

# 1 Emergence of Visual Center-Periphery Spatial Organization in Deep 2 Convolutional Neural Networks

3  
4 Yalda Mohsenzadeh<sup>1,2,3</sup>, Caitlin Mullin<sup>4</sup>, Benjamin Lahner<sup>1</sup>, Aude Oliva<sup>1</sup>

5  
6 <sup>1</sup> Computer Science and Artificial Intelligence Laboratory, MIT, Cambridge, MA, USA

7 <sup>2</sup> Department of Computer Science, The University of Western Ontario, London, ON,  
8 Canada

9 <sup>3</sup> The Brain and Mind Institute, The University of Western Ontario, London, ON, Canada

10 <sup>4</sup> Department of Psychology, Center for Vision Research, York University, Toronto, ON,  
11 Canada

12

13

14 Corresponding author: Yalda Mohsenzadeh (ymohsenz@uwo.ca)

15

## 16 **Abstract**

17 Research at the intersection of computer vision and neuroscience has revealed  
18 hierarchical correspondence between layers of deep convolutional neural networks  
19 (DCNNs) and cascade of regions along human ventral visual cortex. Recently, studies  
20 have uncovered emergence of human interpretable concepts within DCNNs layers  
21 trained to identify visual objects and scenes. Here, we asked whether an artificial neural  
22 network (with convolutional structure) trained for visual categorization would  
23 demonstrate spatial correspondences with human brain regions showing  
24 central/peripheral biases. Using representational similarity analysis, we compared  
25 activations of convolutional layers of a DCNN trained for object and scene  
26 categorization with neural representations in human brain visual regions. Results reveal  
27 a brain-like topographical organization in the layers of the DCNN, such that activations  
28 of layer-units with central-bias were associated with brain regions with foveal tendencies  
29 (e.g. fusiform gyrus), and activations of layer-units with selectivity for image  
30 backgrounds were associated with cortical regions showing peripheral preference (e.g.  
31 parahippocampal cortex). The emergence of a categorical topographical  
32 correspondence between DCNNs and brain regions suggests these models are a good  
33 approximation of the perceptual representation generated by biological neural networks.

34

35 **Keywords:** central peripheral biases; deep convolutional neural networks;  
36 representational similarity analysis; fMRI; topographical maps

## 37 **Introduction**

38 Cortical regions along the ventral visual stream of the human brain (extending from  
39 occipital to temporal lobe) have been shown to preferentially activate to specific image

40 categories<sup>1</sup>. For instance, while the fusiform gyrus shows specialization for faces<sup>2</sup>, the  
41 parahippocampal cortex (PHC) is more selective to spatial layout, places<sup>3,4</sup> and large-  
42 size objects<sup>5,6</sup>. In characterizing the functional properties of these regions, Levy and  
43 colleagues (2001) discovered distinct topographical response patterns, such that face  
44 selective regions of the fusiform gyrus showed a strong preference for central visual  
45 field, while the building selective regions of PHC exhibited a peripheral selectivity bias to  
46 images of scene and large spaces<sup>7</sup>. Thus, while these regions show categorical  
47 selectivity to scenes or faces, their response patterns are strongest when their preferred  
48 category is presented in a topographically favorable location in the visual field. More  
49 specifically, the face selective voxels in the fusiform gyrus have a stronger response  
50 when faces are presented centrally; whereas scene-selective voxels show stronger  
51 activity to space features in the periphery<sup>7-13</sup>.

52

53 These topographical preferences raise questions regarding the origin of this functional  
54 organizing principles: does the way we look at faces and scenes in our natural visual  
55 world account for this bias? We most often fixate on faces bringing face-related  
56 information into our central, high acuity fovea to extract subtle visual features like facial  
57 expressions<sup>14-16</sup>. Places, on the other hand, are used for navigation, extending all  
58 around the visual field, thus we more readily perceive them with the peripheral visual  
59 information<sup>10-13</sup>.

60

61 Recently, a class of computational models, termed deep convolutional neural networks  
62 (DCNNs), inspired by the hierarchical architectures of ventral visual streams  
63 demonstrated striking similarities with the cascade of processing stages in the human  
64 visual system<sup>17-25</sup>. In particular, it has been shown that internal representations of these  
65 models are hierarchically similar to neural representations in early visual cortex<sup>26</sup>, mid-  
66 level (area v4), and high-level (area IT) cortical regions along ventral stream<sup>23,27</sup> in  
67 primates and to functional magnetic resonance imaging (fMRI) and  
68 magnetoencephalography measurements in humans<sup>19,22</sup>. Recent efforts to look into the  
69 features learned by the artificial units of DCNNs have revealed the emergence of  
70 human interpretable concepts<sup>28,29</sup>. For example, Bau et al. (2017) showed while units in  
71 the earlier layers of the network learn patterns of edges, curves and texture, depending  
72 on the categorization task (e.g. object or scene categorization), units in the subsequent  
73 layers show selectivity to shapes and object parts or whole objects and spatial layout  
74 patterns that differentiate scenes<sup>28</sup>. Furthermore, they discovered that the networks  
75 trained on scene or object categorizations spontaneously learned concepts like face,  
76 people or body parts, that they never were trained on them explicitly. This work points to  
77 DCNNs as a useful model of the human visual system and motivates broader  
78 examination of the correspondence between human brain and layered-models.

79

80 These similarities motivated increasing applications of these models in hypothesis-  
81 testing of brain computations in vision<sup>24,30–33</sup>. In the current study, we asked whether  
82 these simplified artificial networks might show a topographical organization similar to  
83 human visual system raised from the statistics of our natural visual world. To test this  
84 hypothesis, we compared the representations of units in a deep neural network trained  
85 on both object and scene categorization<sup>34</sup> (Hybrid-CNN) with representations from  
86 several category-selective areas of the visual hierarchy in the human brain. Given that  
87 this DCNN is trained on natural images representing the statistical distribution of visual  
88 features in the world, with a bias during learning similar to human visual experience (i.e.  
89 most faces are image-centered), we would expect activations of spatial selective nodes  
90 in DCNNs demonstrate category-specific topographical correspondences with human  
91 brain visual regions. Indeed, here we show that model units with central selectivity show  
92 stronger representational similarity to visual brain regions with a strong central-bias (e.g.  
93 fusiform gyrus), while model units with selectivity for image background are highly  
94 correlated with brain regions with peripheral bias (e.g. PHC).

## 95 **Results**

96 To investigate the topographical representations of DCNNs, we probed an 8-layer deep  
97 neural network model (AlexNet<sup>35</sup>, Figure 1A) with high performance in categorizing  
98 scene and object images<sup>34</sup>. The network architecture consists of five convolutional  
99 layers followed by three fully connected layers (see Figure 1A). The model, termed  
100 Hybrid-CNN<sup>34</sup>, is well suited for the purposes of examining topographical  
101 representations as it is trained on both ImageNet, an object-centered dataset<sup>36</sup> and  
102 Places, a scene-centered dataset<sup>34</sup> making it proficient in both object and scene  
103 recognition. We fed this network a stimulus set consisting of 156 natural images of  
104 faces, animates bodies (animals and people), faces, objects and scenes. These 156  
105 images were not used in training of the Hybrid-CNN.

106  
107 To compare the topographical representations of the Hybrid-CNN to those from the  
108 human visual system, we collected fMRI data while participants (N = 15) viewed the 156  
109 images in a rapid event related design and performed a vigilance task (detecting color  
110 changes in the fixation cross). The fMRI data of this study has been published in  
111 Mohsenzadeh et al. (2019)<sup>37</sup>.

112  
113 We employed multivariate pattern analysis to resolve neural representations in  
114 functional MRI data<sup>38–41</sup>. To probe the topographical architecture of the visual system,  
115 we targeted four regions of interest along ventral stream which represent a range of  
116 feature and category selectivity, namely; early visual cortex (EVC), fusiform gyrus  
117 (Fusiform), inferior temporal area (IT), and parahippocampal cortex (PHC). These

118 regions were defined anatomically according to Wang et al. (2014)<sup>42</sup> and Tzourio-  
119 Mazoyer et al. (2002)<sup>43</sup>.

120 A useful tool for comparing data from different modalities has been representational  
121 similarity analysis<sup>18,37,39,40,44–49</sup>. This technique can quantitatively relate brain-activity  
122 measurements, such as fMRI, with computational models, such as DCNNs, by  
123 abstracting the data into representational space using matrices of pairwise similarity (a  
124 representational dissimilarity matrix- RDM). Once in the same space, we can measure  
125 the consistency of information between these two systems.

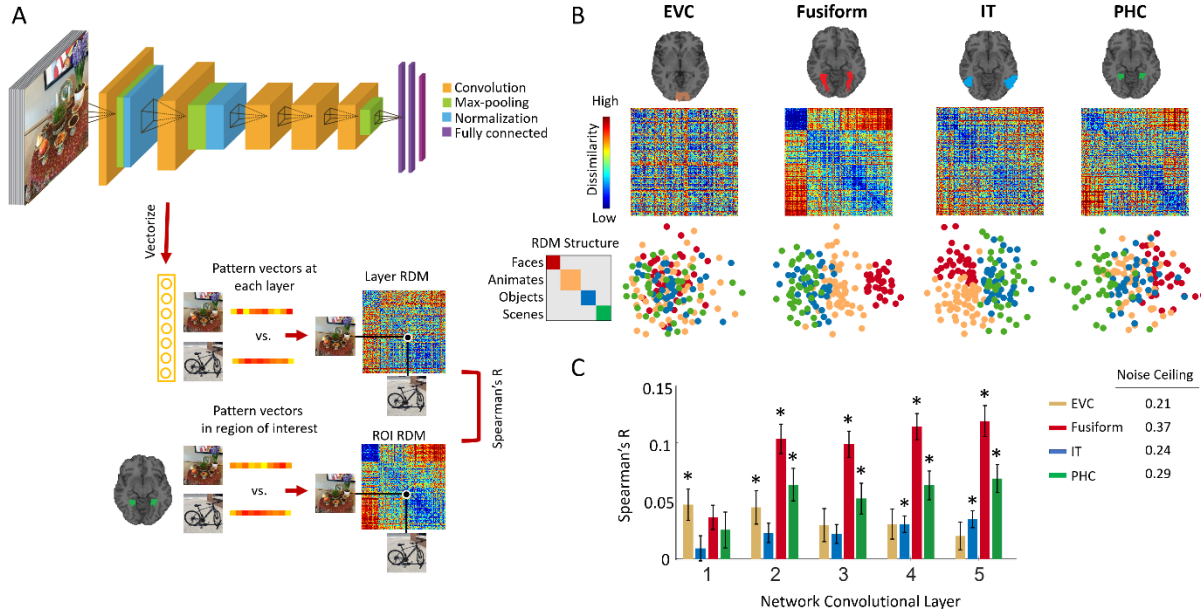
126  
127 Thus, for the DCNN we abstracted the data into the representational space<sup>18,19,22</sup> by  
128 extracting the layer activations, vectorized them and creating a representational  
129 dissimilarity matrix (RDM) of pairwise distances (1-Pearson corr) for each convolutional  
130 layer (1-5) as illustrated in Figure 1A. For each participant's fMRI data, stimulus-specific  
131 voxel responses in each ROI were extracted, noise normalized<sup>48</sup>, arranged into pattern  
132 vectors and then the pairwise distances (1-Pearson corr) were computed and populated  
133 a 156 x 156 distance matrix, also known as representational dissimilarity matrix (RDM).  
134 In this way, the data from each source now exist in a common space.

135  
136 Figure 1B depicts the subject-averaged ROI RDMs and their 2D multidimensional  
137 scaling visualizations. In line with previous literature<sup>2–4,50</sup>, neural representation in EVC  
138 depicts a random pattern consistent with low level visual feature processing in this area,  
139 fusiform dominantly clusters faces, IT shows a clear animate/inanimate distinction, and  
140 PHC discriminates scene images from others.

141  
142 **Hierarchical correspondences between layers of Hybrid-CNN and brain regions of**  
143 **interest along the ventral visual pathway**

144  
145 Here, we first compare the hierarchical correspondences between layers of Hybrid-CNN  
146 and the ventral stream regions of interests. For this, we compute the Spearman's rho  
147 correlations between the network layer RDMs and the individual's ROI RDMs. Figure  
148 1C shows comparison of fMRI representations in EVC, Fusiform, IT and PHC with the  
149 network layer-specific representations. As depicted, earlier layers of the network show  
150 significant correlation with EVC (all stats in Figure 2C, N=15; P<0.05, Bonferoni  
151 corrected). Fusiform gyrus shows progressively higher correlations along the layers of  
152 the network. IT demonstrates significant correlation with later layers (layer 4-5),  
153 consistent with a high-level object category processing in IT (for reviews see<sup>50–52</sup>).  
154 Finally, PHC representation is significantly correlated with mid to later layers of the  
155 network (layer 2 and 5) confirming low to high level scene semantic processes in these  
156 layers. This illustrates the distinct low to high level features of faces across the layers of  
157 this hybrid network.

158



159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

182

183

## Figure 1. Hierarchical correspondences between layers of DCNN and brain

### regions of interest along ventral visual pathway. (A) For each image, the activation

of units in each of the 5 convolutional layers are vectorized. RDM representation for

each layer is created by computing the pairwise distance of these image specific vector

patterns (1-Pearson Corr). Then fMRI RDM representations in EVC, Fusiform, IT and

PHC areas are compared with the RDM representations of each convolutional layer of

Hybrid-CNN by computing Spearman's correlations. (B) Neural representations along

ventral visual pathway. RDM matrices, and 2D multidimensional scaling visualization of

stimuli depicted for early visual cortex (EVC), fusiform gyrus (Fusiform), inferior

temporal cortex (IT) and parahippocampal cortex (PHC). (C) The correlation values for

brain ROIs and layers of DCNN are depicted with bar plots. The error bars indicate the

standard error of the mean and the stars above each bar indicates significant correlation

above zero (N = 15, P < 0.05, Bonferroni-corrected). The noise ceiling for each brain

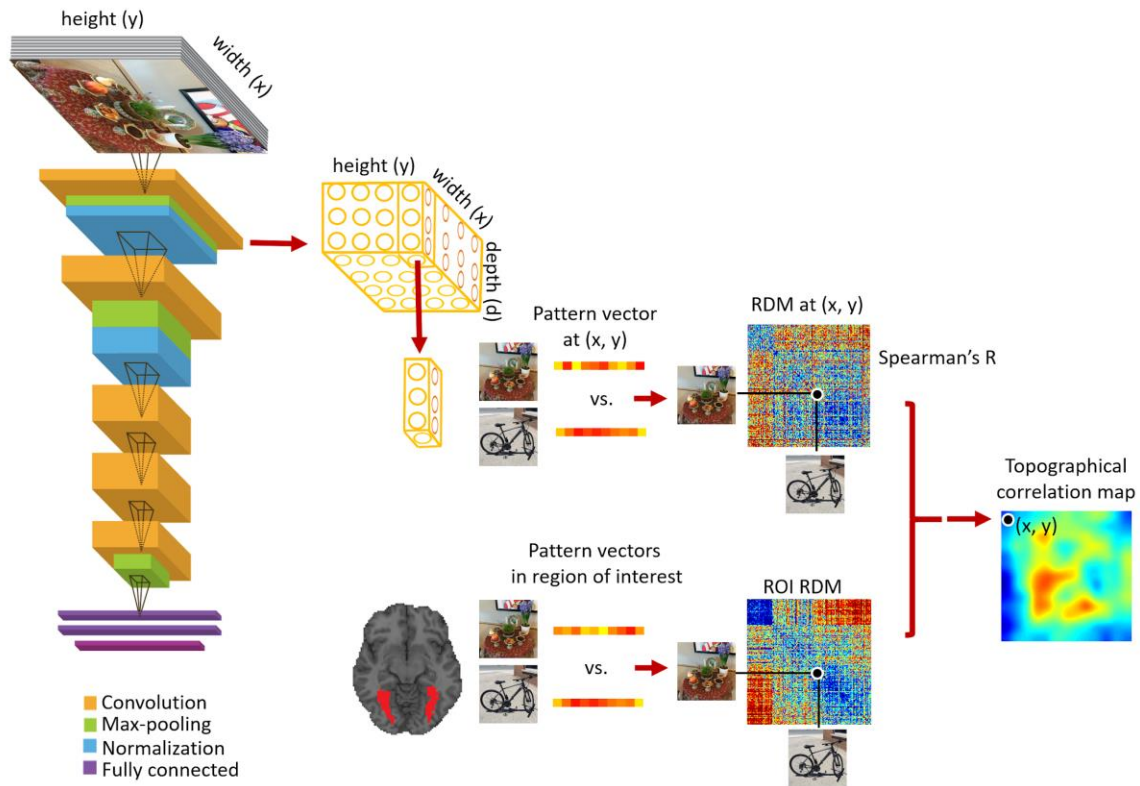
area is reported on the right side of the panel. The pictures used in this figure are not

examples of the stimulus set due to copyright.

184 convolutional layers of deep nets are spatially related with the 2D image space. Here,  
185 we took advantage of this correspondence and extracted the pattern vectors within each  
186 layer for pattern vectors associated with (x,y) positions in this 2D space . We used the  
187 image specific activation pattern and created RDM matrices at each (x, y) position  
188 within each layer. Then compared these RDM matrices with ROI RDMs by computing  
189 Spearman's rho correlations resulting in 2D correlation maps as illustrated in Figure 2.  
190 We call these brain-DCNN maps *topographical correlation maps*.

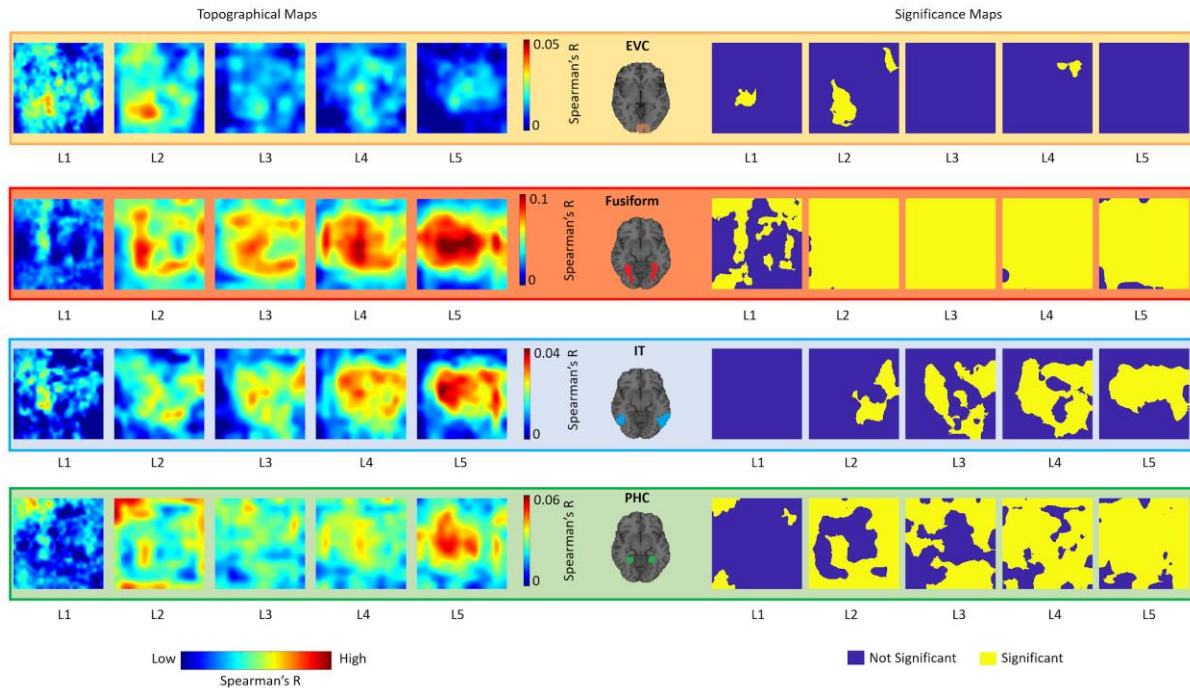
191  
192 Figure 3 illustrates the topographical correspondence between the five convolutional  
193 layers of the Hybrid-CNN model and the four fMRI ROIs and show the significance map  
194 corresponding to each topographical map (N=15; two-sided sign permutation tests,  
195 cluster defining threshold  $P < 0.01$ , and corrected significance level  $P < 0.05$ ). Results are  
196 aligned with the hypothesis that the artificial network with convolutional architecture  
197 demonstrates spatial representations of center-periphery image bias highly correlated  
198 with brain regions showing similar bias: as expected, EVC, which is not category-  
199 dependent, shows a randomly distributed significant topographical pattern in the first  
200 two layers. This is consistent with previous studies showing earlier layers of network  
201 share similar representation with EVC. Furthermore, the depicted patterns indicate  
202 these low-level features are scattered over the image. Correlation maps of fusiform and  
203 layers of the network demonstrate strong center-selective patterns, consistent with the  
204 foveal-bias representations in fusiform gyrus.

205 Correlation maps in IT show a dispersed pattern in mid-level layers which becomes  
206 more centralized over the layers, illustrating a mid to high-level representation  
207 transformation. Lastly, topographical maps of PHC transforms from a clear  
208 background/surrounding organization to a distributed pattern from layer 1 to 5. This  
209 suggests a transformation of low-level periphery features to scene semantics across the  
210 layers.



211

212 **Figure 2. Creating topographical correlation maps.** We extract the 3D activation  
213 patterns from the network convolutional layers. The first 2 Dimensions have a spatial  
214 relation with the image space (width and height). At each (x,y) position in feature maps,  
215 we extract a pattern vector with the length equivalent to the depth and construct the  
216 RDM matrix from the neural network activity patterns at each (x,y) location. Comparison  
217 of these RDM matrices with a brain ROI RDM results in a 2D correlation map which we  
218 then up-sample it to the image size (topographical map). The pictures used in this figure  
219 are not examples of the stimulus set due to copyright.



220  
221 **Figure 3. Topographical correspondence between convolutional layers of DCNNs**  
222 **and human ventral visual regions.** For each brain-model mapping (EVC, Fusiform,  
223 IT, PHC), the first five maps show the correlational topographical maps between each  
224 convolutional layer and the brain ROI; the second five maps show the corresponding  
225 significance maps (two-sided sign permutation tests, cluster defining threshold  $P < 0.01$ ,  
226 and corrected significance level  $P < 0.05$ ). The topographical correlation maps in this  
227 figure are computed following the method depicted in Figure 2. For detailed description  
228 of RDM computations and correlations please see the Method section.

## 229 Discussion

### 230 Summary

231 Hierarchical correspondences have been established between primate ventral visual  
232 pathway and layers of DCNNs<sup>19,22,23</sup>. In the current study, using representational  
233 similarity analysis, we first replicated the previous findings showing hierarchical  
234 correspondence from low to high level ventral visual areas and layers of the deep nets  
235 (Figure 1C). Importantly, we demonstrated for the first time a topographical  
236 correspondence (central/periphery biases) between ventral brain regions and unit  
237 activations of the Hybrid-CNN (Figure 3). Our results revealed foveally biased fusiform  
238 highly correlated with unit activations of the network with the center selective visual field  
239 and peripherally biased PHC strongly correlated with unit activations of the network with  
240 periphery selective receptive fields.

241

### 242 Topographical similarities of ventral visual stream and Hybrid-CNN



243 Visual neuroscience has focused on functional localization in establishing a set of brain  
244 areas along ventral stream that specialize in visual functions such as scene, object,  
245 face, body recognition<sup>2-4,53-56</sup>. On the other hand, the fields of computer vision and  
246 artificial intelligence pursue a different goal, function optimization. While artificial  
247 networks are trained on simple optimization objectives, e.g. object categorization  
248 through minimizing classification errors, they have shown emergence of similar  
249 characteristics akin to the brain areas along ventral visual pathway. Recent works have  
250 suggested that neural networks trained on natural images with categorization  
251 objectives, learn visual features along a hierarchy which matches human visual system  
252 in space and time<sup>19,22,23,28</sup>. For instance, deep nets trained on scene categorization  
253 showed the spontaneous emergence of object/face representations in higher layers of  
254 the network<sup>28,29,34</sup>. Here, we further showed that a network with convolutional constraint  
255 trained on object and scene categorization naturally demonstrates topographical  
256 correspondences with brain areas in the ventral stream showing periphery/central  
257 biases. This raises the question whether these characteristics are developed in our  
258 brain due to the position of these visual features in our visual experience (faces at the  
259 center of visual field while scenes are at the background of our vision). While  
260 neuroscientist addressed some aspects of this question<sup>7-13</sup>, these topographical  
261 correspondences motivate future experiments with deep neural networks to shed light  
262 on the computational principles behind these properties. Nevertheless, our findings in  
263 the current study revealing more similarities between deep nets and ventral pathway,  
264 supports two main hypotheses in vision neuroscience: first the human visual pathway  
265 optimizes the cost function of visual recognition and second the characteristics of neural  
266 tuning, internal neural representations and brain area functions along this pathway are  
267 most likely the result of this cost function optimization<sup>57</sup>.

268

### 269 **Broader implications of brain and CNN similarities**

270 Our data reveals emergence of previously unknown similarities between the visual brain  
271 and a DCNN trained on object and scene categorization. Our results revealed a  
272 topographical correspondence between unit activations of a convolutional neural  
273 network and categorical selective brain regions. The convolutional architecture of the  
274 network limits us to investigate whether the topographical bias in the brain (and the  
275 model) is due to the statistics of the visual input or the neurons (network units) learn  
276 these characteristics. Future computational models with un-tied weights are motivated  
277 to investigate this question. We predict training such network with real world  
278 categorization tasks (objects and scenes) shapes the representation of the network  
279 units. The emergence of hierarchical as well as topographical similarities between the  
280 network and the brain implies that these characteristics of the ventral stream are most  
281 likely the result of computational objective of this pathway being visual categorization.

282 Thus, further understanding the function, computations and connectivity architectures of  
283 visual cortex can potentially guide and give insight into brain-inspired models of vision.

## 284 **Methods**

285 The fMRI data of this study has been published previously<sup>37</sup>. Here, we briefly describe  
286 the necessary information related to experiment design and data acquisition and  
287 analysis.

## 288 **Participants**

289 Fifteen healthy individuals (right handed; 9 females, age: mean  $\pm$  SD = 27.87  $\pm$  5.17  
290 years) with normal or corrected to normal participated in this study. All participants  
291 signed an informed written consent form and were compensated for their time. This  
292 study was conducted in accordance with the Declaration of Helsinki and approved by  
293 the local ethics committee (Institutional Review Board of the Massachusetts Institute of  
294 Technology).

## 295 **Experiment design, task and stimulus set**

296 Participants viewed sequence of images presented for 0.5s with 2.5s inter stimulus  
297 intervals in MRI scanner. Image trials were randomly mixed with 25% null trials during  
298 which fixation cross changed color for 100 msec and participants reported the color  
299 change by a button press. The stimulus set in our study consisted of 156 natural images  
300 of faces, animates bodies (animals and people), objects and scenes. The images were  
301 presented at the center of the screen at 6° visual angle. We acquired fMRI data in two  
302 separate sessions (each session included 5-8 runs) and images were presented once  
303 per run in random order.

## 305 **fMRI data acquisition and analysis**

306 The fMRI data were acquired using a 3T Siemens Trio scanner with 32-channel  
307 phased-array head coil at the Athinoula A. Martinos Imaging Center at the MIT  
308 McGovern Institute for Brain Research. In each session, structural images were  
309 acquired using a standard T1-weighted sequence (176 sagittal slices, FOV = 256 mm<sup>2</sup>,  
310 TR = 2530 ms, TE = 2.34 ms, flip angle = 9°) and then 5-8 runs of 305 volumes of  
311 functional data were acquired for each participant (11-15 runs across the two sessions).  
312 For the functional data, gradient-echo EPI sequence was used (TR = 2000 ms, TE = 29  
313 ms, flip angle = 90°, FOV read = 200 mm, FOV phase = 100%, bandwidth 2368 Hz/Px,  
314 gap = 20%, resolution = 3.1 mm isotropic, slices = 33, ascending interleaved  
315 acquisition).

316

317 We used SPM software to preprocess fMRI data. Functional Data of each participant  
318 were slice-time corrected, realigned and co-registered to the first session T1 structural  
319 scan, and normalized to the standard MNI space. For multivariate analysis, we used  
320 unsmoothed data. We estimated the fMRI responses to the 156 image conditions using  
321 a general linear model (GLM). We modeled the events (image conditions and nulls) with  
322 event onsets and impulse response function (duration of zero), furthermore, we included  
323 motion and run regressors in the GLM. The defined regressors were convolved with the  
324 hemodynamic response function. We then estimated the beta-values for each image  
325 condition and also the residuals from the first-level GLM. The residuals were used as an  
326 estimation of noise structure which then was employed for multivariate noise  
327 normalization (Walther et al., 2016).

328  
329 In the current study, we defined four regions of interest (ROIs) along the ventral stream,  
330 early visual cortex (EVC), fusiform gyrus (Fusiform), inferior temporal cortex (IT), and  
331 parahippocampal cortex (PHC). All defined anatomically based on Wang et al. (2015)<sup>42</sup>  
332 (EVC and PHC) and Tzourio-Mazoyer et al (2002)<sup>43</sup> (IT and Fusiform).

### 333 **Deep convolutional neural network architecture and training**

334 We used a deep convolutional neural network (DCNN) with architecture similar to  
335 Krizhevsky et al. (2012)<sup>35</sup>. This DCNN was trained both on object categories (ImageNet  
336 dataset) and scene categories (Places dataset) and called Hybrid-CNN<sup>34</sup>. The rationale  
337 for this choice is that our stimulus set consists of both objects and scene images and  
338 this network is trained on both ILSVRC 2012<sup>36</sup> and Places dataset<sup>34</sup>. The network  
339 architecture includes five convolutional layers followed by three fully connected layers.  
340 For the purpose of topography comparison, we used the convolutional layers here.

### 341 **Representational similarity analysis of fMRI and DCNN units**

342 We used representational similarity analysis to map fMRI responses and CNN units'  
343 activations into a common space<sup>19,22,39,41</sup>. In this framework, the assumption is that the  
344 pairwise relationships of images are similar in the brain and the model. These pairwise  
345 relationships between 156 image-specific model/brain responses are measured in terms  
346 of dissimilarity distances (here we used 1- Pearson correlation) and summarized in a  
347 156 X 156 representational dissimilarity matrix (RDM). We created subject-specific fMRI  
348 RDMs per region of interest (EVC, Fusiform, IT, and PHC) and model RDMs per layer  
349 or per spatial position within a convolutional layer.

350  
351 In detail, in each fMRI ROI and for each of 156 image conditions we extracted the beta-  
352 value activation patterns, arranged them into vector patterns, normalized them based on  
353 the covariance of the estimated noise<sup>48</sup> and then computed the pairwise dissimilarity of

354 these 156 vector patterns by calculating 1 minus Pearson correlations. This yielded a  
355 156x156 representational dissimilarity matrix (RDM) for each individual and ROI.

### 356 **Brain and DCNN topographical maps**

357 In this study, we used a deep neural network by Zhou et al. (2014) called Hybrid-CNN<sup>34</sup>.  
358 We chose this network as it was trained both on object and scene categories and thus a  
359 suitable model to explain categories in our fMRI data. The network architecture consists  
360 of five convolutional and three fully connected layers similar to Krizhevsky et al.  
361 (2012)<sup>35</sup>. We extracted 3D activation maps from convolutional layers of this network for  
362 each image in our stimulus set. The first two dimensions are spatially related to the  
363 image space (width and height). As depicted in Figure 2, at each (x, y) position on the  
364 activation map, we extracted a pattern vector with the length of the depth. Then we  
365 constructed the RDM matrices from pairwise distances of these image-specific pattern  
366 vectors (1-Pearson corr). Next, we compared these neural network RDM  
367 representations with brain ROI RDMs simply by calculating the Spearman's rho  
368 between them. This process results in a 2D correlation map for each convolutional layer  
369 which we then up-sample it to the input image dimension (topographical maps).

### 370 **Multidimensional scaling (MDS) visualization**

371 Multidimensional scaling<sup>58</sup> is an unsupervised approach to visualize the similarity  
372 relationships of conditions represented in a complex distance matrix, such that similar  
373 image conditions are visualized clustered together and different ones are depicted  
374 apart. Here, our distance matrices are 156 x 156 fMRI RDMs capturing the relations of  
375 neural patterns corresponding to 156 image stimuli in four regions of interest. The first 2  
376 dimensions of MDS are used to visualize these images organized in 5 categories  
377 (Figure 2).

### 378 **Statistical tests**

379 For statistical tests we used nonparametric methods with no assumption on the data  
380 distributions<sup>59,60</sup>. Permutation-based cluster-size inference was used for statistical  
381 inference on the topographical maps (1000 permutations, 0.01 cluster definition  
382 threshold and 0.05 cluster size threshold) with null hypothesis of zero correlation.

### 383 **Acknowledgments**

384 Funding from the Vannevar Bush Faculty Fellowship program by the ONR (N00014-16-  
385 1-3116) and NSF award 1532591 in Neural and Cognitive Systems (to A.O.). Study  
386 conducted at the Athinoula A. Martinos Imaging Center, MIBR, MIT.

387 **Data availability**

388 The data and analysis tools are available from the corresponding author upon request.

389 **Author contributions**

390 Conceptualization, Y.M., A.O.; investigation, Y.M., A.O.; methodology, Y.M.; data  
391 collection, Y.M., C.M.; data analysis, Y.M., B.L.; visualization, Y.M;  
392 writing and editing, Y.M., C.M., A.O.;  
393

394 **Competing Interests**

395 The authors declare no competing interests.

396 **References**

- 397 1. Grill-Spector, K. & Weiner, K. S. The functional architecture of the ventral temporal cortex  
398 and its role in categorization. *Nature Reviews Neuroscience* **15**, 536–548 (2014).
- 399 2. Kanwisher, N., McDermott, J. & Chun, M. M. The Fusiform Face Area: A Module in Human  
400 Extrastriate Cortex Specialized for Face Perception. *The Journal of Neuroscience* **17**, 4302–  
401 4311 (1997).
- 402 3. Epstein, R. & Kanwisher, N. A cortical representation of the local visual environment. *Nature*  
403 **392**, 598–601 (1998).
- 404 4. Epstein, R., Harris, A., Stanley, D. & Kanwisher, N. The Parahippocampal Place Area:  
405 Recognition, Navigation, or Encoding? *Neuron* **23**, 115–125 (1999).
- 406 5. Konkle, T. & Caramazza, A. Tripartite Organization of the Ventral Stream by Animacy and  
407 Object Size. *Journal of Neuroscience* **33**, 10235–10242 (2013).
- 408 6. Konkle, T. & Oliva, A. A Real-World Size Organization of Object Responses in  
409 Occipitotemporal Cortex. *Neuron* **74**, 1114–1124 (2012).
- 410 7. Levy, I., Hasson, U., Avidan, G., Hendler, T. & Malach, R. Center–periphery organization of  
411 human object areas. *Nature Neuroscience* **4**, 533–539 (2001).

- 412 8. Hasson, U., Levy, I., Behrmann, M., Hendler, T. & Malach, R. Eccentricity Bias as an  
413 Organizing Principle for Human High-Order Object Areas. *Neuron* **34**, 479–490 (2002).
- 414 9. Malach, R., Levy, I. & Hasson, U. The topography of high-order human object areas. *Trends*  
415 *in Cognitive Sciences* **6**, 176–184 (2002).
- 416 10. Arcaro, M. J., McMains, S. A., Singer, B. D. & Kastner, S. Retinotopic Organization of  
417 Human Ventral Visual Cortex. *Journal of Neuroscience* **29**, 10638–10652 (2009).
- 418 11. Larson, A. M. & Loschky, L. C. The contributions of central versus peripheral vision to scene  
419 gist recognition. *Journal of Vision* **9**, 6–6 (2009).
- 420 12. Thibaut, M., Tran, T. H. C., Szaffarczyk, S. & Boucart, M. The contribution of central and  
421 peripheral vision in scene categorization: A study on people with central vision loss. *Vision*  
422 *Research* **98**, 46–53 (2014).
- 423 13. Baldassano, C., Fei-Fei, L. & Beck, D. M. Pinpointing the peripheral bias in neural scene-  
424 processing networks during natural viewing. *Journal of Vision* **16**, 9 (2016).
- 425 14. Peters, C. Direction of Attention Perception for Conversation Initiation in Virtual  
426 Environments. in *Intelligent Virtual Agents* (eds. Panayiotopoulos, T. et al.) **3661**, 215–228  
427 (Springer Berlin Heidelberg, 2005).
- 428 15. Treue, S. Visual attention: the where, what, how and why of saliency. *Current Opinion in*  
429 *Neurobiology* **13**, 428–432 (2003).
- 430 16. Yarbus, A. L. *Eye Movements and Vision*. (1967).
- 431 17. Bonner, M. F. & Epstein, R. A. Computational mechanisms underlying cortical responses to  
432 the affordance properties of visual scenes. *PLOS Computational Biology* **14**, e1006111  
433 (2018).
- 434 18. Cichy, R. M., Khosla, A., Pantazis, D. & Oliva, A. Dynamics of scene representations in the  
435 human brain revealed by magnetoencephalography and deep neural networks. *NeuroImage*  
436 **153**, 346–358 (2017).

- 437 19. Cichy, R. M., Khosla, A., Pantazis, D., Torralba, A. & Oliva, A. Comparison of deep neural  
438 networks to spatio-temporal cortical dynamics of human visual object recognition reveals  
439 hierarchical correspondence. *Scientific Reports* **6**, (2016).
- 440 20. Eickenberg, M., Gramfort, A., Varoquaux, G. & Thirion, B. Seeing it all: Convolutional  
441 network layers map the function of the human visual system. *NeuroImage* **152**, 184–194  
442 (2017).
- 443 21. Guclu, U. & van Gerven, M. A. J. Deep Neural Networks Reveal a Gradient in the  
444 Complexity of Neural Representations across the Ventral Stream. *Journal of Neuroscience*  
445 **35**, 10005–10014 (2015).
- 446 22. Khaligh-Razavi, S.-M. & Kriegeskorte, N. Deep Supervised, but Not Unsupervised, Models  
447 May Explain IT Cortical Representation. *PLoS Computational Biology* **10**, e1003915 (2014).
- 448 23. Yamins, D. L. K. *et al.* Performance-optimized hierarchical models predict neural responses  
449 in higher visual cortex. *Proceedings of the National Academy of Sciences* **111**, 8619–8624  
450 (2014).
- 451 24. Rajaei, K., Mohsenzadeh, Y., Ebrahimpour, R. & Khaligh-Razavi, S.-M. Beyond core object  
452 recognition: Recurrent processes account for object recognition under occlusion. *PLOS*  
453 *Computational Biology* **30** (2019).
- 454 25. Cox, D. D. & Dean, T. Neural Networks and Neuroscience-Inspired Computer Vision.  
455 *Current Biology* **24**, R921–R929 (2014).
- 456 26. Cadena, S. A. *et al.* Deep convolutional models improve predictions of macaque V1  
457 responses to natural images. 27
- 458 27. Yamins, D. L., Hong, H., Cadieu, C. & DiCarlo, J. J. Hierarchical Modular Optimization of  
459 Convolutional Networks Achieves Representations Similar to Macaque IT and Human  
460 Ventral Stream. *NIPS* **9** (2013).
- 461 28. Bau, D., Zhou, B., Khosla, A., Oliva, A. & Torralba, A. Network Dissection: Quantifying  
462 Interpretability of Deep Visual Representations. in *2017 IEEE Conference on Computer*

- 463        *Vision and Pattern Recognition (CVPR)* 3319–3327 (IEEE, 2017).  
464        doi:10.1109/CVPR.2017.354
- 465    29. Zhou, B., Bau, D., Oliva, A. & Torralba, A. Interpreting Deep Visual Representations via  
466        Network Dissection. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 1–1  
467        (2018). doi:10.1109/TPAMI.2018.2858759
- 468    30. Walker, E. Y. *et al.* Inception in visual cortex: in vivo-silico loops reveal most exciting  
469        images. *bioRxiv* (2018). doi:10.1101/506956
- 470    31. Bashivan, P., Kar, K. & DiCarlo, J. J. Neural population control via deep image synthesis.  
471        *Science* **364**, eaav9436 (2019).
- 472    32. Jaegle, A. *et al.* A neural correlate of image memorability. *bioRxiv* (2019).  
473        doi:10.1101/535468
- 474    33. Ponce, C. R. *et al.* Evolving Images for Visual Neurons Using a Deep Generative Network  
475        Reveals Coding Principles and Neuronal Preferences. *Cell* **177**, 999-1009.e10 (2019).
- 476    34. Zhou, B., Lapedriza, A., Xiao, J., Torralba, A. & Oliva, A. Learning Deep Features for Scene  
477        Recognition using Places Database. *NIPS* 9 (2014).
- 478    35. Krizhevsky, A., Sutskever, I. & Hinton, G. E. ImageNet classification with deep convolutional  
479        neural networks. *NIPS* **60**, 84–90 (2012).
- 480    36. Deng, J. *et al.* ImageNet: A Large-Scale Hierarchical Image Database. *IEEE Conference on*  
481        *Computer Vision and Pattern Recognition* 8 (2009). doi:10.1109/CVPR.2009.5206848
- 482    37. Mohsenzadeh, Y., Mullin, C., Lahner, B., Cichy, R. & Oliva, A. Reliability and  
483        Generalizability of Similarity-Based Fusion of MEG and fMRI Data in Human Ventral and  
484        Dorsal Visual Streams. *Vision* **3**, 8 (2019).
- 485    38. Haynes, J.-D. & Rees, G. Decoding mental states from brain activity in humans. *Nature*  
486        *Reviews Neuroscience* **7**, 523–534 (2006).
- 487    39. Kriegeskorte, N. Representational similarity analysis – connecting the branches of systems  
488        neuroscience. *Frontiers in Systems Neuroscience* (2008). doi:10.3389/neuro.06.004.2008



- 489 40. Mur, M., Bandettini, P. A. & Kriegeskorte, N. Revealing representational content with  
490 pattern-information fMRI—an introductory guide. *Social Cognitive and Affective*  
491 *Neuroscience* **4**, 101–109 (2009).
- 492 41. Kriegeskorte, N. & Kievit, R. A. Representational geometry: integrating cognition,  
493 computation, and the brain. *Trends in Cognitive Sciences* **17**, 401–412 (2013).
- 494 42. Wang, L., Mruczek, R. E. B., Arcaro, M. J. & Kastner, S. Probabilistic Maps of Visual  
495 Topography in Human Cortex. *Cerebral Cortex* **25**, 3911–3931 (2015).
- 496 43. Tzourio-Mazoyer, N. *et al.* Automated Anatomical Labeling of Activations in SPM Using a  
497 Macroscopic Anatomical Parcellation of the MNI MRI Single-Subject Brain. *NeuroImage* **15**,  
498 273–289 (2002).
- 499 44. Cichy, R. M., Pantazis, D. & Oliva, A. Resolving human object recognition in space and  
500 time. *Nature Neuroscience* **17**, 455–462 (2014).
- 501 45. Cichy, R. M., Pantazis, D. & Oliva, A. Similarity-Based Fusion of MEG and fMRI Reveals  
502 Spatio-Temporal Dynamics in Human Cortex During Visual Object Recognition. *Cerebral*  
503 *Cortex* **26**, 3563–3579 (2016).
- 504 46. Mohsenzadeh, Y., Qin, S., Cichy, R. M. & Pantazis, D. Ultra-Rapid serial visual presentation  
505 reveals dynamics of feedforward and feedback processes in the ventral visual pathway.  
506 *eLife* **7**, 1–23 (2018).
- 507 47. Khaligh-Razavi, S.-M., Cichy, R. M., Pantazis, D. & Oliva, A. Tracking the Spatiotemporal  
508 Neural Dynamics of Real-world Object Size and Animacy in the Human Brain. *Journal of*  
509 *Cognitive Neuroscience* **30**, 1559–1576 (2018).
- 510 48. Walther, A. *et al.* Reliability of dissimilarity measures for multi-voxel pattern analysis.  
511 *NeuroImage* **137**, 188–200 (2016).
- 512 49. Pantazis, D. *et al.* Decoding the orientation of contrast edges from MEG evoked and  
513 induced responses. *NeuroImage* **180**, 267–279 (2018).

- 514 50. DiCarlo, J. J., Zoccolan, D. & Rust, N. C. How Does the Brain Solve Visual Object  
515 Recognition? *Neuron* **73**, 415–434 (2012).
- 516 51. Logothetis, N. K. & Sheinberg, D. L. Visual Object Recognition. 45 (1996).
- 517 52. Tanaka, K. Inferotemporal Cortex and Object Vision. 31 (1996).
- 518 53. Malach, R. *et al.* Object-related activity revealed by functional magnetic resonance imaging  
519 in human occipital cortex. *Proceedings of the National Academy of Sciences* **92**, 8135–8139  
520 (1995).
- 521 54. Downing, P. E. A Cortical Area Selective for Visual Processing of the Human Body. *Science*  
522 **293**, 2470–2473 (2001).
- 523 55. Grill-Spector, K. *et al.* Differential Processing of Objects under Various Viewing Conditions  
524 in the Human Lateral Occipital Complex. *Neuron* **24**, 187–203 (1999).
- 525 56. Dilks, D. D., Julian, J. B., Paunov, A. M. & Kanwisher, N. The Occipital Place Area Is  
526 Causally and Selectively Involved in Scene Perception. *Journal of Neuroscience* **33**, 1331–  
527 1336 (2013).
- 528 57. Marblestone, A. H., Wayne, G. & Kording, K. P. Toward an Integration of Deep Learning  
529 and Neuroscience. *Frontiers in Computational Neuroscience* **10**, (2016).
- 530 58. Shepard, R. N. Multidimensional Scaling, Tree-Fitting, and Clustering. *Science* **210**, 390–  
531 398 (1980).
- 532 59. Maris, E. & Oostenveld, R. Nonparametric statistical testing of EEG- and MEG-data. *Journal*  
533 *of Neuroscience Methods* **164**, 177–190 (2007).
- 534 60. Pantazis, D., Nichols, T. E., Baillet, S. & Leahy, R. M. A comparison of random field theory  
535 and permutation methods for the statistical analysis of MEG data. *NeuroImage* **25**, 383–394  
536 (2005).
- 537