152}$ ). 215 TDANN RSMs closely mirror those in human VTC ( $\tau$  = [.69, .73]), significantly better than DNN-SOM ( $\tau$  = [.31, .35]; 216 post-hoc Tukey's HSD  $p < 10^{-13}$ ), ITN ( $\tau = [.46, .56]; p < 10^{-13}$ ), Task Only ( $\tau = [.65, .68]; p = .001$ ) and Unoptimized 217  $(\tau = [.11, .14]; p < 10^{-13})$  models (Figure 3b). The similarity between human and TDANN RSMs also depends 218 strongly on the training data being naturalistic. Training on artificial stimuli such as white noise and sine gratings 219 yields RSMs that significantly deviate from the human data (Supplementary Figure S9b). 220

The TDANN predicts category-selectivity maps To compare models against topographic benchmarks, we generate 221 selectivity maps for each of the five object categories (Figure 3c), then quantify their structure by measuring the 222 pairwise difference in selectivity as a function of pairwise cortical distance (Figure 3d). We find that for all categories, 223 the curve computed for TDANN is similar to human VTC, whereas the DNN-SOM and ITN are abiologically smooth, 224 and maps in the Unoptimized and Task Only models lack structure. We summarize category selectivity map structure 225 with the same smoothness metric used in V1 (Figure 3e), and find that TDANN maps were as smooth as those in 226 human VTC (permutation test: p = .30). In contrast, VTC maps were significantly smoother than Task Only or 227 Unoptimized models (ps < .001) and less smooth than the DNN-SOM (p < .001). ITN category selectivity maps were 228





Figure 3. The TDANN predicts the functional organization of higher visual cortex. (a) Representational similarity matrices (RSMs) for the TDANN and human VTC, computed across selectivity maps of the five object categories. Diagonal is blank to indicate trivially perfect correlation. (b) Functional similarity between the TDANN, human VTC, and alternative models, measured as the similarity of RSMs. Green: mean of pairwise human-to-human similarity values. (c) Selectivity (t-value), for each category plotted on the simulated cortical sheet of the VTC-like layer in an example TDANN. Black star: unit whose responses to images in each of the five categories are plotted directly below (individual dots: single images, bar height: mean across images). Scale bar: 1cm. (d) Difference in pairwise selectivity as a function of pairwise cortical distance for units in each of five candidate model types: the TDANN (purple), deep neural network self-organizing map (DNN-SOM; plus markers), interactive topographic network ("ITN", Blauch et al. [20]; circles), Unoptimized ("x" markers), and Task Only (diamond markers). Curves are normalized to the chance level obtained by random sampling of unit pairs. Green: Human data averaged over the eight subjects in the NSD data. Shaded regions: 95% confidence interval across different subsets of units from models trained with different random initial seeds. (e) Smoothness of selectivity maps for each category and each candidate model. Dashed green: mean of human data. (f) Category-selective patches for an example hemisphere in human ventral temporal cortex (VTC), TDANN, a Task Only model (no patches detected), a DNN-SOM, and a reproduction of the "ITN" simulated cortical sheet from [20]. Object categories are indexed by color as in (a) and (c). Examples from different initial random seeds are shown in Supplementary Figure S10. (g) Number of category-selective patches (averaged across categories) for the TDANN, DNN-SOM, and ITN. Dashed green: average of human data. ANOVA for difference in patch count:  $F(5, 179) = 32.7, p < 10^{-22}$ . Post-hoc Tukey's tests: significant difference between VTC and ITN ( $p = 1.2 \times 10^{-5}$ ). (h) Average surface area of category-selective patches. Same plotting conventions as in (f). ANOVA for difference in patch area:  $F(5, 187) = 15.4, p < 10^{-11}$ . Post-hoc Tukey's tests: significant difference between VTC and DNN-SOM ( $p < 10^{-10}$ ). (i) Each human subject and model instance compared to the mean patch area (y-axis) and patch number (x-axis) in the human data. (i) Overlap between face-selectivity and body-selectivity vs. overlap between face-selectivity and place-selectivity, for each human hemisphere (green dots), each TDANN instance (purple dots), the ITN (gray dot), each DNN-SOM (gray plus signs), and Task Only models (gray diamonds).

230 For the remaining topographic benchmarks, we follow the literature by thresholding selectivity maps to find

strongly-selective units (Supplementary Figure S6a-d). Clusters of selective units are identifiable in human VTC, 231 TDANN, the SOM and ITN models, but not in Task Only or Unoptimized models. We use a data-driven approach 232 to automatically identify large contiguous clusters of selective units as "patches" (Figure 3f). We find similar sets 233 of patches in VTC and the TDANN: both contain a small number of patches selective for each category (except for 234 object-selective patches, which are not found in VTC), and the patches are similar in size. Quantitative comparison 235 supports the similarity of human VTC and TDANN: there is no significant difference in patch count (p = 0.99, Figure 236 3g) or patch area (p = 0.67; Figure 3h). In contrast, we find that the ITN has more than twice as many patches as 237 VTC ( $p = 1.2 \times 10^{-5}$ ), although the patches are as large on average as those in VTC (p = 0.99). The DNN-SOM 238 fails to match VTC in the other extreme: while the number of patches in the DNN-SOM is similar to that in VTC 239 (p = 0.15), the patches are too large  $(p < 10^{-10})$ . Joint comparison of models and humans on both patch count and 240 size (Figure 3i) highlights the stronger correspondence between TDANN and human VTC than alternative models. 241

An important hallmark of the functional organization of higher visual cortex is the reproducible spatial arrangement 242 of units selective for different categories. A prominent example is the close proximity of face-selective and 243 body-selective regions [49, 77] and the separation between face- and place-selective regions. A measure of proximity 244 between face- and body-selective regions was previously introduced in Lee et al. [78]. Here we measured the 245 co-occurrence of face-selective and body-selective units (and face-selective and place-selective units) in human 246 VTC with an overlap score that ranges between 1 (face-selectivity perfectly predicts body-selectivity) to 0.5 (no 247 relationship), to 0 (face- and body-selectivity perfectly anti-correlated). As expected, Face-Body overlap scores 248 are high in human VTC (95% CI across subjects and hemispheres: [.66, .72]), whereas Face-Place overlap 249 was significantly lower (95% CI: [.40, .45], Wilcoxon signed-rank test against one-sided alternative W = 136; p =250  $1.5 \times 10^{-5}$ ; Figure 3j). The same pattern is apparent in the TDANN: Face-Body Overlap ([.63, .71]) is significantly 25 higher than Face-Place Overlap ([.14, .26]; W = 15; p = .03). In the ITN, the Face-Body overlap score was lower 252 than in human VTC (.52), but still higher than the Face-Place overlap score (.36). Neither the the DNN-SOM nor the 253 Task Only models had higher Face-Body overlap than Face-Place overlap (Figure 3j; ps > 0.5). 254

To further gain intuition for the tuning profiles of model units, we synthesized images that optimally drive each region of the VTC-like layer. We find that the VTC-like layer smoothly maps object feature space onto the two-dimensional simulated cortical sheet; e.g., face-patches are optimally driven by stimuli with apparent eyes (Supplementary Figure S7). We also tested how the nature of the training dataset affects the accuracy of topographic maps in the TDANN (see Lee et al. [78], Figure 7 for a similar analysis). We find that training the TDANN on natural images (either ImageNet [64] or Ecoset [79]) produces accurate V1-like and VTC-like maps, whereas training on noise or simpler hand-crafted stimuli fails to provide a unified account of ventral stream topography (Supplementary Figure S9).

Together, these results demonstrate that TDANN is the only model to exhibit spatially structured category selectivity that is consistent with a large battery of benchmarks comparing models to human VTC.

### <sup>264</sup> Multiple signatures of functional organization emerge at the same spatial constraint strength

The TDANN optimization framework requires the selection of a single free parameter,  $\alpha$ , the weight of the spatial loss in the training objective. When  $\alpha = 0$  ("Task Only"), spatial information is ignored during training, whereas setting  $\alpha$ too high may encourage pathologically strong correlations that interfere with representation learning. In the results above,  $\alpha$  is set to 0.25. Here, we validate this choice by demonstrating that many benchmarks of neural similarity are simultaneously satisfied by low-to-intermediate values of  $\alpha$ .

Comparison of OPMs in the V1-like layer and category-selectivity maps in the VTC-like layer (Figure 4a) in models 270 trained at 7 different levels of  $\alpha$  shows that functional organization is absent when  $\alpha = 0$ , structured at intermediate 271 values of  $\alpha$ , and deteriorates at the highest values of  $\alpha$ . We quantify the dependence of functional organization 272 on  $\alpha$  with three kinds of benchmarks: functional similarity (Figure 4b), map smoothness (Figure 4c), and presence 273 of topographic phenomena (i.e. pinwheels and patches; Figure 4d). First considering functional similarity, we find 274 that the fraction of V1-like layer units that are orientation selective is closest to macaque V1 when  $\alpha$  is low, and 275 representational similarity between the VTC-like layer and human VTC is maximized at  $\alpha = 0.25$  (Figure 4b). The 276 smoothness of topographic maps is most brain-like at  $\alpha = 0.1$  for OPMs in the V1-like layer and at  $\alpha = 0.25$  for 277 category-selectivity maps in the VTC-like layer (Figure 4c). Finally, we find that the density of pinwheels in the 278 V1-like layer and category-selectivity maps in the VTC-like layer are most similar to measurements in macaque V1 279 and human VTC, respectively, at  $\alpha = 0.25$  (Figure 4d). 280

<sup>281</sup> A specific range of  $\alpha$  values ( $0.1 \le \alpha \le 0.25$ ) thus produces experimentally-observed outcomes across a variety of <sup>282</sup> independent functional and topographic benchmarks in multiple brain areas, suggesting that the  $\alpha$  parameter may <sup>283</sup> provide insights into biophysical mechanisms underlying the emergence of functional organization.



Figure 4. Convergence of multiple benchmarks indicates a balancing between functional and spatial constraints. (a) Topographic maps in the V1-like (top row) and VTC-like layer (bottom row) of TDANN models trained at different levels of the spatial weight  $\alpha$ . Top: Orientation map structure and pinwheels become apparent at  $\alpha > 0.1$  and persist until  $\alpha = 1.25$ . Dots: estimated pinwheel locations; black: clockwise, white: counterclockwise. Bottom: Category selectivity maps, with selective units (t > 12) colored according to their preferred category. (b) Functional correspondence to neural data as a function of  $\alpha$ . Top: Fraction of units strongly orientation selective (circular variance  $\leq 0.6$ ) in the V1-like layer. Dashed green: value measured in macaque V1 (from Ringach et al. [34]). Dashed gray: mean value for Unoptimized models. Shaded regions: 95% CI across multiple initial random seeds. Bottom: Representational similarity between the VTC-like layer and human VTC (as in Figure 3). Error region indicates 95% CI across model seeds and human hemispheres. In both plots, the vertical line at  $\alpha = 0.25$  marks the default value used in prior figures. (c) Topographic map smoothness is a function of  $\alpha$ . Top: OPM smoothness in the V1-like layer. Dashed green: value in macaque V1. Dashed gray: smoothness in an Unoptimized model. Bottom: Category selectivity map smoothness in the VTC-like layer. Dashed lines indicate means across human subjects and hemispheres from the NSD data; one line per category. (d) Density of topographic phenomena of interest as a function of  $\alpha$ . Top: Pinwheel density in OPMs from the V1-like layer, as a function of  $\alpha$ . Bottom: Number of category selective patches for each category in the VTC-like layer, as a function of  $\alpha$ . Human data in dashed lines.

# Two key factors underlying functional organization: self-supervised learning and a scalable spatial constraint

Having established that specific TDANN models accurately predict the functional organization of the ventral visual
 stream, we consider what key factors enable the emergence of this functional organization. We reasoned that if
 some combinations of optimization objectives yield brain-like functional organization and others do not, it will shed
 light on the constraints underlying the observed functional organization. Thus, we train models with alternative task
 and spatial objectives, then apply our benchmarks to evaluate which models are most consistent with empirical data.

For the "task component" of its loss function, the TDANN uses contrastive self-supervision [61, 63], a framework for learning representations that transfer easily to many downstream tasks. These self-supervised algorithms have been shown to generalize to many downstream computer vision tasks despite being trained only on a large set of unlabeled



Figure 5. Self-supervision and scalable spatial constraints underly the emergence of functional organization. In each panel, TDANN shown in purple, Categorization-trained in gold, Absolute SL in red, and ventral stream measurements in green. (a) Left: comparison of task objectives. The TDANN uses contrastive self-supervision (top) which encourages similarity between representations of different views of the same image while increasing distance between representations of views of other images. Categorization (bottom) compares predicted class probabilities to the human-labeled correct class. Right: comparison of spatial objectives.  $S_{ij}$ : response similarity of units *i* and *j*.  $d_{ij}$ : cortical distance between units *i* and *j*. TDANN uses the Relative SL (top), which correlates the population of response similarities and pairwise inverse distances. Prior work [78] used the Absolute SL (bottom), which directly subtracts inverse cortical distance from response similarity magnitude. (b) Smoothed orientation preference maps (OPMs) in the V1-like layer of the TDANN (left), a Categorization trained model (middle), and a model trained with the Absolute SL (right). Dots: detected pinwheels.  $\alpha = 0.25$  for models shown in each panel. (c) Category selective units in the VTC-like layer of the TDANN (left), a categorization trained model (middle) and a model trained with the absolute SL (right). (d) Right: Smoothness of OPMS in the V1-like layer of each model type. Green line: value computed macaque V1. (e) Density of detected pinwheels. Green: estimated value in macaque V1. (f) Right: Smoothness of face selectivity maps in the VTC-like layer of each model type. Green line: value from human VTC. (g) Average number of category-selective patches, in the VTC-like layer in each model. Green: average value in human VTC.

natural images [80]. However, most studies comparing neural networks to the brain have used a supervised object
 categorization ([26, 25, 78]; Figure 5a-bottom left). Thus, we tested whether training with an object categorization
 objective produces different functional organization than self-supervision, and if so, which is more similar to the
 observed functional organization of the ventral visual stream.

We also investigate how the form of the spatial objective function affects emergent functional organization. The 298 spatial component of the TDANN loss function is generally intended to capture the constraints on unit-to-unit 299 correlations within cortical neighborhoods, but the specifics of its functional form embody conceptually distinct 300 mechanistic ideas about how a hypothetical cortical development circuit might measure functional correlations and 301 compare them to cortical distances. In prior work, Lee et al. [78] introduced a spatial loss function that subtracts 302 the inverse of pairwise cortical distances from the magnitude of pairwise response correlations (Figure 5a-bottom 303 right), such that nearby units develop similar responses. That loss function was developed to match empirical 304 measurements in macaque IT, but was not intended to generalize to other regions of the human ventral visual 305 stream. We refer to it as the Absolute Spatial Loss (or SLAbs), because minimizing it requires an absolute match 306 between response correlations and the inverse of cortical distances. While Lee et al. [78] found that training models 307 with SLAbs produced clustering of category-selective units in a late model layer, we discovered a critical flaw when 308 training with SLAbs in all model layers: in layers with shorter cortical distances, SLAbs can only be minimized if 309 response correlations are pathologically high. The TDANN instead uses a more flexible spatial loss function that we 310 term the Relative Spatial Loss ( $SL_{Rel}$ ; Figure 5a-top right). This SL requires that inverse cortical distances will be 31 correlated with response similarity (see Methods for mathematical details). SL<sub>Rel</sub> effectively enforces response 312 similarity between pairs of units that are relatively close together. Thus, the Relative SL allows the distance 313

over which local correlations extend to depend on the total size of the cortical area. Interestingly, we find that switching from  $SL_{Abs}$  to  $SL_{Rel}$  slightly increased the model's capacity for object categorization at all levels of  $\alpha$ (Supplementary Figure S12). How do models trained for different objectives differ on topographic benchmarks?

We compare the TDANN (self-supervised and Relative SL) to categorization-trained models (differing only in task 317 objective) and Absolute SL models (differing only in spatial objective) on our battery of topographic and functional 318 benchmarks: (i) evaluating the smoothness of OPMs and face-selectivity maps in the V1-like and VTC-like layers, 319 respectively, and (ii) counting the number of pinwheel-like discontinuities and category-selective patches in those 320 layers, respectively. Categorization-trained models were slightly but significantly less smooth than the TDANN (mean 321 smoothness = 0.56, U = 25, p = 0.008), but with an equal density of pinwheels (2.07 pinwheels / column spacing<sup>2</sup>; 322 U = 10, p = 0.69). Absolute SL models generally resemble those in the TDANN (Figure 5b), but with significantly 323 lower smoothness (TDANN mean: 0.71, Absolute SL: 0.40; U = 25, p = 0.008; Figure 5d) and slightly lower pinwheel 324 density (TDANN: 2.14 pinwheels / column spacing  $^2$ , Absolute SL: 0.89; U = 21, p = 0.09; Figure 5e). 325

Strikingly, however, category-selectivity maps in the VTC-like layer were much less organized in the 326 Categorization-trained models than in the self-supervised TDANNs. At the same spatial weight of  $\alpha = 0.25$ , 327 clear clusters of category-selective units are observed in the self-supervised but not the categorization-trained 328 model (Figure 5c). The Absolute SL models also fail to form organized category-selectivity maps at this level 329 of  $\alpha$ . Quantitative comparison reveals smoother category selectivity maps in the TDANN (mean smoothness of 330 face-selectivity maps = 0.44) than in either categorization-trained models (0.09; Mann-Whitney U = 25, p = 0.008; 331 Figure 5f) or in Absolute SL models (0.13). The TDANN also has a significantly higher number of identified category 332 selective patches (mean = 1.2) than either categorization-trained (mean = 0) or Absolute SL alternatives (mean 333 = 0.08; U = 25, p = 0.008; Figure 5g). Thus, the nature of the training objective strongly constrains the emergent 334 functional organization, with self-supervised learning and relative spatial loss objectives producing the most brain-like 335 functional organization. 336

#### 337 Spatial constraints make learned representations more brain-like by reducing intrinsic dimensionality

A natural question is whether training for spatial objectives also has an effect on the *non-topographic* properties of learned representations. Because the TDANN allows the network's features to be influenced by the spatial constraint during training, we can directly address this guestion.

A powerful way to test if spatially-constrained models learn different features than standard DANNs is to measure 341 how well model unit responses can predict neural responses to large set of naturalistic images in primate visual 342 cortex [30, 28, 25, 81]. A popular approach to predicting neuronal firing rates is to fit the responses of individual 343 neural units with a linear combination of many hundreds or thousands of model units. Consistent with prior work 344 involving non-spatial models [61], we find that models trained with different objectives are largely indistinguishable 345 in their ability to predict neural firing rates when using this standard linear-regression method for mapping model 346 units to neural firing rates [25, 61, 29] (Figure 6a). The linear-regression mapping is thus insensitive to the dramatic 347 differences between models trained with different objectives and spatial constraint magnitudes that are apparent in 348 our analysis of functional organization. A possible explanation for this apparent discrepancy is that linear regression 349 is too permissive of mapping: even if a model lacks individual units that resemble recorded neurons, a combination 350 of units might still allow for accurate prediction of neural responses. We tested this prediction by performing a more 351 stringent one-to-one mapping, in which individual VTC-like layer model units - not a linear mixture of units - are 352 assigned to individual VTC voxels in a one-to-one fashion. Intriguingly, we found that this one-to-one assignment 353 resulted in much stronger matches between TDANN model units and voxels recorded in the Natural Scenes Dataset 354 (NSD) [75] than models trained with other objectives (i.e. categorization or Absolute SL, Figure 6b). This correlation 355 peaks at  $\alpha = 0.25$ , the same value identified by topographic benchmarks (Figure 4), providing more evidence that 356 the constraints driving brain-like functional organization also make learned representations more brain-like. 357

Many factors might contribute to the differences in representation between the TDANN and those of poorer-fitting 358 models. Because the TDANN's spatial constraint encourages units to respond more similarly to one another, we 359 hypothesized that the intrinsic dimensionality of the population might decrease as  $\alpha$  increases. Relatedly, recent 360 work has demonstrated that spatially unconstrained DANN responses to natural images have substantially higher 36 intrinsic dimension than real macaque and rodent V1, and that models with lower dimensionality better predict 362 neural responses [82]. Thus, we tested whether decreased intrinsic dimensionality might explain why the TDANN 363 representations are more brain-like than representations from other models. Consistent with our hypothesis, we 364 find that the addition of the spatial constraint decreases intrinsic dimensionality in the VTC-like layer regardless of 365 the training objective (Figure 6c; see Supplementary Figure S13a for eigenspectra in all layers). When  $\alpha = 0$ , all 366 models have higher effective dimensionality (ED; Elmoznino and Bonner [83], Del Giudice [84]; see methods) than 367 human VTC (mean across subjects = 16.7), although the dimensionality of the VTC-like layer in categorization-trained 368

models (76.8) is nearly three times higher than in the self-supervised models (TDANN and Absolute SL: 27.8). At 369 the spatial weight magnitude  $\alpha = 0.25$ , at which the TDANN best matches neural data, the TDANN's VTC-like layer 370 approaches the dimensionality of human VTC (TDANN mean = 13.2). However, the dimensionality of models trained 37 with  $SL_{Abs}$  decreases too guickly (mean = 6.5), and categorization-trained models remain higher than human VTC 372 at this level of  $\alpha$  (mean = 42.7). 373

We conclude that the close match between the TDANN and human VTC, on both topographic and non-topographic benchmarks, may be due in part to an alignment of their intrinsic dimensionality. Similar results are observed when summarizing the response eigenspectrum with power law fits, as in Stringer et al. [85], Kong et al. [82] (Supplementary Figure S13c). Intriguingly, we find that the effective dimensionality of the TDANN roughly converges to a common value of approximately 15 across model layers at  $\alpha = 0.25$  (Figure 6d), raising the possibility that a similar dimension stabilization phenomenon occurs across brain areas in the ventral stream. These results provide new evidence that the computational constraints generating cortical topography strongly influence non-topographic

features, making them more brain-like by virtue of decreasing the dimensionality of population responses. 381



Figure 6. Spatial constraints make learned representations more brain-like and reduce intrinsic dimensionality (a) Variance explained under a linear regression mapping between model units and macaque IT neurons, as a function of the spatial loss weight  $\alpha$  and the training objective. (b) Mean correlation between model units and VTC voxels under a one-to-one mapping as a function of  $\alpha$ . Green: mean human-to-human correlation under the same one-to-one mapping. (c) Estimated effective dimensionality (cf. Elmoznino and Bonner [83], Del Giudice [84]) of the population response in the VTC-like layer of models trained at different levels of  $\alpha$  and with different objectives. Green: mean value in human VTC from the NSD dataset. (d) Effective dimensionality in the TDANN across all layers and levels of  $\alpha$ . In all panels, shaded vertical bar indicates value of  $\alpha$ demonstrated in prior analyses to best match topographic phenomena.

#### The TDANN minimizes inter-layer wiring length 382

374

375

376

377

378

379

380

Identifying the optimization paradigm that is most consistent with neural data provides insight into the constraints 383 underlying neural development, but prompts a deeper question: why would these constraints be favored by 384 evolutionary selection? A natural hypothesis is that cortical networks with strong functional organization also 385 minimize wiring length, and thus reduce brain size, weight, and power consumption [86, 18]. We test this hypothesis 386 by asking whether the optimization paradigm that generated a functional organization that best fit neural benchmarks 387 - intermediate spatial weight  $\alpha$ , self-supervised learning, and spatial costs that scale with cortical surface area -388 also reduces between-layer wiring length. In feedforward networks that lack intra-layer connectivity, such as the 389 TDANN, any gains in wiring efficiency must be between layers. Accordingly, we measure inter-layer wiring length by 390 identifying populations of co-activated units in adjacent layers, then estimating the length of fibers needed to connect 391 those populations. We first present natural images to the network and record the locations of the most responsive 392 units in each layer, then simulate fiber bundles that originate in an earlier "source" layer and terminate in the following 393 "target" layer, adding inter-layer fibers until the total squared distance between each activated unit and its nearest 394 fiber is below a specified threshold (see Methods, Figure 7a). The total wiring length is taken as the sum of the 395 lengths of each fiber. 396

Presenting the TDANN with natural images leads to clustered responses in the VTC-like layer of all models trained 397 with  $\alpha > 0$ , with multiple clusters apparent at higher levels of  $\alpha$  (Supplementary Figure S14). Does the increase in 398 clustering within layers result in shorter wiring length between layers? We find that inter-layer wiring length is indeed 399 minimized at higher levels of  $\alpha$  (Figure 7b). However, we also find that object categorization performance decreases 400 as wiring efficiency improves (Figure 7c), indicating that models at low-to-intermediate levels of  $\alpha$  optimally balance 401 performance with inter-layer wiring efficiency. This coincidence of optimal  $\alpha$  values suggests that the functional 402 organization of the ventral visual stream balances inter-area wiring costs with performance. Critically, we find 403 that wiring is most efficient for the optimization objectives that yield the most brain-like functional organization: 404 wiring length is higher in both categorization-trained models and those trained with the Absolute SL (Figure 7d). 405



Figure 7. Minimization of inter-layer (feedforward) wiring length in models with brain-like functional organization. (a) Example wiring length computation between adjacent layers. Units in brown are the top 5% most active units in the Source layer for an arbitrarily-selected natural image, while units in green are the top 5% most active in the Target layer. Black dots show the origination and termination points of fibers that would be required to connect populations of active units across layers. (b) Wiring length between layers 4 and 5 ("V1"; left), and layer 8 and 9 ("VTC", right) as a function of  $\alpha$ . Shaded regions: 95% CI of measurements from different cortical neighborhoods, model seeds, and input images. (c) Accuracy on object categorization vs total wiring length, for models trained at different levels of  $\alpha$ . (d) Wiring length in both early and later model layers for models trained with different task and spatial objectives ( $\alpha = 0.25$  for all). Error bar: 95% CI over different image presentations and model seeds.

Thus, wiring length minimization provides a normative explanation for the superiority of self-supervised learning and area-normalized spatial constraints.

#### <sup>408</sup> Proof-of-principle: Using the TDANN as a digital twin for experimental design

A quantitatively accurate and mechanistically grounded model of functional organization, such as the TDANN,
 enables a spectrum of applied use cases that rely on estimating the effects of spatially-modulated neural
 perturbations. Here we apply TDANN as a digital twin of visual cortex and demonstrate two novel applications:
 1) performing an *in silico* microstimulation experiment, and 2) proof-of-principle for prototyping a simple cortical
 prosthetic device.

Simulated microstimulation reveals functional similarity of conneted unit populations Microstimulation experiments 414 in the macaque [87] found that stimulating neurons in a face patch selectively drives activity in other face patches, 415 and prior work with topographic models of macaque IT [78] found a similar result. We tested if the TDANN also 416 captures this connectivity by stimulating local populations of units in the penultimate model layer and recording 417 evoked responses in the following VTC-like layer. Mirroring results in macaque IT, we find that stimulating units 418 in a TDANN face patch drives localized activity in a face patch in the following layer (Figure 8a). We repeat the 419 stimulation for 99 other sites equally spaced on the simulated cortex, and find that the selectivity of a stimulated unit 420 in the source layer strongly predicts the selectivity of activated units in the target layer (Figure 8b), especially for 421 stimulation sites closer to the center of the simulated cortical tissue. 422

Simulation of cortical prosthetic devices with TDANNs A unique advantage of a unified topographic model such as TDANN is that it can be used to prototype the effects of simultaneous stimulation of multiple cortical areas, experiments which are challenging to perform *in vivo*. Based on recent advances in machine learning and visual cortical prostheses [88, 89, 90], we introduce a framework using TDANNs to prototype multi-region cortical stimulation devices. The framework has two components (Figure 8c, d): 1) a Stimulation Simulator that transforms desired activity patterns on the cortical sheet into *device-achievable* patterns, and 2) a Percept Synthesizer that estimates the percept evoked by stimulation with those patterns.

The Stimulation Simulator takes an input image, uses the TDANN to predict the precise pattern of responses in each layer, and then constrains that pattern into one that is physically achievable by a specific hypothetical stimulation device (Figure 8c). We model two kinds of physical constraints: spatial precision – the resolution at which the device can create activity patterns, and regional access – the subset of cortical areas that are accessible to the device. Spatial precision is modeled as a Gaussian blur of the desired activity pattern and regional access by restricting the model layers that participate in the simulation.

To synthesize percepts from device achievable patterns, we use an approach inspired by Granley et al. [90] and Shahbazi et al. [91] to synthesize the input image which generates the target activity pattern – i.e., a neural metamer. Figure 8e illustrates predicted percepts for hypothetical cortical stimulation devices with variable precision and access. Unsurprisingly, a device with infinitely high spatial stimulation precision yields sharp percepts even when only early cortical areas are stimulated (Figure 8e, top left). However, the percepts quickly deteriorate as the spatial precision of the device decreases (Figure 8e lower left). Notably, our simulation suggests that, at lower spatial

442 precision, the quality of percepts can be improved by adding stimulation of higher cortical areas (Figure 8e, middle 443 rows).

While we have neglected many critical details here, including spatiotemporal processing, cortical magnification, and the need to validate percepts, we hope that this proof of principle motivates the use of TDANN to make testable predictions about the nature of percepts elicited by various cortical stimulation devices.



Figure 8. Using TDANNs to simulate spatial stimulation devices. (a) Stimulation of a local population of units in the second to last convolutional layer drives spatially-localized responses in the final convolutional layer. Responses are functionally aligned, such that stimulating face-selective units (Site 1) drives activity in face-selective units in the following layer. Right: Results for a second stimulation site, at the intersection of place-, body-, and character-selective patches. (b) Similarity in tuning of stimulated units in the source layer and responding units in the target layer for 100 evenly-spaced stimulation sites. Each dot compares tuning similarity for the true distribution of activated units (x-axis) and a randomly shuffled selection of units (y-axis). Dot color: distance of the stimulation site from the center of the cortical tissue. (c-d) Conceptual framework for applying the TDANN to the prototyping of visual cortical prostheses. (c) Stimulation Simulator: the TDANN is used to generate predicted activity patterns from a given visual input (top row). Patterns are then degraded according to the limitations on a hypothetical stimulation device: reduced spatial precision results in blurring of the target activity pattern (bottom row), and limits to regional access restrict the set of layers that participate. Here, Layer 8 is faded-out to show that this particular hypothetical device cannot reach that cortical area. (d) Given a device-achievable stimulation pattern produced by the Stimulation Simulator in (c), we synthesize the image that could evoke that pattern: the predicted percept. To build intuition for the fidelity of predicted percepts, we use an example input image of the the first four lines of a Snellen eye chart. (e) Predicted percepts for 25 theoretical cortical stimulation devices with different capabilities. Devices vary in the precision with which they are able to produce desired activity patterns (full-width at half-maximum (FWHM) of the spread of activity on cortex increases with rows) and the number of cortical areas that can be simultaneously simulated (columns).

# 447 Discussion

In this work, we leveraged the neural network modeling framework to seek the principles of functional organization
 in the primate ventral visual stream. We found that training a spatially-augmented deep neural network for a specific
 combination of objectives results in a model, the TDANN, that captures topographic properties throughout the ventral
 stream, from the pinwheels of V1 to the category-selective patches of higher-level visual cortex.

We identified two specific factors critical to the emergence of brain-like functional organization. First, we found 452 that self-supervised learning of task-general representations yields better organization than the more common 453 alternative of supervising on the singular task of visual object recognition. Recent work has suggested that functional 454 specialization in the brain - e.g., one population of units responsible for discrimination of different faces and another 455 for recognition of different objects - arises under joint training for two different supervised recognition tasks, one for 456 faces and one for objects [92]. Our results demonstrate that functional specialization can emerge under a single 457 unsupervised learning objective on a single training set, suggesting that general mechanisms can produce the 458 kinds of functional specialization that is typically assumed to require multiple objectives or multiple distinct datasets. 459 Second, we found that the spatial constraint in our model should compare response similarity and physical similarity 460 according to a metric that scales with the size of each cortical area, rather than being fixed for all cortical areas. 461 This finding suggests that the actual circuits responsible for shaping the structure of local response correlation in 462 cortical neighborhoods should scale with the surface area of each cortical region. Our identification of these two 463 critical factors demonstrates that a goal-driven modeling approach to understanding neural sensory systems can 464 vield concrete and specific insights into their underlying principles. 465

Critically, the two factors that we found are essential for brain-like functional organization in the visual system are 466 not specific to the visual modality, and might extend to predict the abundant, yet largely unexplained, functional 467 organization in other sensory systems. For example, neurons in primary auditory cortex are arranged according 468 to the frequency they respond most strongly to (tonotopy [2]), and in secondary auditory areas, neurons cluster 469 according to their preference for speech and music [23, 93]. It is possible that the representations carried by these 470 neurons are also learned by contrastive self-supervision, and that their topographic organization is explained by 471 scalable spatial constraints of the forms described here. Likewise, the functional organization of somatosensory [4], 472 entorhinal [6, 5] and parietal cortices [3] may be explained by the specific yet general principles for representation 473 learning and spatial smoothness that we have identified. Under this hypothesis, it is only the structure of the input 474 data (e.g., auditory experience, somatosensory input) that changes, but the cortical mechanisms for learning and 475 organization remain universal across cortical systems. Future work can directly test that hypothesis by training 476 TDANN variants to learn spatially-organized representations specific to each system. 477

The TDANN is the first model to predict functional organization in multiple cortical areas by learning features and 478 topography, from scratch, in and end-to-end optimization framework trained directly on image inputs. As such, it 479 represents an improvement over a number of related prior approaches. For example, hand-crafted self-organizing 480 maps (SOMs) [94, 8, 10, 11, 9] have simplified the problem of topographic map formation by modeling a limited set 481 of fixed feature dimensions (e.g., orientation preference and spatial frequency tuning), then modifying the tuning of 482 model units along these dimensions such that nearby units develop similar selectivity. While such SOMs produce 483 qualitatively smooth V1-like orientation maps, we find that they fail to quantitatively predict the topographic properties 484 485 of V1 orientation maps (Figure 2). Recent attempts to abandon hand-crafted feature dimensions have trained SOMs to smoothly map the outputs of categorization-pretrained DCNNs [12, 13]. While these DNN-SOMs have 486 the advantage of operating on images rather than predefined features, we find that they are quantitatively less 487 accurate than the TDANN (Figure 3) at explaining the functional organization of VTC, and fail to reproduce the 488 topography of V1 (Figure 2). Another recent approach, the ITN [20], appended topographic layers to a pretrained 489 DCNN backbone and trained for supervised categorization under an additional wiring length minimization constraint. 490 While the ITN reproduces many features of VTC topography, it does not predict the size, number, and geometry 491 of category-selective patches as accurately as the TDANN, and cannot predict the functional organization of areas 492 outside VTC. Prior work from our groups also followed the TDANN optimization framework, but used a supervised 493 categorization task, a spatial constraint that did not scale with cortical area, and applied only to the VTC-like layer 494 [78]. While this model was able to predict many properties of the functional organization of macaque IT, it incapable of 495 predicting the organization of other ventral stream regions. Our present results (Figure 5) demonstrate that different 496 spatial and task objectives are required for a TDANN to accurately match the functional organization of multiple areas 497 of the ventral visual stream. 498

That the TDANN is trained end-to-end provides two interesting opportunities for understanding the interaction between learned representations and functional organization during development. First, our preliminary analyses suggest that trajectories of TDANN functional architecture throughout training roughly match the faster development of earlier vs higher cortical regions (Figure S17) and the emergence of V1-like topography from retinal wave-like stimuli (Figure S18). Rigorously testing those predictions would be most interesting when the TDANN is optimized using naturalistic movie streams that match the visual statistics and acuity limitations of human development [95, 96]. Second, we found that the presence of the spatial constraint during training modulated the nature of learned representations, making them more brain-like and stabilizing their intrinsic dimensionality (Figure 6).

While the TDANN is the first unified model of ventral stream functional organization, it has a number of important 507 limitations. Because the core DCNN architecture used in this work is strictly feedforward, there are no direct 508 connections between different units in the same laver. Thus, we are only able to draw inferences about how the 509 spatial constraint affects wiring length between layers. A more complex architecture could include both intra-layer 510 recurrence and long-range feedback connections [97], although our results demonstrate that explicitly modeling 511 these recurrent connections is not necessary to produce accurate topographic maps (see Figure 6 of Blauch et al. 512 [20]), raising the possibility that minimization of the length of long-range fibers may be the key determinant of the 513 functional organization of visual cortex. 514

We also note that our model, like all convolutional neural networks, uses the same filter weights across the entire 515 visual field (termed "weight sharing"). This short-cut makes large-scale network training feasible; however, it is 516 biologically implausible and potentially interferes with topographic map formation, since changing input weights to a 517 unit in one part of the cortical sheet will also change the weights of many other distant units in a non-local fashion. 518 Some topographic models avoid this issue by forgoing the use of convolutional layers altogether, but in doing so 519 forfeit the ability to model retinotopically-organized cortical areas. In contrast, our approach is to pre-optimize unit 520 positions (see Methods) in a way that allows the learning of locally-smooth topographic maps even with convolutional 521 layers (see Methods). In the brain, a similar pre-optimization may be achieved by chemical gradients [98] and 522 experience-independent refinement of neural circuits during embryonic development[99, 100, 101, 102]. 523

<sup>524</sup> Finally, an exciting application of the TDANN is the simulation of experiments with spatial manipulations and readouts <sup>525</sup> (Figure 8). Virtually every experiment that uses topographic structure as a dependent variable, including controlled

rearing and task learning paradigms, could first prototype experiments with TDANNs. In addition, experiments that involve inactivation or stimulation of local populations of neurons (e.g. Rajalingham and DiCarlo [103], Shahbazi et al. [91]) could use the TDANN to predict the downstream behavioral impact of those manipulations prior to collecting data. The tools to perform stimulation or inactivation of neural populations have become commonplace in systems neuroscience in the past decade, but their engagement with the strongest models of neuronal function – task-optimized neural networks – has been limited due to the lack of image-computable models that not only explain the responses of individual neurons [104, 105, 106, 25] but that are also mapped to cortical tissue. As a unified

model of functional organization, the TDANN is well-suited to bridge this gap.

## 534 Methods

## 535 Code and data availability

<sup>536</sup> Code for model training and analyses is available at https://github.com/neuroailab/TDANN.

### 537 Neural network architecture and training

Model training. We build off of the torchvision implementation of ResNet-18 [59] and train models with modifications 538 to the VISSL framework [107]. All models were trained for 200 epochs of the ILSVRC-2012 (ImageNet Large-Scale 539 Visual Recognition Challenge; Deng et al. [64]) training set. Unless otherwise indicated, models were each trained 540 from five different random initial seeds. Network parameters were optimized with stochastic gradient descent with 541 momentum ( $\gamma = 0.9$ ), a batch size of 512, and a learning rate initialized to 0.6 then decaying according to a 542 cosine learning schedule [108]. Models were trained either for supervised 1000-way object categorization or on the 543 self-supervised contrastive objective "SimCLR" [63]. Following training, categorization accuracy for self-supervised 544 models was assessed by freezing the parameters of the model and training a linear readout from the outputs of the 545 final layer. The linear readout is trained for 28 epochs with a batch size of 1,024 and a learning rate initialized to 0.04 546 and decreasing by a factor of 10 every eight epochs. 547

Initialization of model unit positions. Prior to training, model units in each layer are assigned fixed positions in a 548 two-dimensional cortical sheet that is specific to that layer. For efficiency, we do not embed the units of the very 549 first convolutional layer. The size of the cortical sheet in each layer depends on a mapping between model layers 550 and regions in the human ventral visual pathway, as well as a commitment to the extent of the visual field being 551 modeled. For example, because we map model Layer 4 to human V1, the surface area of the cortical sheet in that 552 layer is set to  $13cm^2$ : the mean value reported by Benson et al. [109] for the surface area of the section of human 553 V1 that is sensitive to the central 7 degrees of visual angle. Another critical parameter in our framework is the size 554 of a "cortical neighborhood": during training, computation of the spatial loss is restricted to units within the same 555 cortical neighborhood. We set the neighborhood width to match measurements made of the spatial extent of lateral 556 connections in different cortical areas of the macaque (from Yoshioka et al. [110]), then scale up to achieve estimates 557 that might match the human ventral visual pathway. Table 1 details the sizes of simulated cortical sheets and cortical 558 neighborhoods in all layers. 559

Layer	# Units	Size of Cortical sheet	Neighborhood Size	Region
Layer 2	200704	$5.7mm^{2}$	$47 \mu m$	Retina
Layer 3	200704	$5.7mm^{2}$	$47 \mu m$	Retina
Layer 4	100352	$13.5 cm^2$	$1.6mm^{*}$	V1
Layer 5	100352	$13.5 cm^2$	$1.6mm^*$	V1
Layer 6	50176	$12cm^2$	4mm	V2
Layer 7	50176	$5cm^2$	2.5mm	V4
Layer 8	25088	$49cm^2$	31mm	VTC
Layer 9	25088	$49cm^2$	31mm	VTC

**Table 1.** Parameters for layer positions. \*the value of 1.6mm used in the V1-like layer is known to be inaccurate, but matching the proper value yields too few units in each cortical neighborhood to compute pairwise distances. See Supplementary Figure S5 for a solution to this problem.

<sup>560</sup> Positions are assigned in a two-stage process:

Stage 1: Naive Retinotopic Initialization Because each layer performs a convolution over the previous layer's outputs, responses are organized into spatial grids. We preserve this intrinsic organization by assigning each model unit to a region of the simulated cortical sheet that corresponds to its spatial receptive field.

Stage 2: Pre-optimization of positions Convolutional networks share filter weights between units at different 564 locations; thus, local updates to a single unit entail updates to all units with the same filter weights. It is highly 565 unlikely that an arbitrary configuration of unit positions will permit local smoothness under this global coordination 566 constraint. Thus, we perform pre-optimization of unit positions to identify a set of unit positions for which learning 567 smooth cortical maps is possible. Specifically, we spatially shuffle the units of a pre-trained DCNN on the cortical 568 sheet such that nearby units have correlated responses to a set of sine grating images. The choice of sine gratings 569 here is inspired by observations that edge-like propagating retinal waves drive experience-independent organization 570 of the visual system in primates and other mammals [99, 100, 101, 102]. 571

The spatial shuffling works as follows: 1) Select a cortical neighborhood at random. 2) Compute the pairwise response correlations of all units in the neighborhood. 3) Choose a random pair of units, and swap their locations in the cortical sheet. 4) If swapping positions decreases local correlations (measured as an increase in the Spatial Loss function described below), undo the swap. 5) Repeat steps 3-4 500 times. 6) Repeat steps 1-5 10,000 times.

*Loss functions.* We use two kinds of loss functions: spatial losses that encourage topographic structure, and task losses that encourage the learning of visual representations. We detail each in turn below:

**Spatial loss** The spatial loss (SL) function encourages nearby pairs of units to have response profiles that are more correlated with one another than those of distant of units. Consider a neighborhood with N units. The vector of pairwise Pearson's response correlations,  $\vec{r}$ , has length  $M = \binom{N}{2}$ , the number of unique pairs. Let the corresponding

vector of pairwise Euclidean cortical distances be denoted  $\vec{d}$ .

582 We define two SL variants:

$$SL_{Abs} = \frac{1}{M} \sum_{i=1}^{M} |r_i - D_i|,$$
 (2)

$$SL_{\rm Rel} = 1 - Corr(\vec{r}, \vec{D}),$$
 (3)

where Corr is the Pearson's correlation function and  $\vec{D}$  is the inverse distance:

$$D_i = \frac{1}{d_i + 1} \tag{4}$$

Task loss The task loss is computed from the output of the final model layer. We use two task losses: the
 object categorization cross-entropy loss used in supervised object recognition (e.g. Krizhevsky et al. [111]) and
 the self-supervised SimCLR objective [63].

Combination of losses during training
 On each batch, model weights are updated to minimize a weighted sum of
 the task loss and the spatial loss contributed by each layer:

$$TDANN Loss = L_{task} + \sum_{k \in layers} \alpha_k SL_k$$
(5)

where  $\alpha$  is the weight of the spatial loss.

<sup>590</sup> Overview of Training In summary, models are trained in 6 steps:

- <sup>591</sup> 1. ResNet-18 is trained on the task loss only.
- <sup>592</sup> 2. Positions in each layer are initialized to preserve coarse retinotopy (Stage 1).
- Positions are further pre-optimized in an iterative process that preserves retinotopy while bringing together
   units with correlated responses to sine gratings images (Stage 2).
- <sup>595</sup> 4. Positions are frozen and never again modified.
- 596 5. All network weights are randomly re-initialized.
- <sup>597</sup> 6. The network is trained to minimize a weighted combination of the spatial and task loss components.

### 598 Benchmarks comparing macaque V1 to model V1-like layers

#### 599 Stimuli and Tuning Curves.

Sine Grating Images Tuning to low-level image properties such as orientation, spatial frequency, and chromaticity was assessed by constructing  $224 \times 224$  pixel sine grating images that span 8 orientations evenly spaced between 0 and 180 degrees, 8 spatial frequencies between 0.5 and 12 cycles per degree, 5 spatial phases, and two chromaticities: black/white gratings and red/cyan gratings.

**Tuning Curves** We evaluated tuning for orientations and spatial frequencies by constructing tuning curves for each unit. Color-responsiveness is assessed by comparing the mean response to all black and white gratings to the mean response to all red/cyan gratings. The distribution of model unit activations for a given layer was rescaled to match the minimum and maximum firing rates reported in [34]. We quantify the orientation tuning strength of model units using circular variance (CV), where values closer to 0 correspond to sharper tuning. As in Ringach et al. [34], CV is defined as:

$$CV = 1 - \left| \frac{\sum_{k} r_k e^{i2\theta_k}}{\sum_{k} r_k} \right|$$
(6)

<sup>610</sup> Where  $\theta_k$  is the *k*th orientation, in radians, and  $r_k$  is the scaled response to that orientation. Orientation tuning <sup>611</sup> curves are additionally fit with a von Mises function whose peak is taken as the preferred orientation.

#### 612 Models.

Hand-Crafted Self-Organizing Map Our hand-crafted self-organizing map (SOM) implementation uses the *MiniSom* library [112], with parameters adapted from Swindale and Bauer [11]. We instantiate the SOM as a 128 x 128 grid
 of model units.

<sup>616</sup> 10,000 training samples were randomly constructed by selecting a random (x, y) location, orientation ([0,  $\pi$ ], spatial <sup>617</sup> frequency ([0, 1]), and chromaticity (black/white, colorful).

As in Swindale and Bauer [11], SOM weights were initialized retinotopically with randomly-selected initial preferred orientations.

<sup>620</sup> The SOM is trained by presenting training examples for a total of 700,000 updates. After each example, the "winning"

unit (i.e. the one with the highest response) is updated with a learning rate of  $\epsilon=0.02$  to be more strongly aligned

with the input stimulus, and its neighbors are updated in proportion to their proximity to the winner, as determined by a Gaussian neighborhood function parameterized by  $\sigma = 2.5$ .

Following training, each sine grating in the set of probe stimuli is presented to the SOM by projecting it into the six-dimensional space of SOM unit tuning and computing the response of each SOM unit to the stimulus. Once responses to each stimulus are obtained, tuning curves are constructed as usual.

**DNN-SOM** The DNN-SOM is identical to the hand-crafted SOM, except that 1) the inputs are derived from the outputs of the first layer of an AlexNet model pretrained for ImageNet object categorization and 2) the learning rate is increased, which we found helps convergence. Following the approach of Zhang et al. [12], we take the responses of the first AlexNet layer to all 50,000 natural images in the ImageNet dataset, reduce their dimensionality with principal components analysis, and train the SOM on those examples.

Response Benchmarks. Model responses are compared to macague V1 by considering preferred orientations and 632 orientation tuning strength. Orientation tuning strength is computed as circular variance (CV) and compared 633 between the population of model units and the empirical distribution provided by Ringach et al. [34] with the 634 Kolmogorov-Smirnov distance. To filter out noisy units, we compute CV for model units with a mean response 635 magnitude of at least 1.0. The distribution of preferred orientations is also compared to empirical data collected by 636 De Valois et al. [35] by counting the number of units preferring each of four orientations: 0, 45, 90, and 135 degrees. 637 In Figure S3b we compute a "Cardinality Index": the fraction of preferred orientations that include, 0, 90, and 180 638 degrees. 639

Topographic Benchmarks. Orientation preference maps (OPMs) are compared to empirical measurements in two
 ways: counting pinwheels and quantifying map smoothness.

We interpolate the OPM onto a two-dimensional grid by computing the circular mean of the **Pinwheel Detection** 642 preferred orientation of units near a given location. If the population of model units near a grid location has 643 high heterogeneity in preferred orientation, we disqualify that pixel for having an unreliable estimate of preferred 644 orientation. Each grid location is assigned a "winding number" [17], computed by considering the preferred 645 orientations of the eight pixels directly bordering the pixel under consideration. Moving clockwise around the 646 bordering eight pixels, the change in preferred orientation from pixel to pixel is summed. A high winding number 647 indicates a clockwise pinwheel, and a low winding number indicates a counterclockwise pinwheel, where the 648 thresholds for "high" and "low" are selected to be consistent with manual annotation of clear pinwheels. 649

Pairwise Tuning Difference We compute the smoothness of orientation preference maps by constructing a curve 650 relating pairwise difference in preferred orientation to pairwise cortical distance. First, we restrict the population of model units to those with the highest 25% peak-to-peak tuning curve magnitudes. This filtering step removes units 652 with weak responses or responses that would be indistinguishable from a "cocktail blank" background activity level. 653 and we consider it equivalent to neuron selection in electrophysiological and optical imaging studies [34, 43]. As in 654 similar approaches to quantifying OPM structure (e.g. Chang et al. [68]), pairs of units are binned according to their 655 distance, and the average absolute different in preferred orientation is plotted for each distance bin. Because there 656 can be hundreds of thousands of units in a given layer, we restrict this analysis to randomly-selected neighborhoods 657 of a fixed width, then sample many neighborhoods from each map. Finally, we divide the pairwise distance by the 658 chance value obtained by random resampling of unit pairs, such that a values < 1 indicate more similar tuning than 659 would be expected by chance. 660

The OPM curves are compared to reconstructed macague V1 data from Nauhaus et al. [43]. 661

We adopt an identical approach for the construction of a neural spatial frequency preference map, where data are 662 also provided for the same imaging window in Nauhaus et al. [43]. A similar strategy was used to recover data on 663

cytochrome oxidase (CO) uptake from Livingstone and Hubel [38]. 664

Smoothness We define a smoothness score for a given map by comparing the tuning similarity for the nearest 665 model unit pairs to the tuning similarity of the least similar pairs. Concretely, given a vector x of pairwise tuning 666 similarity values, sorted in order of increasing cortical distance: 667

$$S(x) = \frac{\max(x) - x_0}{x_0}$$
(7)

#### Benchmarks comparing human VTC to model VTC-like layers 668

Stimuli. We evaluate the selectivity of neurons and model units to visual object categories using the "fLoc" functional 669 localizer stimulus set [76]. fLoc contains five categories, each with two subcategories consisting of 144 images 670 each. The categories are faces (adult and child faces), bodies (headless bodies and limbs), written characters 671 (pseudowords and numbers), places (houses and corridors), and objects (string instruments and cars). Selectivity 672 was assessed by computing the t-statistic over the set of functional localizer stimuli and defining a threshold above 673 which units were considered selective. 674

$$t = \frac{\mu_{\rm on} - \mu_{\rm off}}{\sqrt{\frac{\sigma_{\rm on}^2}{N_{\rm on}} + \frac{\sigma_{\rm off}^2}{N_{\rm off}}}},\tag{8}$$

where  $\mu_{on}$  and  $\mu_{off}$  are the mean responses to the "on" categories (e.g., adult and child faces) and "off" categories 675 (e.g., all non-face categories), respectively,  $\sigma^2$  are the associated variances of responses to exemplars from those 676 categories, and N is the number of exemplars being averaged over.

677

Human Data. We compare models to human data from the Natural Scenes Dataset (NSD) [75], a high-resolution 678 fMRI dataset of responses to 10,000 natural images in each of eight individuals (see Allen et al. for details). Models 679 are compared to two aspects of this dataset: single-trial responses to the main set of natural images per participant 680 (see "One-to-one mapping") and selectivity in response to the "fLoc" stimuli. Single-trial responses were z-scored 681 across images for each voxel and session and then averaged across three trial repeats. Selectivity was computed 682 on the "fLoc" experiment as described in the previous section, generating t-maps for each of the five categories for 683 each individual subject. 684

The VTC region of interest (ROI) was drawn based on anatomical landmarks to follow the convention in the literature 685 [113] and is provided in the NSD data release as the "Ventral" ROI in the "streams" parcellation. 686

Models. 687

65

Interactive Topographic Network (ITN) We reconstruct maps from a variant of the ITN in Blauch et al. [20] that was 688 trained and evaluated on the same images as the remaining models. 689

**DNN-SOM** Two related approaches for building SOM models of higher visual cortex have recently been published 690 [12, 13]. Because neither paper evaluates the resulting topographic maps with the fLoc stimuli, we reimplement 691 the approach of Zhang et al. [12] as follows. We extract the responses of each unit in the final layer of a pretrained 692 AlexNet to all 50,000 images in the ImageNet validation set. The responses are then reduced to the first four principal 693 components. The SOM is initialized as a 200 x 200 grid of model units with a Gaussian neighborhood function set 694 to  $\sigma = 6.2$ . The learning rate is set to 1.0 and the SOM is trained for 200,000 total iterations. The fLoc images are 695

presented to the pretrained AlexNet model and projected into the space spanned by the four principal components computed previously. The response of each model unit to each fLoc image is computed by taking the dot product of the unit weight matrix with the projected fLoc images. The SOM is then treated identically to the VTC-like layer of TDANN.

#### 700 Response Benchmarks.

**Representational similarity analysis** We compare functional properties of human VTC and models with representational similarity analysis (RSA) [72]. For any given model or human hemisphere, we compute a representational similarity matrix (RSM) as the pairwise Pearson's correlation between patterns of selectivity for each of the five fLoc categories. The diagonal of the RSM is trivially 1.0 and is ignored in further analysis. The similarity of two RSMs is computed as Kendall's  $\tau$ .

#### 706 **Topographic Benchmarks.**

Pairwise Tuning Difference We measure pairwise difference in VTC-like layer unit tuning as a function of cortical 707 distance. We draw 25 randoms samples of 500 units each. Each sample is filtered to include only units with a mean 708 response of at least 0.5 a.u.. For each fLoc category, the absolute pairwise difference in selectivity is computed 709 for pairs of units separated by different cortical distances. Curves are normalized by the chance value obtained by 710 randomly shuffling unit positions. Smoothness of maps is computed from these curves, same as in our analysis of V1. 711 To compare a model to a human hemisphere, we compute the mean category-by-category difference in smoothness, 712 e.g. comparing model face map smoothness to human face map smoothness, model body map smoothness to 713 human body map smoothness, etc. Permutation tests randomly assigning category-by-category smoothness profiles 714 to either "model" or "human" were used to assess the statistical significance of the mean difference in smoothness. 715

Patch Count and Size Patches are automatically detected in maps of category selectivity by identifying contiguous regions of highly-selective units (or voxels, for human VTC). Patch identification has a small number of parameters that can be adjusted for maps of different sizes and with different dynamic ranges of selectivity values. The first step in identifying patches is to smooth and interpolate discrete selectivity maps. The selectivity map is then thresholded, and contiguous islands surviving the threshold are retained as candidate patches. Each candidate patch is further filtered for reasonable size: patches must be at least  $100mm^2$  and no larger than  $45cm^2$ . Finally, the 2D geometry of the patch is constructed by fitting the concave hull of the points within the patch.

The following table identifies the relevant parameters for patch identification in human VTC and for each candidate model class.

Model	Selectivity Threshold	Smoothing $\sigma$	Minimum Size square mm	Maximum Size square mm
Human VTC	4	None	100	None
TDANN	2	2.4	100	4500
ITN	8	0.7	100	4500
DNN-SOM	10	2.4	100	4500

724

 Table 2. Patch detection parameters for human VTC and each model.

Selectivity Overlap We determine if units (or voxels, for human VTC) that are selective for a pair of categories overlap with one another as follows. First, we bin the cortical sheet into discrete square neighborhoods of width 10mm. In each neighborhood, the fraction of units selective for Category X and Category Y are recorded. We consider two populations as overlapping if there is a strong correlation between the proportions recorded across neighborhoods, i.e., if the frequency of Category 1 selectivity is predictive of Category Y selectivity and vice-a-versa. The X-Y Overlap score is computed as

$$Overlap = \frac{1 - RankCorr(X, Y)}{2},$$
(9)

where RankCorr is the Spearman's rank correlation coefficient and  $\vec{X}$  is the proportion of units selective for Category X in each cortical neighborhood. The category selectivity threshold was set at t > 4.

Linear regression. Neural predictivity is computed against a given dataset as the mean variance explained across 733 neurons and splits of the data. In practice we follow the parameters and design decisions made by the BrainScore 734 team [30]; they are repeated here for completeness. We use partial least squares (PLS) regression to predict the 735 activity of a given neuron as a linear weighted sum of model units in a given layer. Model activations are preprocessed 736 by first projecting unit responses to ImageNet images onto the first 1000 principal components, i.e. each component 737 is a linear mixture of model units. This projection is used when fitting on the stimuli that were shown to the animal. 738 When fitting IT, we use data from Majaj, Hong, et al., 2015 [32], which consists of multi-electrode array data in 739 responses to quasi-naturalistic scenes with a variety of objects on a variety of backgrounds. Variance explained is 740 corrected by dividing raw predictivity by the internal noise ceiling, a measure of the consistency of each recorded 741 742 neuron.

#### One-to-one mapping of visual cortical responses 743

A direct, one-to-one mapping between units and voxels is computed by assigning each unit in a layer of the network 744 to a single voxel based on responses to a given dataset. In practice, we correlate individual model unit activations to 745 the natural images from the Natural Scenes Dataset [75] with responses to these same images on the single voxel 746 level for a given subject. Unit-to-voxel assignments are determined using a polynomial-time optimal assignment 747 algorithm [114] which maximizes the overall average correlation between unit and voxel pairs, on a given training 748 set. The 515 shared images that all eight subjects viewed three times were held out as a test set and all reported 749 one-to-one correlations are calculated on this test set, using the unit-to-voxel assignments determined from training. 750 Each unit-to-voxel correlation is normalized by the individual voxel noise ceiling of that assigned voxel (see Allen et al. 75 for information on the calculation of the intra-individual voxel noise ceilings in NSD). One-to-one correlations were 752 calculated on an individual subject basis for each of the self-supervised and supervised models trained at each level 753 of the spatial weight  $\alpha$ . The inter-individual, or subject-to-subject, noise ceiling, was calculated in the same manner, 754 this time assigning voxels from one subject to voxels from another subject based on how correlated responses to 755 the shared 515 images were for each potential voxel pair. For the subject-to-subject assignment, we used an 80/20 756 train/test split and averaged results for each subject combination across 5 splits. A similar analysis will appear in a 757 forthcoming publication by Finzi et al. 758

#### Wiring Length 759

We measure the functional wiring length between two adjacent layers, the "source" layer and the "target" layer by 760 first identifying the units with the highest responses in each layer, then computing the length of inter-layer fibers that 761 would be required to connect them. First, for a given natural image input, we identify the top p% most responsive 762 units in each of two adjacent layers. We set p to 5% in the V1-like layers and 1% in the VTC-like layers. We note 763 that for computational tractability, we restrict our analysis to small neighborhoods in the V1-like layers and average 764

results across many random neighborhood selections. 765

Next, inter-layer fibers are added one by one, until all activated units in the earlier "source" layer are sufficiently 766 close to the location at which a fiber originates. In practice, we find the optimal fiber origination sites using the 767 k-means clustering algorithm, and continue adding fibers until the total "inertia" of the k-means clustering falls below 768 a specified threshold,  $k_{\rm thresh}$ . Inertia is computed as the sum of the squared distances between each activated unit 769 and its nearest fiber, and  $k_{\rm thresh}$  is set such that the mean distance from each unit to its nearest fiber is not greater 770 than  $d_{\rm thresh}$ .  $d_{\rm thresh}$  is set to 10.0mm in the VTC-like layer pairs, and is reduced to 0.9mm in the V1-like layer pairs 771 to reflect the smaller cortical neighborhood. Having established the number of inter-layer fibers required and their 772 origination sites in the "source" layer, we identify optimal termination sites for those fibers in the "target" layer as 773 follows. The set of target layer termination sites is identified as the centroids from k-means clustering, with k set to 774 the number of fibers. Finally, fibers are assigned between origination sites and termination sites with the linear sum 775 assignment algorithm, and the total wiring length is computed as the sum of the lengths of each individual inter-layer 776 fiber. 777

A critical decision when measuring wiring length in this way is how to situate units from two layers in a common 778 physical space. By design, each TDANN layer occupies a unique two-dimensional sheet, leaving the spatial 779 relationships between units in different cortical sheets undefined. Here, we assume that the "source" cortical sheet 780 and "target" cortical sheet lie in the same 2D plane, joined at one edge. Concretely, we can position the "target" 78 sheet to the left, right, above, or below the "source" layer. Without reason to choose one of these strategies, we 782 compute the optimal wiring length for each of the four options and report the average across all shift directions. 783

## 784 Dimensionality

In our analyses of dimensionality, we consider the responses of the full population of model units in each layer to a set of 10,112 natural images from the NSD [75]. Following [83], we perform spatial max-pooling on the convolutional feature maps, then compute the eigenspectrum of these responses. We summarize the dimensionality of the responses by their effective dimensionality (ED; Del Giudice [84]):

$$ED = \frac{\left(\sum_{i=1}^{N} \lambda_i\right)^2}{\sum_{i=1}^{N} \lambda_i^2},$$
(10)

where  $\lambda_i$  is the *i*th eigenvalue, and N is the number of eigenvectors.

#### 790 Microstimulation of model units on the simulated cortical sheet

We simulate the microstimulation of local populations of model units to 1) gain insight into the functional properties of local populations, and 2) measure effective connectivity between groups of units in adjacent layers. In all analyses, stimulation is performed by fixing the activity of units to values determined by a 2D Gaussian function. Units near the center of the Gaussian have their activity set to the maximal value, and activity falls off with distance from the center. We consider the top 5% of units, ranked by activity level, as being responsive in either the "Source" layer, where activity is set according to the 2D Gaussian, or in the following "Target" layer, where unit activity is determined by the network architecture and learned weights.

Functional Alignment In VTC-like layers, we measure functional alignment between layers by comparing the category selectivity of activated units in the Source layer (Layer 8) with the selectivity of responsive units in the Target layer (Layer 9). For each stimulation site, we compute the mean selectivity (*t*-statistic) of the top 5% most activated units for each of the following categories: faces, bodies, characters, cars, and places. This five-element "selectivity profile" can then be compared to the profile of the top 5% most strongly responding units in the Target layer by computing  $\chi^2$  distance between selectivity profiles. Similarity is then taken as the negative log distance and compared to a shuffle-control in which a random subset of units is compared instead of the top 5% most active units.

#### 805 Simulation of a Visual Cortical Prosthesis

In Figure 8, we demonstrate a proof of concept for using topographic DCNNs to prototype visual cortical prosthetic 806 devices. This proof of concept consists of two distinct stages: 1) generating device-achievable stimulation patterns 807 with a Stimulation Simulator, and 2) generating the estimated percept (Percept Synthesizer) that would result 808 by stimulating cortical areas with those patterns. To generate stimulation patterns, we feed a target image into 809 TDANN and record the precise activation magnitude of each model unit in each layer. If an infinitely high-precision 810 stimulation device with absolute coverage of the cortical sheet in all cortical areas were available, we would stimulate 811 cortex with this set of precise activation patterns. However, real stimulation devices are limited in many ways, 812 including limits to their spatial precision and the set of cortical areas they can access. Thus, we use TDANN 813 to produce device-achievable stimulation patterns, i.e., those that are consistent with the limitations of cortical 814 stimulation devices. Here we take a simple approach by considering degradation of high-precision patterns into 815 device-achievable patterns by Gaussian blurring. In each layer, we first interpolate the precise activity patterns onto 816 a high-resolution grid  $(2500 \times 2500 \text{ px})$ , then blur the resulting pattern with a 2D Gaussian kernel whose  $\sigma$  parameter 817 is set according to the desired blur level. Because different layers have different cortical sheet sizes (e.g. 70mm 818 on an edge in the VTC-like layer and 37mm on an edge in the V1-like layer), the width of the Gaussian in pixels is 819 variable, even though the width of the Gaussian in mm is constant. Finally, we perform a nearest-neighbor lookup 820 such that each model unit adopts the activity level of the pixel closest to its location. This set of activity patterns is the 821 final "device-achievable" pattern. The Stimulation Simulator also allows any specific subset of layers to be included; 822 e.g. the first two layers only, or all eight layers. We consider this restriction comparable to the limited access a neural 823 stimulation device might be restricted to. 824

Given a set of device-achievable activity patterns, we seek to determine the estimated percept that would be evoked 825 if that pattern were written into cortex, i.e., the visual input that is most consistent with those patterns. To this end, we 826 follow the example of Granley et al. [90] and use gradient-ascent image optimization methods to synthesize an image 827 such that the activity pattern produced by presenting that image is as close as possible to the device-achievable 828 target pattern. We use the *lucent* Python package to iteratively optimize an image to minimize the total mean squared 829 error, summed across layers, between the target activity patterns and the current evoked patterns at that iteration. 830 We optimize the image for 3000 steps at a learning rate of 0.05; further optimization has little effect on reducing the 83 mean squared error. The optimized result is the predicted percept for a given input image and theoretical cortical 832 stimulation device. 833

# **Author Contributions**

E.M. and D.F. performed analyses. E.M., K.G.-S., and D.L.K.Y. wrote the paper. H.L., J.J.D., and D.L.K.Y. originally

<sup>836</sup> conceived the approach.

## **Acknowledgements**

This work was supported by a National Science Foundation Graduate Research Fellowship awarded to E.M., a
 National Institutes of Health grant (RO1 EY 022318) awarded to K.G.-S., a Simons Foundation grant (543061)
 awarded to D.L.K.Y., a National Science Foundation CAREER grant (1844724) awarded to D.L.K.Y., and an Office
 of Naval Research grant (S5122) awarded to D.L.K.Y. We also thank the NVIDIA corporation and the Google TPU

Research Cloud group for hardware grants. We are grateful to Ben Sorscher for helpful discussions.

REFERENCES

## **References**

- 1. D H Hubel and T N Wiesel. Receptive fields, binocular interaction and functional architecture in the cat's visual cortex. *J. Physiol.*, 160:106–154, January 1962.
- Colin Humphries, Einat Liebenthal, and Jeffrey R Binder. Tonotopic organization of human auditory cortex.
   *Neuroimage*, 50(3):1202–1211, April 2010.
- 3. B M Harvey, B P Klein, N Petridou, and S O Dumoulin. Topographic representation of numerosity in the human parietal cortex. *Science*, 341(6150):1123–1126, September 2013.
- 4. Y C Wong, H C Kwan, W A MacKay, and J T Murphy. Spatial organization of precentral cortex in awake primates. I. Somatosensory inputs. *J. Neurophysiol.*, 41(5):1107–1119, September 1978.
- <sup>852</sup> 5. Horst A Obenhaus, Weijian Zong, R Irene Jacobsen, Tobias Rose, Flavio Donato, Liangyi Chen, Heping
   <sup>853</sup> Cheng, Tobias Bonhoeffer, May-Britt Moser, and Edvard I Moser. Functional network topography of the medial
   <sup>854</sup> entorhinal cortex. *Proc. Natl. Acad. Sci. U. S. A.*, 119(7), February 2022.
- 6. Yi Gu, Sam Lewallen, Amina A Kinkhabwala, Cristina Domnisoru, Kijung Yoon, Jeffrey L Gauthier, Ila R Fiete, and David W Tank. A Map-like Micro-Organization of Grid Cells in the Medial Entorhinal Cortex. *Cell*, 175(3):
   736–750.e30, October 2018.
- 7. Harry G Barrow, Alistair J Bray, and Julian M L Budd. A Self-Organizing Model of "Color Blob" Formation.
   *Neural Comput.*, 8(7):1427–1448, October 1996.
- 8. Teuvo Kohonen. Self-organized formation of topologically correct feature maps. *Biol. Cybern.*, 43(1):59–69, 1982.
- 9. K Obermayer, H Ritter, and K Schulten. A principle for the formation of the spatial structure of cortical feature maps. *Proc. Natl. Acad. Sci. U. S. A.*, 87(21):8345–8349, November 1990.
- 10. Richard Durbin and Graeme Mitchison. A dimensionality reduction framework for understanding cortical maps.
   *Letters to nature*, 343:644–647, 1990.
- N V Swindale and H Bauer. Application of Kohonen's self–organizing feature map algorithm to cortical maps
   of orientation and direction preference. *Proceedings of the Royal Society of London B: Biological Sciences*,
   265(1398):827–838, May 1998.
- 12. Yiyuan Zhang, Ke Zhou, Pinglei Bao, and Jia Liu. Principles governing the topological organization of object selectivities in ventral temporal cortex. September 2021.
- Fenil R Doshi and Talia Konkle. Visual object topographic motifs emerge from self-organization of a unified
   representational space. September 2022.
- R Linsker. From basic network principles to neural architecture: emergence of orientation columns. *Proc. Natl. Acad. Sci. U. S. A.*, 83(22):8779–8783, November 1986.
- <sup>875</sup> 15. K D Miller, J B Keller, and M P Stryker. Ocular dominance column development: analysis and simulation.
   <sup>876</sup> Science, 245(4918):605–615, August 1989.
- 16. K D Miller. A model for the development of simple cell receptive fields and the ordered arrangement of orientation columns through activity-dependent competition between ON- and OFF-center inputs. *J. Neurosci.*, 14(1):409–441, January 1994.
- 17. Miguel A Carreira-Perpiñán, Richard J Lister, and Geoffrey J Goodhill. A computational model for the development of multiple maps in primary visual cortex. *Cereb. Cortex*, 15(8):1222–1233, August 2005.
- 18. R A Jacobs and M I Jordan. Computational Consequences of a Bias toward Short Connections. *J. Cogn. Neurosci.*, 4(4):323–336, 1992.
- 19. A A Koulakov and D B Chklovskii. Orientation preference patterns in mammalian visual cortex: a wire length minimization approach. *Neuron*, 29(2):519–527, February 2001.
- Nicholas M Blauch, Marlene Behrmann, and David C Plaut. A connectivity-constrained computational account
   of topographic organization in primate high-level visual cortex. *Proc. Natl. Acad. Sci. U. S. A.*, 119(3), January
   2022.

- <sup>889</sup> 21. A Hyvärinen, P O Hoyer, and M Inki. Topographic independent component analysis. *Neural Comput.*, 13(7): <sup>890</sup> 1527–1558, July 2001.
- <sup>891</sup> 22. T Anderson Keller, Qinghe Gao, and Max Welling. Modeling Category-Selective Cortical Regions with <sup>892</sup> Topographic Variational Autoencoders. October 2021.
- 23. Alexander J E Kell, Daniel L K Yamins, Erica N Shook, Sam V Norman-Haignere, and Josh H McDermott. A
   Task-Optimized Neural Network Replicates Human Auditory Behavior, Predicts Brain Responses, and Reveals
   a Cortical Processing Hierarchy. *Neuron*, 98(3):630–644.e16, May 2018.
- 24. Daniel L Yamins, Ha Hong, Charles Cadieu, and James J DiCarlo. Hierarchical modular optimization of
   convolutional networks achieves representations similar to macaque IT and human ventral stream. *Adv. Neural Inf. Process. Syst.*, 26, 2013.
- 25. Daniel L K Yamins, Ha Hong, Charles F Cadieu, Ethan A Solomon, Darren Seibert, and James J DiCarlo.
   Performance-optimized hierarchical models predict neural responses in higher visual cortex. *Proc. Natl. Acad.* Sci. U. S. A., 111(23):8619–8624, 2014.
- 26. Seyed-Mahdi Khaligh-Razavi and Nikolaus Kriegeskorte. Deep supervised, but not unsupervised, models may explain IT cortical representation. *PLoS Comput. Biol.*, 10(11):e1003915, November 2014.
- 27. Umut Güçlü and Marcel A J van Gerven. Deep Neural Networks Reveal a Gradient in the Complexity of Neural Representations across the Ventral Stream. *J. Neurosci.*, 35(27):10005–10014, July 2015.
- 28. Santiago A Cadena, George H Denfield, Edgar Y Walker, Leon A Gatys, Andreas S Tolias, Matthias Bethge,
   and Alexander S Ecker. Deep convolutional models improve predictions of macaque V1 responses to natural
   images. *PLoS Comput. Biol.*, 15(4):e1006897, April 2019.
- 29. Jonas Kubilius, Martin Schrimpf, Kohitij Kar, Rishi Rajalingham, Ha Hong, Najib Majaj, Elias Issa, Pouya Bashivan, Jonathan Prescott-Roy, Kailyn Schmidt, Aran Nayebi, Daniel Bear, Daniel L Yamins, and James J
  <sup>910</sup> DiCarlo. Brain-Like Object Recognition with High-Performing Shallow Recurrent ANNs. In H Wallach, H Larochelle, A Beygelzimer, F Alche-Buc, E Fox, and R Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 12805–12816. Curran Associates, Inc., 2019.
- 30. Martin Schrimpf, Jonas Kubilius, Michael J Lee, N Apurva Ratan Murty, Robert Ajemian, and James J DiCarlo.
   Integrative Benchmarking to Advance Neurally Mechanistic Models of Human Intelligence. *Neuron*, September 2020.
- 31. Martin Schrimpf, Idan Asher Blank, Greta Tuckute, Carina Kauf, Eghbal A Hosseini, Nancy Kanwisher, Joshua B Tenenbaum, and Evelina Fedorenko. The neural architecture of language: Integrative modeling converges on predictive processing. *Proc. Natl. Acad. Sci. U. S. A.*, 118(45), November 2021.
- 32. N J Majaj, H Hong, E A Solomon, and J J DiCarlo. Simple Learned Weighted Sums of Inferior Temporal
   Neuronal Firing Rates Accurately Predict Human Core Object Recognition Performance. *Journal of Neuroscience*, 35(39):13402–13418, 2015.
- 33. Beth L Chen, David H Hall, and Dmitri B Chklovskii. Wiring optimization can relate neuronal structure and function. *Proc. Natl. Acad. Sci. U. S. A.*, 103(12):4723–4728, March 2006.
- <sup>925</sup> 34. Dario L Ringach, Robert M Shapley, and Michael J Hawken. Orientation selectivity in macaque V1: diversity and laminar dependence. *J. Neurosci.*, 22(13):5639–5651, July 2002.
- 35. R L De Valois, E W Yund, and N Hepler. The orientation and direction selectivity of cells in macaque visual cortex. *Vision Res.*, 22(5):531–544, 1982.
- 36. R L De Valois, D G Albrecht, and L G Thorell. Spatial frequency selectivity of cells in macaque visual cortex.
   *Vision Res.*, 22(5):545–559, 1982.
- <sup>931</sup> 37. S Zeki. Colour coding in the cerebral cortex: the reaction of cells in monkey visual cortex to wavelengths and colours. *Neuroscience*, 9(4):741–765, August 1983.
- 38. M S Livingstone and D H Hubel. Anatomy and physiology of a color system in the primate visual cortex. J.
   *Neurosci.*, 4(1):309–356, January 1984.

- 39. G G Blasdel and G Salama. Voltage-sensitive dyes reveal a modular organization in monkey striate cortex.
   *Nature*, 321(6070):579–585, 1986.
- 40. A Grinvald, E Lieke, R D Frostig, C D Gilbert, and T N Wiesel. Functional architecture of cortex revealed by optical imaging of intrinsic signals. *Nature*, 324(6095):361–364, 1986.
- 41. Tobias Bonhoeffer and Amiram Grinvald. Iso-orientation domains in cat visual cortex are arranged in pinwheel-like patterns. *Nature*, 353(6343):429–431, 1991.
- 42. M Hübener, D Shoham, A Grinvald, and T Bonhoeffer. Spatial relationships among three columnar systems in cat area 17. *J. Neurosci.*, 17(23):9270–9284, December 1997.
- 43. Ian Nauhaus, Kristina J Nielsen, Anita A Disney, and Edward M Callaway. Orthogonal micro-organization
   of orientation and spatial frequency in primate primary visual cortex. *Nat. Neurosci.*, 15(12):1683–1690,
   December 2012.
- 44. Shu-Chen Guan, Nian-Sheng Ju, Louis Tao, Shi-Ming Tang, and Cong Yu. Functional organization of spatial
   frequency tuning in macaque V1 revealed with two-photon calcium imaging. *Prog. Neurobiol.*, 205:102120,
   October 2021.
- 45. R Desimone, T D Albright, C G Gross, and C Bruce. Stimulus-selective properties of inferior temporal neurons in the macaque. *J. Neurosci.*, 4(8):2051–2062, August 1984.
- 46. C G Gross, C E Rocha-Miranda, and D B Bender. Visual properties of neurons in inferotemporal cortex of the Macaque. *J. Neurophysiol.*, 35(1):96–111, January 1972.
- 47. Mark A Pinsk, Kevin DeSimone, Tirin Moore, Charles G Gross, and Sabine Kastner. Representations of faces and body parts in macaque temporal cortex: a functional MRI study. *Proc. Natl. Acad. Sci. U. S. A.*, 102(19):
   6996–7001, May 2005.
- 48. Doris Y Tsao, Winrich A Freiwald, Roger B H Tootell, and Margaret S Livingstone. A Cortical Region Consisting
   Entirely of Face-Selective Cells. *Science*, 311(5761):670–674, February 2006.
- 49. Mark A Pinsk, Michael Arcaro, Kevin S Weiner, Jan F Kalkus, Souheil J Inati, Charles G Gross, and Sabine
   Kastner. Neural representations of faces and body parts in macaque and human cortex: a comparative FMRI
   study. J. Neurophysiol., 101(5):2581–2600, May 2009.
- 50. N Kanwisher, J McDermott, and M M Chun. The fusiform face area: a module in human extrastriate cortex specialized for face perception. *J. Neurosci.*, 17(11):4302–4311, June 1997.
- <sup>963</sup> 51. R Epstein and N Kanwisher. A cortical representation of the local visual environment. *Nature*, 392(6676):
   <sup>964</sup> 598–601, April 1998.
- 52. P E Downing, Y Jiang, M Shuman, and N Kanwisher. A cortical area selective for visual processing of the human body. *Science*, 293(5539):2470–2473, September 2001.
- <sup>967</sup> 53. Bruce D McCandliss, Laurent Cohen, and Stanislas Dehaene. The visual word form area: expertise for reading
   <sup>968</sup> in the fusiform gyrus. *Trends Cogn. Sci.*, 7(7):293–299, July 2003.
- <sup>969</sup> 54. Tanya Orlov, Tamar R Makin, and Ehud Zohary. Topographic representation of the human body in the <sup>970</sup> occipitotemporal cortex. *Neuron*, 68(3):586–600, November 2010.
- 55. Kevin S Weiner and Kalanit Grill-Spector. Not one extrastriate body area: using anatomical landmarks,
   hMT+, and visual field maps to parcellate limb-selective activations in human lateral occipitotemporal cortex.
   *Neuroimage*, 56(4):2183–2199, June 2011.
- 56. Kalanit Grill-Spector and Kevin S Weiner. The functional architecture of the ventral temporal cortex and its role in categorization. *Nat. Rev. Neurosci.*, 15(8):536–548, August 2014.
- 57. Blake A Richards, Timothy P Lillicrap, Philippe Beaudoin, Yoshua Bengio, Rafal Bogacz, Amelia Christensen, Claudia Clopath, Rui Ponte Costa, Archy de Berker, Surya Ganguli, Colleen J Gillon, Danijar Hafner, Adam Kepecs, Nikolaus Kriegeskorte, Peter Latham, Grace W Lindsay, Kenneth D Miller, Richard Naud, Christopher C Pack, Panayiota Poirazi, Pieter Roelfsema, João Sacramento, Andrew Saxe, Benjamin Scellier, Anna C Schapiro, Walter Senn, Greg Wayne, Daniel Yamins, Friedemann Zenke, Joel Zylberberg, Denis Therien, and Konrad P Kording. A deep learning framework for neuroscience. *Nat. Neurosci.*, 22(11): 1761–1770, November 2019.

- 58. Daniel L K Yamins and James J DiCarlo. Using goal-driven deep learning models to understand sensory cortex. *Nat. Neurosci.*, 19(3):356–365, 2016.
- 59. Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In Proceedings of the IEEE conference on computer vision and pattern recognition, pages 770–778, 2016.
- 60. Michael J Arcaro and Margaret S Livingstone. A hierarchical, retinotopic proto-organization of the primate visual system at birth. *Elife*, 6, July 2017.
- 61. Chengxu Zhuang, Siming Yan, Aran Nayebi, Martin Schrimpf, Michael C Frank, James J DiCarlo, and Daniel
   L K Yamins. Unsupervised neural network models of the ventral visual stream. *Proc. Natl. Acad. Sci. U. S. A.*,
   118(3), January 2021.
- 62. Talia Konkle and George A Alvarez. A self-supervised domain-general learning framework for human ventral stream representation. *Nat. Commun.*, 13(1):491, January 2022.
- 63. Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A Simple Framework for Contrastive Learning of Visual Representations. February 2020.
- <sup>996</sup> 64. J Deng, W Dong, R Socher, L J Li, Kai Li, and Li Fei-Fei. ImageNet: A large-scale hierarchical image database.
   <sup>997</sup> In 2009 IEEE Conference on Computer Vision and Pattern Recognition, pages 248–255, June 2009.
- 65. Soumya Chatterjee, Kenichi Ohki, and R Clay Reid. Chromatic micromaps in primary visual cortex. *Nat. Commun.*, 12(1):2315, April 2021.
- Matthias Kaschube, Michael Schnabel, Siegrid Löwel, David M Coppola, Leonard E White, and Fred Wolf.
   Universality in the evolution of orientation columns in the visual cortex. *Science*, 330(6007):1113–1116,
   November 2010.
- <sup>1003</sup> 67. Margaret Henderson and John T Serences. Biased orientation representations can be explained by experience <sup>1004</sup> with nonuniform training set statistics. *J. Vis.*, 21(8):10, August 2021.
- 68. Jeremy T Chang, David Whitney, and David Fitzpatrick. Experience-Dependent Reorganization Drives Development of a Binocularly Unified Cortical Representation of Orientation. *Neuron*, May 2020.
- Bardo N Ferreiro, Sergio A Conde-Ocazionez, João H N Patriota, Luã C Souza, Moacir F Oliveira, Fred Wolf,
   and Kerstin E Schmidt. Spatial clustering of orientation preference in primary visual cortex of the large rodent
   agouti. *iScience*, 24(1):101882, January 2021.
- To. Dario L Ringach, Patrick J Mineault, Elaine Tring, Nicholas D Olivas, Pablo Garcia-Junco-Clemente, and
   Joshua T Trachtenberg. Spatial clustering of tuning in mouse primary visual cortex. *Nat. Commun.*, 7:12270,
   August 2016.
- <sup>1013</sup> 71. Anupam K Garg, Peichao Li, Mohammad S Rashid, and Edward M Callaway. Color and orientation are jointly <sup>1014</sup> coded and spatially organized in primate primary visual cortex, 2019.
- <sup>1015</sup> 72. Nikolaus Kriegeskorte, Marieke Mur, and Peter Bandettini. Representational similarity analysis connecting
   <sup>1016</sup> the branches of systems neuroscience. *Front. Syst. Neurosci.*, 2(November):4, 2008.
- 73. Eshed Margalit, Keith W Jamison, Kevin S Weiner, Luca Vizioli, Ru-Yuan Zhang, Kendrick N Kay, and
   Kalanit Grill-Spector. Ultra-high-resolution fMRI of Human Ventral Temporal Cortex Reveals Differential
   Representation of Categories and Domains, 2020.
- 74. J V Haxby, M I Gobbini, M L Furey, A Ishai, J L Schouten, and P Pietrini. Distributed and overlapping
   representations of faces and objects in ventral temporal cortex. *Science*, 293(5539):2425–2430, September
   2001.
- Total Science And Science And
- 76. Anthony Stigliani, Kevin S Weiner, and Kalanit Grill-Spector. Temporal Processing Capacity in High-Level
   Visual Cortex Is Domain Specific. *J. Neurosci.*, 35(36):12412–12424, September 2015.

- <sup>1029</sup> 77. Kevin S Weiner and Kalanit Grill-Spector. Sparsely-distributed organization of face and limb activations in <sup>1030</sup> human ventral temporal cortex. *Neuroimage*, 52(4):1559–1573, 2010.
- T8. Hyodong Lee, Eshed Margalit, Kamila M Jozwik, Michael A Cohen, Nancy Kanwisher, Daniel L K Yamins, and
   James J DiCarlo. Topographic deep artificial neural networks reproduce the hallmarks of the primate inferior
   temporal cortex face processing network. July 2020.
- Johannes Mehrer, Courtney J Spoerer, Emer C Jones, Nikolaus Kriegeskorte, and Tim C Kietzmann. An
   ecologically motivated image dataset for deep learning yields better models of human vision. *Proc. Natl. Acad. Sci. U. S. A.*, 118(8), February 2021.
- 80. Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre
   Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Mahmoud Assran, Nicolas Ballas, Wojciech
   Galuba, Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael Rabbat, Vasu Sharma, Gabriel
   Synnaeve, Hu Xu, Hervé Jegou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. DINOv2:
   Learning Robust Visual Features without Supervision. April 2023.
- 81. U Guclu and M A J van Gerven. Deep Neural Networks Reveal a Gradient in the Complexity of Neural
   Representations across the Ventral Stream, 2015.
- Nathan C L Kong, Eshed Margalit, Justin L Gardner, and Anthony M Norcia. Increasing neural network
   robustness improves match to macaque V1 eigenspectrum, spatial frequency preference and predictivity.
   *PLoS Comput. Biol.*, 18(1):e1009739, January 2022.
- 83. Eric Elmoznino and Michael F Bonner. High-performing neural network models of visual cortex benefit from
   high latent dimensionality. February 2023.
- <sup>1049</sup> 84. Marco Del Giudice. Effective Dimensionality: A Tutorial. *Multivariate Behav. Res.*, 56(3):527–542, 2021.
- 85. Carsen Stringer, Marius Pachitariu, Nicholas Steinmetz, Matteo Carandini, and Kenneth D Harris.
   High-dimensional geometry of population responses in visual cortex. *Nature*, 571(7765):361–365, July 2019.
- 1052
   86. Dmitri B Chklovskii, Thomas Schikorski, and Charles F Stevens. Wiring optimization in cortical circuits.
   1053
   Neuron, 34(3):341–347, 2002.
- 87. Sebastian Moeller, Winrich A Freiwald, and Doris Y Tsao. Patches with links: a unified system for processing
   faces in the macaque temporal lobe. *Science*, 320(5881):1355–1359, June 2008.
- 1056
   88. Michael S Beauchamp, Denise Oswalt, Ping Sun, Brett L Foster, John F Magnotti, Soroush Niketeghad, Nader
   1057
   Pouratian, William H Bosking, and Daniel Yoshor. Dynamic Stimulation of Visual Cortex Produces Form Vision
   1058
   in Sighted and Blind Humans. *Cell*, 181(4):774–783.e5, May 2020.
- Maureen van der Grinten, Jaap de Ruyter van Steveninck, Antonio Lozano, Laura Pijnacker, Bodo Rückauer,
   Pieter Roelfsema, Marcel van Gerven, Richard van Wezel, Umut Güçlü, and Yağmur Güçlütürk. Biologically
   plausible phosphene simulation for the differentiable optimization of visual cortical prostheses. December
   2022.
- 90. Jacob Granley, Alexander Riedel, and Michael Beyeler. Adapting Brain-Like Neural Networks for Modeling
   Cortical Visual Prostheses. September 2022.
- 91. Elia Shahbazi, Timothy Ma, Martin Pernus, Walter J Scheirer, and Arash Afraz. The causal role of the inferior
   temporal cortex in visual perception. January 2023.
- Matharina Dobs, Julio Martinez, Alexander J E Kell, and Nancy Kanwisher. Brain-like functional specialization
   merges spontaneously in deep neural networks. *Science Advances*, 8(11):eabl8913, 2022.
- 93. Sam V Norman-Haignere, Jenelle Feather, Dana Boebinger, Peter Brunner, Anthony Ritaccio, Josh H
   McDermott, Gerwin Schalk, and Nancy Kanwisher. A neural population selective for song in human auditory
   cortex. *Curr. Biol.*, 32(7):1470–1484.e12, April 2022.
- 94. Rosemary A Cowell and Garrison W Cottrell. What evidence supports special processing for faces? A cautionary tale for fMRI interpretation. *J. Cogn. Neurosci.*, 25(11):1777–1793, November 2013.
- <sup>1074</sup> 95. Lukas Vogelsang, Sharon Gilad-Gutnick, Evan Ehrenberg, Albert Yonas, Sidney Diamond, Richard Held,
   <sup>1075</sup> and Pawan Sinha. Potential downside of high initial visual acuity. *Proc. Natl. Acad. Sci. U. S. A.*, 115(44):
   <sup>1076</sup> 11333–11338, October 2018.

- 96. Omisa Jinsi, Margaret M Henderson, and Michael J Tarr. Early experience with low-pass filtered images facilitates visual category learning in a neural network model. *PLoS One*, 18(1):e0280145, January 2023.
- 97. A Nayebi, J Sagastuy-Brena, D M Bear, K Kar, J Kubilius, S Ganguli, D Sussillo, J J DiCarlo, and D L K Yamins.
   Recurrent Connections in the Primate Ventral Visual Stream Mediate a Tradeoff Between Task Performance and Network Size During Core Object Recognition. *Neural Computation*, 34(8):1652–1675, July 2022.
- 98. Jianhua Cang, Megumi Kaneko, Jena Yamada, Georgia Woods, Michael P Stryker, and David A Feldheim.
   Ephrin-as guide the formation of functional maps in the visual cortex. *Neuron*, 48(4):577–589, November
   2005.
- M Meister, R O Wong, D A Baylor, and C J Shatz. Synchronous bursts of action potentials in ganglion cells of
   the developing mammalian retina. *Science*, 252(5008):939–943, May 1991.
- 1087 100. Jinwoo Kim, Min Song, Jaeson Jang, and Se-Bum Paik. Spontaneous Retinal Waves Can Generate 1088 Long-Range Horizontal Connectivity in Visual Cortex. *J. Neurosci.*, 40(34):6584–6599, August 2020.
- 101. Todd McLaughlin, Christine L Torborg, Marla B Feller, and Dennis D M O'Leary. Retinotopic map refinement requires spontaneous retinal waves during a brief critical period of development. *Neuron*, 40(6):1147–1160, December 2003.
- 102. Xinxin Ge, Kathy Zhang, Alexandra Gribizis, Ali S Hamodi, Aude Martinez Sabino, and Michael C Crair. Retinal
   waves prime visual motion detection by simulating future optic flow. *Science*, 373(6553), July 2021.
- 103. Rishi Rajalingham and James J DiCarlo. Reversible Inactivation of Different Millimeter-Scale Regions of
   Primate IT Results in Different Patterns of Core Object Recognition Deficits. *Neuron*, 102(2):493–505.e5,
   April 2019.
- <sup>1097</sup> 104. Tiago Marques, Martin Schrimpf, and James J DiCarlo. Multi-scale hierarchical neural network models that <sup>1098</sup> bridge from single neurons in the primate primary visual cortex to object recognition behavior. March 2021.
- 105. Santiago A Cadena, Konstantin F Willeke, Kelli Restivo, George Denfield, Fabian H Sinz, Matthias Bethge,
   Andreas S Tolias, and Alexander S Ecker. Diverse task-driven modeling of macaque V4 reveals functional
   specialization towards semantic tasks. May 2022.
- 102 106. Pouya Bashivan, Kohitij Kar, and James J DiCarlo. Neural population control via deep image synthesis.
   1103 Science, 364(6439), May 2019.
- 107. Priya Goyal, Quentin Duval, Jeremy Reizenstein, Matthew Leavitt, Min Xu, Benjamin Lefaudeux, Mannat Singh, Vinicius Reis, Mathilde Caron, Piotr Bojanowski, and Others. Vissl, 2021.
- 108. Ilya Loshchilov and Frank Hutter. SGDR: Stochastic Gradient Descent with Warm Restarts. August 2016.
- 109. Noah C Benson, Jennifer M D Yoon, Dylan Forenzo, Stephen A Engel, Kendrick N Kay, and Jonathan Winawer.
   Variability of the Surface Area of the V1, V2, and V3 Maps in a Large Sample of Human Observers. *J. Neurosci.*, 42(46):8629–8646, November 2022.
- 110. T Yoshioka, J B Levitt, and J S Lund. Intrinsic lattice connections of macaque monkey visual cortical area V4.
   J. Neurosci., 12(7):2785–2802, July 1992.
- 1112 111. Alex Krizhevsky, Geoffrey E Hinton, and Ilya Sutskever. ImageNet Classification with Deep Convolutional 1113 Neural Networks. *the Neural Information Processing Systems Foundation 2012 conference*, pages 1–9, 2012.
- 1114 112. Giuseppe Vettigli. MiniSom: minimalistic and NumPy-based implementation of the Self Organizing Map, 2018.
- 113. Lior Bugatus, Kevin S Weiner, and Kalanit Grill-Spector. Task alters category representations in prefrontal but 1116 not high-level visual cortex. *Neuroimage*, 155:437–449, July 2017.
- 1117 114. James Munkres. Algorithms for the Assignment and Transportation Problems. *Journal of the Society for* 1118 *Industrial and Applied Mathematics*, 5(1):32–38, March 1957.

## **Supplementary Information**

#### 1120 V1-like maps produced with alternative feature sets

Figure 2 demonstrates that co-training for spatial and task losses is sufficient to generate V1-like topography. 1121 However, we have not ruled out the possibility that generating orientation-selective units and arranging them on 1122 the cortical sheet via other strategies could produce V1-like maps. To address this concern, we derive orientation 1123 prefrence maps (OPMs) from three different strategies for learning and spatially organizing model units. We first 1124 compare the standard TDANN, in which unit positions are fixed prior to training and model weights are optimized to 1125 minimize both task and spatial losses, to a Task Only DCNN whose weights are optimized only for the task loss. To 1126 generate an OPM from the Task Only model, we freeze network weights then iteratively shuffle model units on the 1127 cortical sheet such that the Spatial Loss is minimized post-hoc. Accordingly, we refer to this model as a "Post-hoc" 1128 arrangement of DCNN features. We find that OPM smoothness is nearly identical when co-learning features with 1129 the spatial loss (i.e., TDANN) than when first learning features and then post-hoc arranging units in the cortical sheet 1130 (Supplementary Figure S4). A third alternative is to bypass the learning of features altogether and use a hard-coded 113 Gabor filterbank (GFB) to generate model units, as has been suggested as a model of V1 neuron tuning (e.g. Jones 1132 and Palmer [5], Dapello et al. [3]). Following the same approach as in the Task Only model for deriving OPMs, we 1133 find that the hard-coded GFB features fail to produce a smooth OPM. How can we reconcile the apparent inadequacy 1134 of the Gabor filterbank in generating V1-like topography with its strong orientation selectivity? One possibility is that 1135 a Gabor filterbank lacks the required complexity to form responses to natural images that drive brain-like topography, 1136 but that simpler stimuli may improve the accuracy of its topographic predictions. To test whether the nature of the 1137 images presented to the model matters, we evaluated the same three feature sets (TDANN, Post-hoc, Task Only, 1138 and GFB) on a set of simple sine grating images (Figure S4a, bottom). Interestingly, we find that the TDANN, 1139 Post-hoc Task Only, and GFB feature sets all produce smooth OPMs when their units are organized with respect 1140 to correlations of sine grating responses. We conclude that the TDANN is the only model that, by co-learning 1141 features and topography, is able to produce brain-like OPMs from realistically complex natural inputs. Task Only 1142 and Hand-Crafted feature spaces are capable of producing V1-like OPMs only when presented with simple inputs, 1143 whereas the core advantage of the TDANN is its ability to learn a feature space that produces brain-like functional 1144 organization in the presence of realistically complex natural images. 1145

#### 1146 Natural image inputs are required for the emergence of brain-like functional organization

Work in developmental neuroscience and psychology has called into consideration the influence of visual experience 1147 on the development of structure and function in visual cortex. We leveraged the ability of self-supervised TDANNs 1148 to predict functional organization after learning from unlabeled visual data streams to determine which inputs might 1149 drive the emergence of brain-like topographic maps. We evaluated networks trained on four distinct image datasets, 1150 including the natural image datasets ImageNet and Ecoset [9], and two artificial datasets: sine gratings and white 1151 noise images. We find that for both natural image datasets, there is brain-like functional organization of V1-like 1152 and VTC-like layers. In the V1-like layer, 14% of units in the Ecoset-trained network and 20% of units in the 1153 ImageNet-trained network were strongly orientation selective (circular variance < 0.6), and we observe smooth 1154 OPMs with pinwheels in models trained from both datasets. Further, we found similar numbers of VTC-like layer 1155 units with selectivity t > 5 in both models (12.7% for Ecoset and 14.2% for ImageNet), and we detect patches 1156 selective for all five categories in both models (Figure S9). 1157

While the suitability of naturalistic stimuli for generating brain-like functional organization may not be surprising, we 1158 wanted to test if simpler artificial datasets could succeed in matching neural data for two reasons. First, it has 1159 been demonstrated that patterned endogenous activity prior to eye opening can establish visual cortical circuitry [4]. 1160 Second, if artificial synthetic stimuli were suitable for constructing brain models, we could avoid needing to collect 1161 large natural image datasets. We trained TDANN on two artificial stimulus sets: a set of sine grating stimuli that may 1162 loosely mirror endogenous activity patterns, and Gaussian white noise images. The grating-trained model exhibited 1163 a very high fraction of strongly orientation selective units in the V1-like layer (73%). However, the grating-trained 1164 model had no category selectivity (1.2% of units selective at t > 5, averaged across categories), and no detectable 1165 patches. Thus, simple oriented stimuli may be sufficient to drive V1-like map formation, but natural stimuli are 1166 necessary to develop the remainder of the ventral visual pathway. We next evaluated a model trained on white noise, 1167 which allows us to isolate the effects of the model architecture and loss functions in the absence of structure in the 1168 input data. We find that training on white noise prevents the learning of strongly orientation-selective units in the 1169 V1-like layer (0% of units with circular variance < 0.6) or strongly category-selective units in the VTC-like layer (4% 1170 of units with t > 5). Surprisingly, however, the noise-trained TDANN does learn some weak functional organization. 117 In the V1-like layer, a weak orientation preference map is formed, and in the VTC-like layer, two character-selective 1172 patches and one face-selective patch is observed. These results suggest that the spatial loss is able to produce some 1173

topographic structure even in the absence of patterned inputs, although the strength of the selectivity is extremely
 weak. Taken together, these analyses of the impact of training data on functional organization support the necessity
 and sufficiency of natural images for the emergence of robust V1-like and VTC-like topographic maps.

## 1177 Probing the tuning of unit populations outside of category-selective patches

If the VTC-like layer smoothly encodes a space of objects, we might expect that images synthesized to drive high 1178 responses in nearby regions of the cortical sheet would be perceptually similar. Indeed, we find that optimal image 1179 characteristics smoothly vary across the cortical surface. Higher spatial frequency and rectilinear features dominate 1180 the upper right sides of the map, while curvilinear and lower spatial frequency features best drive the top and bottom 118 edges. We find that these optimal images also align with category-selectivity, e.g., the input images that best drive 1182 units in face-selective patches tend to contain eyes and fall in the more curvilinear regions of feature space. Images 1183 synthesized to maximize regions that fall between patches (sites 5, 10, 11, 15, 16, and 20) lack clearly discernible 1184 object categories, but nonetheless follow the smooth gradients of image features across the cortical surface. Thus, 1185 it appears that the VTC-like layer learns a smooth mapping of object space in two dimensions, and that patches 1186 emerge as regions of that space that align with the category localization stimuli that we use to probe the model. 1187

### 1188 Supplemental Methods

**Dimensionality Summarize by Power Law Exponent.** Following Kong et al. [7], Stringer et al. [11], we summarize the eigenvalues by fitting a line to the log-log plot of eigenvalues against their principal component index, and report the absolute value of the best fit line as the power law exponent. To prevent fitting to nonlinear regions in the earliest and latest parts of the distribution, the line is fit from the 2nd to the 50th principal component.

Linear regression. Neural predictivity is computed against a given dataset as the mean variance explained across 1193 neurons and splits of the data. In practice we follow the parameters and design decisions made by the BrainScore 1194 team [10]; they are repeated here for completeness. We use partial least squares (PLS) regression to predict the 1195 activity of a given neuron as a linear weighted sum of model units in a given layer. Model activations are preprocessed 1196 by first projecting unit responses to ImageNet images onto the first 1000 principal components, i.e. each component 1197 is a linear mixture of model units. This projection is used when fitting on the stimuli that were shown to the animal. 1198 When fitting V1, we use data from Cadena et al. [2], which consists of single-neuron recordings to a set of natural 1199 images. When fitting V4 and IT, we use data from Majaj, Hong, et al., 2015 [8], which consists of multi-electrode 1200 array data in responses to quasi-naturalistic scenes with a variety of objects on a variety of backgrounds. Variance 1201 explained is corrected by dividing raw predictivity by the internal noise ceiling, a measure of the consistency of each 1202 recorded neuron. 1203

Unit Clustering. The degree to which of responses to natural images are clustered is computed by considering the
 locations of the 5% of units that respond most strongly to a given input image. We compute the distribution of
 pairwise distances between these highly-active units, then count the number of pairs that are within 10.0mm of each
 other: if the count is high, then the active units are concentrated into a small number of clusters. Finally, clusterness
 is defined as the ratio between the number of nearby pairs in the true response pattern to the number of nearby pairs
 when locations are randomly shuffled. We compute results for 10 random position shuffles, 64 randomly-selected
 images used as input, and five random initial seeds for each model.

**Gabor Filter Bank (GFB).** In Figure S4, we generate responses from a Gabor filter bank by following the VOneNet implementation in Dapello et al. [3]. For computational tractability and to produce a similar quantity of units as in the TDANN V1-like layer, we reduce the number of simple and complex channels from 256 to 64 each, and increase the stride of the convolution from 4 to 8 pixels. The resulting filter bank is then treated identically to the TDANN V1-like layer when extracting responses and constructing tuning curves.

The orientation preference map (OPM) for the GFB model is produced by assigning GFB outputs to random initial positions, then minimizing the Spatial Loss by iteratively swapping the locations of randomly-selected pairs of units as described above.

**Stimulus Optimization.** We use image synthesis methods, implemented in the *lucent* Python package (https:// github.com/greentfrapp/lucent), to generate images which reproduce patterns of stimulation. Specifically, we synthesize an input image that minimizes the mean squared error between a desired pattern of activity and the actual pattern obtained by presenting the synthesized image to the network. The desired pattern of activity is set according to a two-dimensional Gaussian centered over some region of the cortical sheet. In these experiments we set the parameter of the Gaussian to 3.5mm. For efficiency, we also remove units far from the center of the Gaussian from the computation of the mean squared error: units below 10% of the height of the Gaussian are ignored. All



**Figure S1. Minimization of loss components during training. (a)** Task loss throughout training. **(b)** Spatial loss in each of the eight convolutional layers during training. Shaded area: 95% confidence interval (CI) across random initializations. **(c)** Response correlation decreases as a function of the cortical distance between model unit pairs in each model layer. Shaded region: 95% CI from repeated sampling of different cortical neighborhoods in each layer. **(d)** Portions of the cortical sheet from each of four convolutional layers; units colored according to their correlation with an arbitrarily selected seed unit, marked by the black star.

synthesized images begin as 128 x 128 pixels of white noise and are optimized for 1,024 steps. We retain the default
 settings for image transforms, which include optimization in the Fourier basis, color channel decorrelation, jittering,
 rotation, and padding.

*Retinal Waves.* In Figure S18 we organize DCNN units in the cortical sheet according to their response correlations to a series of simulated retinal waves.

**Creating Retinal Waves** Our simulation of retinal wave activity is heavily inspired by the description in Kim et al. [6]. We simulate the retina as a two-dimensional circle of radius 320px. The retina has three spatially-overlapping cell layers: one for ON-RGCs (retinal ganglion cells), one for OFF-RGCs, and one for amacrine cells. Each cell can



**Figure S2. Selection of V1-like and VTC-like layers**. (a) Variance explained by linear regression of TDANN layer outputs to measurements in macaque V1, V4, and IT. V1 predictivity peaks in the first three layers, whereas IT predictivity peaks in the last two layers. (b) Fraction of units strongly orientation selective in each layer. (c) Fraction of units that are strongly selective for each category (*t*-value > 10) in each layer.



**Figure S3. Topographic and representational benchmarks in the V1-like model layer (a)** Orientation, spatial frequency, and chromatic preference maps for all candidate model types. **(b)** Left: Distribution of preferred orientations for each model type. Right: Cardinality index, computed as the fraction of units selective for cardinal orientations to units selective for the obliques. Dashed green light indicates value in macaque V1. **(c)** Left: Distribution of circular variance for each model and for macaque V1. Vertical line indicates cutoff for strong selectivity. Right: percentage of units strongly selective for orientations in each model type.

<sup>1234</sup> be in one of four states: inhibited, recruitable (but not currently active), refractory (recently active but not recruitable <sup>1235</sup> yet), or active (currently "on"). Cells are connected to each other according to the following rules: 1) ON-RGCs are <sup>1236</sup> connected to one another in an excitatory fashion within a radius of  $r_{on}$ , 2) ON-RGCs are connected to amacrine <sup>1237</sup> cells in an excitatory fashion within a radius of  $r_{on}$ , and amacrine cells inhibit OFF-RGCs within a radius of  $r_{amacrine}$ .

A wave is initiated by setting some subset of the ON-RGCs to the "active" state. The activated subset is determined by picking a random location along the edge of the retina and activating cells along a thin strip at that location. The wave is then propagated for up to t timesteps (propagation is halted if the wave runs off screen and all cells are off). At each timestep, activity is propagated as follows. First, all cells that have been active longer than a specified



Figure S4. Smoothed OPMs from alternative feature spaces (a) Top row: smoothed OPMs from the TDANN, a Task Only model with post-hoc unit organization, and a Gabor Filterbank with post-hoc unit organization, where units are brought closer together if they have similar responses to ImageNet images. Bottom row: same as top, but with unit proximity optimized with respect to sine grating image responses. (b) Pairwise orientation tuning difference over distance, and corresponding smoothness scores, for the maps in (a).



Figure S5. Orientation preference maps (OPMs) and pinwheel density in alternative models For demonstration, all models in this figure had unit positions organized post-hoc to achieve a strong OPM, i.e., they are not proper TDANNs. (a) OPM in a small region of the standard ResNet-18 TDANN. Pinwheels are shown by black and white dots. (b) OPM in the V1-like layer of a categorization-trained ResNet-50, in which the increased number of channels allows a reduction of cortical neighborhood size and, accordingly, a dramatic increase in pinwheel density. (c) OPM in the V1-like layer of a categorization-trained ResNet-18 with twice the number of channels in each layer, in which the increased number of channels allows a reduction of cortical neighborhood size and, accordingly, a dramatic increase in pinwheel density.

"active duration" are set to the refractory state. Second, cell activity levels are updated by multiplying the connectivity 1242 matrices with the previous activity states. Third, we activate all ON-RGCs who are in the recruitable state and whose 1243 activity exceeds an activity threshold of  $t_{\rm ON-active}$ . Fourth, we inhibit all OFF-RGCs whose activity falls below a 1244 threshold of  $t_{OFF-active}$ . OFF-RGCs whose activity passes that threshold are activated if they are currently in the 1245 recruitable state. Finally, amacrine cells whose activity exceeds  $t_{\rm amacrine-active}$  are set to active. The remainder of 1246 the amacrine cells are made recruitable instead. Images of the simulated activity at each stage are produced by 1247 creating binary masks of the locations of active ON-RGCs. Half of the waves are randomly assigned to map the 1248 binary images to a black and white colormap, and the remainder are assigned to a red and green colormap. 1249

In this work, we produce retinal waves with two sets of parameters. The following parameters are common to both sets of retinal waves:  $r_{\rm on} = 15$ ,  $r_{\rm amacrine} = 1.5$ , t = 20,  $t_{\rm ON-active} = 7$ ,  $t_{\rm amacrine-active} = 0.1$ ,  $t_{\rm OFF-active} = 0.1$ . In one of the two sets of waves, the active duration is set to 100ms, and in the other, the active duration is set to 200ms. In practice, the waves produced with the longer active duration are twice as thick.

Measuring Responses to Waves Each wave consists of a number of images, one per timestep of the simulation. Because the simulated retina is circular, the corners of each image never contain simulated activity. To make better use of each image, we take a central square crop of each image of size M x M pixels then resize the image back to 224 x 224 pixels. M is selected such that all regions of the crop contain activity: for an image of size 224px,



**Figure S6. Data in each human subject from the NSD fLoc experiment and patch detection protocol** All scale bars: 2mm. **(a)** Map of face selectivity in the right-hemisphere VTC region of interest (ROI) for one example subject. A: anterior, M: medial, L: lateral, P: posterior. **(b)** Thresholded face selectivity map for the same subject as (a). **(c)** Category selectivity map for all five fLoc categories. **(d)** Patches detected from the category selective clusters in (c). **(e)** Detected patches in each hemisphere (LH = left hemisphere, RH = right hemisphere) for each subject.



**Figure S7. Optimal stimuli throughout the VTC-like layer (a)** VTC-like layer of an example TDANN model. Overlaid numbers correspond to sub-panels in (b). (b) Images synthesized to maximally activate a local population of units centered at the indicated location in (a). Small dot in bottom left of each image indicates the patch membership of that location, e.g., a red dot indicates that the image optimally drives units that happen to be in a face-selective patch.

 $M = \sqrt{2 \times (\frac{224}{2})^2} \approx 158$ . As with all other images presented to the DCNN models, the images are then preprocessed and normalized.

For a wave with *t* timesteps, each model unit produces *t* responses. We integrate responses to each wave by computing the mean response across all waves. Anecdotally, similar results are achieved by computing the maximum response instead of the mean. Unit-to-unit correlations are then computed by considering the vector of integrated responses for each wave. We use the unit-to-unit correlations to perform swap-based organization of units on the cortical surface such that correlated units are moved to be nearby one another.



Figure S8. Topographic maps for models trained with the Relative SL and the Supervised Categorization objective (a) Orientation preference maps (OPMs) in the V1-like layer of models at each level of  $\alpha$  trained from five different random seeds with the categorization objective. A region of each cortical sheet is shown, with black and white dots indicating locations of detected clockwise and counter-clockwise pinwheels, respectively. Gray square covers the Task Only seed 0, which was used during position initialization. (b) Category selectivity maps for the VTC-like layer of each model in (A). Plotting conventions as in Figure 3.



Figure S9. Topographic maps from models trained with different training datasets (a) Orientation, spatial frequency, and chromatic preference maps in the V1-like layer of models trained with ImageNet images, the *Ecoset* training set, a set of hand-selected sine gratings (increased  $\alpha = 10$ ), and Gaussian white noise images. (b) Representational similarity between human VTC and models trained with each dataset. Error bar: 95% CI across human hemispheres. (c) Category selectivity maps for the VTC-like layer of each model in (a). Plotting conventions as in Figure 3.



Figure S10. Topographic maps for models trained with the Relative SL and self-supervision (a) Orientation preference maps (OPMs) in the V1-like layer of models at each level of  $\alpha$  trained from five different random seeds with the Relative Spatial Loss (SL). A region of each cortical sheet is shown, with black and white dots indicating locations of detected clockwise and counter-clockwise pinwheels, respectively. (b) Category selectivity maps for the VTC-like layer of each model in (a). Plotting conventions as in Figure 3.



Figure S11. Topographic maps for models trained with the Absolute SL (a) Orientation preference maps (OPMs) in the V1-like layer of models at each level of  $\alpha$  trained from five different random seeds with the Absolute Spatial Loss (SL). A mm region of each cortical sheet is shown, with black and white dots indicating locations of detected clockwise and counter-clockwise pinwheels, respectively. (b) Category selectivity maps for the VTC-like layer of each model in (a). Plotting conventions as in Figure 3.



Figure S12. Additional comparison of models trained with different task and spatial objectives (a) Spatial loss in the VTC-like layer of TDANN models (purple), categorization-trained models (gold), and models trained with the Absolute SL (red) throughout training. (b) Categorization accuracy (top-1 ImageNet validation set performance) for models trained at each level of  $\alpha$  with either the Relative (purple) or Absolute (red) SL. (c) Categorization accuracy for models trained at each level of  $\alpha$  directly on the supervised categorization objective.



Figure S13. Dimensionality of model unit populations as a function of training objective and  $\alpha$  (a) Variance explained by each principal component (PC) for each layer of TDANNs trained at different levels of the spatial weight magnitude  $\alpha$ . Components computed from responses to 10,000 images from the NSD [1]. (b) Variance explained by each principal component in the VTC-like layer of models trained with  $\alpha = 0.25$  and different objectives. (c) Power law coefficient fit to eigenspectra from the VTC-like layer of models trained with  $\alpha = 0.25$  and different objectives.



Figure S14. Clustering of responses to natural images as a function of  $\alpha$  (a) Strength of activation in the TDANN VTC-like layer to an arbitrarily-selected natural image, for models trained at different levels of the spatial weight ( $\alpha$ ). (b) Probability density function of pairwise distances between pairs of activated units for each model type, computed over repeated presentations of different natural images. Curve color indicates the level of the spatial weight ( $\alpha$ ) that model was trained with. (c) Clusterness, measured as the increase in unit density above the chance of value (dashed line: 1.0). Error bars: 95% Cl over different random initial model seeds and images used to generate responses. ANOVA:  $F(7,32) = 70.5, p < 10^{-16}$ , post-hoc Tukey's tests: significantly lower clusterness for  $\alpha = 0$  and Unoptimized models compared to models with  $\alpha > 0$ , all post-hoc ps < .001.



Figure S15. Prediction of neural firing rates with linear regression compared against topographic map smoothness (a) Models trained at different levels of  $\alpha$  (represented by dot size) and with different objectives compared in their capacity to predict macaque V1 firing rates (Var Exp) and the smoothness of their orientation preference maps. b) As in (b), but for prediction of firing rates in macaque inferotemporal cortex (IT) and smoothness of face selectivity maps. No difference in variance explained when  $\alpha < 25$ , all pairwise *ps* from Mann-Whitney tests p > 0.42.



Figure S16. Topographic maps in each layer of a representative TDANN model. (a) Orientation, spatial frequency, and chromatic preference maps in each layer. Plotting conventions as in Figure 2. (b) Category selectivity map in each layer. (c) Orientation and color selectivity in the V4-like model layer. Units in magenta are selective for color and not orientation, units in green are selective for orientation and not color, and units in black are selective for both orientation and color. Similar data in macaque V4 is shown in the inset at bottom right (from [12]).



**Figure S17. Topographic maps in a representative TDANN throughout training. (a)** OPMs in the V1-like layer at initialization (left), and after 1, 11, 51, 101, and 200 epochs of training. **(b)** Category selectivity maps in the VTC-like model layer at each timepoint. **(c)** Smoothness as a function of training step for orientation, spatial frequency, and color preference maps. Smoothness peaks early then plateaus. **(d)** Selectivity of category selectivity maps for each fLoc category. Smoothness increases throughout training.



**Figure S18.** Simulated retinal waves can drive unit-to-unit correlations comparable to static sine gratings. (a) Five example frames from a simulated retinal wave movie. The responses to each frame are integrated to compute the mean response to each wave. (b) OPMs created by post-hoc organization of units in the V1-like layer of a Task Only SimCLR model, when the unit-to-unit correlations are computed by presenting retinal wave movies (left), a dataset of sine gratings (middle), or natural images (right). Scale bar: 2mm.

## **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., Hutchinson, J. B., Naselaris, T., and Kay, K. A massive 7T fMRI dataset to bridge cognitive
   neuroscience and artificial intelligence. *Nat. Neurosci.*, 25(1):116–126, Jan. 2022.
- Cadena, S. A., Denfield, G. H., Walker, E. Y., Gatys, L. A., Tolias, A. S., Bethge, M., and Ecker, A. S. Deep convolutional models improve predictions of macaque V1 responses to natural images. *PLoS Comput. Biol.*, 15 (4):e1006897, Apr. 2019.
- Dapello, J., Marques, T., Schrimpf, M., Geiger, F., Cox, D. D., and DiCarlo, J. J. Simulating a primary visual cortex at the front of CNNs improves robustness to image perturbations. June 2020.
- 4. Ge, X., Zhang, K., Gribizis, A., Hamodi, A. S., Sabino, A. M., and Crair, M. C. Retinal waves prime visual motion detection by simulating future optic flow. *Science*, 373(6553), July 2021.
- 5. Jones, J. P. and Palmer, L. A. An evaluation of the two-dimensional Gabor filter model of simple receptive fields in cat striate cortex. *J. Neurophysiol.*, 58(6):1233–1258, Dec. 1987.
- 6. Kim, J., Song, M., Jang, J., and Paik, S.-B. Spontaneous Retinal Waves Can Generate Long-Range Horizontal Connectivity in Visual Cortex. *J. Neurosci.*, 40(34):6584–6599, Aug. 2020.
- Kong, N. C. L., Margalit, E., Gardner, J. L., and Norcia, A. M. Increasing neural network robustness improves match to macaque V1 eigenspectrum, spatial frequency preference and predictivity. *PLoS Comput. Biol.*, 18 (1):e1009739, Jan. 2022.
- Majaj, N. J., Hong, H., Solomon, E. A., and DiCarlo, J. J. Simple Learned Weighted Sums of Inferior Temporal Neuronal Firing Rates Accurately Predict Human Core Object Recognition Performance. *Journal of Neuroscience*, 35(39):13402–13418, 2015.
- Mehrer, J., Spoerer, C. J., Jones, E. C., Kriegeskorte, N., and Kietzmann, T. C. An ecologically motivated image dataset for deep learning yields better models of human vision. *Proc. Natl. Acad. Sci. U. S. A.*, 118(8), Feb. 2021.
- 1289 10. Schrimpf, M., Kubilius, J., Lee, M. J., Ratan Murty, N. A., Ajemian, R., and DiCarlo, J. J. Integrative 1290 Benchmarking to Advance Neurally Mechanistic Models of Human Intelligence. *Neuron*, Sept. 2020.
- 1291 11. Stringer, C., Pachitariu, M., Steinmetz, N., Carandini, M., and Harris, K. D. High-dimensional geometry of 1292 population responses in visual cortex. *Nature*, 571(7765):361–365, July 2019.
- 12. Tanigawa, H., Lu, H. D., and Roe, A. W. Functional organization for color and orientation in macaque V4. *Nat. Neurosci.*, 13(12):1542–1548, Dec. 2010.