

CATEGORY VS SHAPE IN DEEP NETWORKS & VISION

Orthogonal Representations of Object Shape and Category in Deep Convolutional Neural Networks and Human Visual Cortex

Astrid A. Zeman*, J. Brendan Ritchie, Stefania Bracci, Hans Op de Beeck

Department of Brain and Cognition, KULeuven

*corresponding author

Correspondence Email: astrid.zeman@kuleuven.be

Correspondence Address: Brain and Cognition Department, KULeuven.

Tiensestraat 102, Leuven 3000 Belgium.

CATEGORY VS SHAPE IN DEEP NETWORKS & VISION

1 Abstract

2 Deep Convolutional Neural Networks (CNNs) are gaining traction as the benchmark model of
3 visual object recognition, with performance now surpassing humans. While CNNs can accurately
4 assign one image to potentially thousands of categories, network performance could be the result
5 of layers that are tuned to represent the visual shape of objects, rather than object category, since
6 both are often confounded in natural images. Using two stimulus sets that explicitly dissociate
7 shape from category, we correlate these two types of information with each layer of multiple
8 CNNs. We also compare CNN output with fMRI activation along the human visual ventral
9 stream by correlating artificial with biological representations. We find that CNNs encode
10 category information independently from shape, peaking at the final fully connected layer in all
11 tested CNN architectures. Comparing CNNs with fMRI brain data, early visual cortex (V1) and
12 early layers of CNNs encode shape information. Anterior ventral temporal cortex encodes
13 category information, which correlates best with the final layer of CNNs. The interaction
14 between shape and category that is found along the human visual ventral pathway is echoed in
15 multiple deep networks. Our results suggest CNNs represent category information independently
16 from shape, much like the human visual system.

17

18 Keywords: deep learning, shape, object categorisation, Convolutional Neural Networks (CNNs),
19 fMRI

20

21

CATEGORY VS SHAPE IN DEEP NETWORKS & VISION

22 Introduction

23
24 In recent years, the performance of Deep Convolutional Neural Networks (CNNs) has
25 improved significantly, such that they are able to meet¹⁻³, and even surpass⁴ human performance
26 in classifying objects. In light of these impressive findings, these artificial networks are
27 increasingly compared to their biological counterparts, resulting in an accumulation of evidence
28 for their use as a benchmark model of visual object recognition^{5, 6}. For example, the internal
29 representations of CNNs show correspondence with human ventral temporal cortex (VTC) as
30 measured by fMRI, as well as with primate inferotemporal cortex (IT) measured using single cell
31 recordings⁷⁻¹². The correspondence between deep networks and neural representations along the
32 visual pathway has even allowed for accurate neural response prediction of single-cell recordings
33 in IT⁹ as well as fMRI¹³. Representational similarities have been further extended from the
34 spatial into the temporal domain, with results showing a corresponding ordering of processing
35 between CNNs and the human visual brain using MEG¹⁴. These accumulating findings showcase
36 the ability of CNNs to model neurons from single unit responses to entire populations, spanning
37 the multiple scales and dimensions used to study neural activity, and making CNNs one of the
38 best models to date for studying vision in the human and primate brain.

39 While these feats are impressive, it is unclear to what extent these results are easily
40 interpretable in terms of category representations. Object category information can often be
41 confounded with low-level visual features, such as colour, texture, and shape¹⁵. In this paper, we
42 highlight the significant interaction between shape and category that is known to occur in natural
43 images¹⁶ and address the possibility that these networks may distinguish between object
44 categories by relying upon visual features, such as shape, rather than high-level category

CATEGORY VS SHAPE IN DEEP NETWORKS & VISION

45 representations. Indeed, the shape similarity of objects has been capitalised on in the machine
46 learning field to improve performance¹⁷. CNNs are proficient at representing the perceived shape
47 of objects, as opposed to their physical shape¹⁸ and there are claims that CNNs rely heavily upon
48 shape information for classification¹⁹. Two-dimensional regular vs irregular shape
49 representations have been found in monkey IT, which are highly comparable to late layers of
50 CNNs¹². Furthermore, CNNs mimic a behavioural bias in humans known as the “shape-bias”,
51 which is the preference to categorise an object based on shape rather than colour²⁰. Given that
52 these networks are adept at representing object shape, it is possible they are taking advantage of
53 shape-based features, instead of category information, to classify object images.

54 Recent neuroimaging studies have begun to de-cofound category from visual features,
55 including shape, in order to investigate their interaction along the visual ventral pathway^{10, 16, 21,}
56 ²². VTC in humans is one of the main category-selective areas²³, distinguishing, for example,
57 between animate and inanimate objects^{24, 25}. To build up this category-related representation,
58 visual information is processed in a series of stages along the ventral visual pathway, from
59 primary visual cortex (area V1) through to VTC²³. In recent years, the exact role of VTC has
60 come under question, in particular whether this area encodes category-specific information, or
61 simply the low-level visual properties associated with category, such as colour, shape, size and
62 texture^{15, 26, 27}. Proklova, Kaiser & Peelen²² found that VTC encodes texture and outline
63 alongside category-specific information that is not present in earlier visual areas. Another higher
64 visual area, lateral occipitotemporal complex (LOTC), was found to encode category-associated
65 shape properties as well as category-selective information²¹. Other category-orthogonal object
66 properties, including size, position and pose, show higher population decoding performance in
67 monkey IT (analogous to human VTC) compared to early visual areas, contrary to what was

CATEGORY VS SHAPE IN DEEP NETWORKS & VISION

68 previously believed¹⁰. Indeed, the majority of visual object representations in IT may be
69 accounted for by object shape, or other low-level visual properties, rather than category²⁸.
70 Nevertheless, studies that explicitly de-confound category from more low-level properties
71 suggest that the category selectivity cannot be fully explained by these other properties^{10, 16, 21},
72 and point towards a so-called feature-dependent categorical code¹⁵.

73 In this paper, we explicitly dissociate shape from category in two stimulus sets to
74 determine: (i) how CNNs represent object shape and category when they are independent from
75 one another; and (ii) how these artificial representations correspond with shape and category
76 representations in human visual cortex. Using two carefully designed stimulus sets, which
77 orthogonalise shape and category, we assess four top-performing CNNs in their ability to
78 represent category independently from shape layer by layer. Taking the same two stimulus sets,
79 we measure human fMRI responses when viewing these images and assess the interaction
80 between shape and category along the visual ventral stream. Finally, we compare artificial
81 representations with human fMRI responses for the same two stimulus sets, to evaluate how
82 closely CNNs reflect biological representations.

83

84

CATEGORY VS SHAPE IN DEEP NETWORKS & VISION

85 Methods

86 We aimed to determine the relationship between models of shape and category, CNNs,
87 and neural responses in the human visual ventral pathway. We tested object shape and category
88 representation in four top-performing CNNs and compared this with behavioural ratings of shape
89 and category as well as human fMRI response patterns from experiments in two previous
90 studies^{16, 29}. Below we describe participants, stimulus sets, CNN architectures, the neuroimaging
91 experiments, and data analysis.

92

93 **Participants**

94 All participants gave written informed consent. All experiments were approved by the
95 Ethics Committee at KU Leuven and the University Hospitals Leuven. All methods were
96 performed in accordance with the relevant guidelines and regulations. For the behavioural ratings,
97 each stimulus set was rated by an independent group of participants (N= 4 for set A; N = 16 for
98 set B). For the neuroimaging experiments, there were 15 participants (8 females, mean age of 30
99 years) scanned in fMRI experiment A, none whom were excluded. There were also 15
100 participants (8 females, mean age of 24 years) scanned for fMRI experiment B, with one person
101 who was excluded due to excessive head motion. All subjects had normal or corrected vision.

102

103 **Stimulus sets**

104 The stimuli in both experiments were designed to dissociate shape from category
105 information. Both stimulus sets are grayscale images of objects on a white or grey background,
106 centred at the origin and presented at a normal viewing angle (see Figure 1). Set A contains 32
107 unique images, divided into 2 equally sized categories (animal vs non-animal) and 2 equally

CATEGORY VS SHAPE IN DEEP NETWORKS & VISION

108 sized groups of shapes (low and high aspect ratio). Set B contains 54 images divided into 6
109 object categories (minerals, animals, fruit/veg, music, sport and tools) and 9 shape types. The
110 model design for each stimulus set, which orthogonalises shape from category, is illustrated in
111 Figure 1. For additional information about the stimulus sets, refer to Ritchie and Op de Beeck²⁹
112 and Bracci and Op de Beeck¹⁶, for Set A and B respectively.

113 To confirm that shape was not predictive of category information for each of the stimulus
114 sets, we analysed the images using low-level GIST descriptors³⁰ and tested how well these visual
115 features predicted shape or category using Linear Discriminant Analysis (LDA). GIST provides
116 a low dimensional representation of an image based on spectral and coarsely localised
117 information. We defined the GIST descriptors to include 8 orientations over 8 scales and
118 combine this with LDA. For Set A, we ran a two-way classification using a leave-one-level out
119 procedure, for example, training on bar stimuli and generalising to blob stimuli to test for
120 animacy classification. For Set B, we followed a six-way classification using a leave-one-level
121 out test procedure, permuting across all possible groups of train and test combinations and
122 averaging across results. For example, we selected six shape clusters of the total nine, trained an
123 LDA on GIST descriptors from five clusters ($5 \times 6 = 30$ images) and tested whether the algorithm
124 could predict the 6 different categories from the held out images. All six-way shape and category
125 combinations were tested and averaged.

126

127 **Behavioural data**

128 Each stimulus set was rated on object category and shape properties by means of the
129 multiple object arrangement method³¹. Participants rated similarity in two task contexts: for
130 *object category*, “arrange the images based on the semantic similarity among objects”; for *object*

CATEGORY VS SHAPE IN DEEP NETWORKS & VISION

131 *shape*, “arrange the images based on perceived object shape similarity”. These models, based on
132 behavioural data, are meant to better represent the stimulus psychological space relative to the
133 stricter *design-based* models (2 categories x 2 shape types in set A; 6 categories x 9 shape types
134 in set B). For example, in Set B, the *design-based* shape model represents the 9 different shape
135 types as equidistant from one other, whereas the *behaviour-based* shape model is sensitive to
136 further variation between the 9 shape types in terms of between-type similarity. The behaviour-
137 based model for Set B illustrates that elongated objects (the final 3 shape types), regardless of
138 their orientation, are perceived as being more similar to each other relative to round objects (the
139 first 3 shape types), which is not visible in the design-based model. Figure 1A and 1B depicts
140 both *design-based* and *behaviour-based* models.

141

142 **fMRI Experiments**

143 Here we provide a summary of the fMRI procedures and analyses, the full details are
144 provided in Ritchie and Op de Beeck²⁹ for experiments using Set A and Bracci and Op de
145 Beeck¹⁶ for Set B.

146 **Preprocessing and Analysis**

147 All imaging data was pre-processed and analysed using SPM and MATLAB. For each
148 participant, fMRI data was slice-time corrected, motion corrected (using spatial realignment to
149 the first image), coregistered to each individual’s anatomical scan, segmented and spatially
150 normalised to the standard MNI template. Functional images were resampled to 3 x 3 x 3 mm
151 voxel size and spatially smoothed by convolving with a Gaussian kernel of 6mm FWHM for Set
152 A and 4mm FWHM for Set B³². After pre-processing, a GLM was used to model the BOLD
153 signal for each participant, for each stimulus, at each voxel. Regressors for the GLM included

CATEGORY VS SHAPE IN DEEP NETWORKS & VISION

154 each stimulus condition of interest (32 for A, 54 for B) and 6 motion correction parameters (x, y
155 and z coordinates for translation and rotation). Each predictor had its time course modelled as a
156 boxcar function convolved with the canonical haemodynamic response function, producing a
157 single estimate for each voxel per predictor for every run. The beta weights fitted to each GLM
158 were used to create Representational Dissimilarity Matrices (RDMs) for each participant
159 (defined below).

160

161 **Regions of Interest (ROIs)**

162 Neural representational content was investigated in three main ROIs in visual cortex: primary
163 visual cortex (V1), and ventral temporal cortex (VTC), which was split into posterior (VTC post)
164 and anterior (VTC ant) halves. These ROIs were chosen for their relevance in both object shape
165 and category information processing²³. VTC is bounded laterally by the occipitotemporal sulcus
166 (OTS), posteriorly by the posterior transverse collateral sulcus (ptCoS) and anteriorly by the
167 anterior tip of the mid-fusiform sulcus (MFS)²³. ROIs were defined at the group level by
168 combining the anatomical criteria above (using the Neuromorphometrics atlas in SPM) with
169 functional criteria (all active voxels for the contrast of all conditions versus baseline that
170 responded to visual information exceeding the statistically uncorrected threshold of $p < 0.001$ in
171 a second-level analysis). For further details on ROI definition, please refer to Bracci, Kalfas &
172 Op de Beeck³³ where the exact same ROI criteria were applied. We used a two-factor repeated-
173 measures Analysis of Variance Model (ANOVA) to assess the interaction between two within-
174 participant factors: conditions (shape, category) and area (V1, VTC post and VTC ant).

175

CATEGORY VS SHAPE IN DEEP NETWORKS & VISION

176 **Deep Neural Network Architectures**

177 Each architecture consists of multiple convolutional layers followed by pooling
178 operations and fully-connected layers. For each CNN, which was pre-trained on the ImageNet
179 dataset³⁴, we ran a forward pass of each image in the stimulus set through the network. We
180 output the activation of weights in each layer, resulting in a matrix with size of the *nodes per*
181 *layer* times *the stimulus set* (32 for A, 54 for B). We calculated *I - correlation* for each activation
182 pattern of one stimulus with another to obtain an RDM with size $N \times N$, where N = the number
183 of stimulus conditions (32 x 32 for A, 54 x 54 for B). We did not include final softmax
184 classification layers in our analysis, since we were interested in the structure of layer
185 representations and not classification performance per se.

186 **CaffeNet**

187 CaffeNet is an implementation of AlexNet¹ in the Caffe deep learning framework³⁵.
188 CaffeNet is an 8-layer convolutional neural network (CNNs) with five convolutional layers and
189 three fully connected layers.

190 **VGG-19**

191 VGG-19³ was the top ranking CNN for single object localisation in ILSVRC 2014, and
192 second-running in image classification³⁴. VGG-19 consists of 19 weighted layers with an
193 additional softmax read-out layer for classification. The architecture contains 16 convolutional
194 layers separated by five max pooling layers, with the final 3 layers being fully-connected.

195 **GoogLeNet**

196 GoogLeNet², also known as InceptionNet, was the top-performing architecture for image
197 classification in ILSVRC 2014³⁴. GoogLeNet is a 22-layer deep network, when counting only
198 parameterised layers, or 27 layers deep if including pooling operations. All convolution,

CATEGORY VS SHAPE IN DEEP NETWORKS & VISION

199 reduction and projection layers use rectified linear activation. The bottom layers of the network
200 follow conventional convolutional neural network architecture, consisting of chained
201 convolutional operations followed by max pooling. The top layers of the network replace
202 multiple fully-connected layers with an average pooling layer, a single fully connected layer and
203 a classification layer. The middle layers of the network differ substantially from traditional
204 convolutional neural network structure, consisting of stacked “inception” modules, which are
205 miniature networks containing one max pooling and 3 multi-sized convolution operations (1 x 1,
206 3 x 3 and 5 x 5 convolutions) in parallel configuration. Convolution operations inside inception
207 modules are optimised with dimensionality reduction, by preceding expensive 3 x 3 and 5 x 5
208 convolution operations with 1 x 1 convolutions. Inception modules allow for increased width of
209 the network, as well as depth, while maintaining a constant computational budget.

210 **ResNet50**

211 ResNets are a family of extremely deep architectures that won the ILSVRC classification
212 task in 2015³⁶. ResNet50 contains 50 stacked “residual units”, which use a split-transform-merge
213 strategy to perform identity mappings in parallel to 3x3 convolutions with rectification. ResNets,
214 like GoogLeNet², are multi-branch architectures, containing only 2 branches (performing identity
215 projection and 3x3 convolutions) instead of GoogLeNet’s maximum 4 branch inception modules
216 (performing multi-size convolutions). Identity mappings perform a key role in the architecture’s
217 success, forcing the network to preserve features, rather than learn entirely new representations
218 at every layer, as is the case with conventional CNNs³⁷. The final 3 layers of ResNet50 are
219 identical in design to GoogleNet, performing average pooling, transformation to 1000
220 dimensions using full connections and softmax classification (not included in our analysis).

221

CATEGORY VS SHAPE IN DEEP NETWORKS & VISION

222 **Representational Similarity Analysis**

223 We used Representational Similarity Analysis (RSA) to quantitatively compare CNN
224 representations per layer with design models, behavioural ratings, and with fMRI neuroimaging
225 data. RSA compares RDMs, which characterise the representational information in a brain or
226 model³⁸. Given a set of activity patterns (biological, behavioural or artificial) for a set of
227 experimental conditions, the dissimilarity between patterns is computed as 1 minus the
228 correlation across the units that compose the patterns. RDMs are symmetrical about a zero
229 diagonal, where 0 denotes perfect correlation. RSA assesses second-order isomorphism, which is
230 the shared similarity in structure between dissimilarity matrices³⁹. Spearman rank order
231 correlation was used to compare dissimilarity matrices, since the relationship between RDMs
232 cannot be assumed to be linear³⁸. In cases where there was any dependency relationship between
233 shape and category RDMs (visible in the Set A behavioural data), we used partial correlation.
234 We determined the significance of every correlation by comparing it with a null distribution
235 obtained by randomly permuting the RDM labels and then calculating dissimilarity relationships
236 1000 times.

237

238 **Results**

239 **Behavioural Data**

240 For each stimulus set, participants provided similarity judgments for the shape and
241 category dimension (see Figure 1, right column). For Set A, we found a significant correlation
242 between the behavioural models for shape and category (Spearman's $\rho = 0.4753$, $p < 0.001$
243 permutation test with 10000 randomisations of stimulus labels) and so partial correlations when
244 carrying out RSA with Set A behavioural models. For Set A, as expected, behavioural and

CATEGORY VS SHAPE IN DEEP NETWORKS & VISION

245 design category models strongly correlate with one another ($\rho = 0.8555$, $p < 0.001$) and design
246 shape strongly correlates with behavioural shape ($\rho = 0.7849$, $p < 0.001$). For Set B, we found no
247 significant correlation between behavioural models for shape and category ($\rho = 0.006$, $p =$
248 0.8209). Again, as expected, shape behavioural and design models were significantly correlated
249 ($\rho = 0.4145$, $p < 0.001$) and category behavioural and design models were also significantly
250 correlated ($\rho = 0.6195$, $p < 0.001$).

251

252 **Low-level Shape Analysis of Stimuli**

253 Using GIST³⁰ descriptors of each image and combining this with LDA, we confirmed
254 that category could not be predicted based upon these low-level descriptors whereas shape could,
255 demonstrating that our stimulus sets were properly orthogonalised. LDA with GIST predicted
256 shape above chance level, at 87.5% for Set A and 69% for Set B. Category was predicted below
257 chance level, at 37.5% for Set A and 10% for Set B.

258

259 **Shape and category RSA on all CNN layers for Stimulus Sets A and B**

260 Figure 2 illustrates layer-by-layer RSA between the CNN representations and the shape
261 and category models and behavioural data in the two stimulus sets. Note that all RSA using Set
262 A behavioural models involved partial correlations (explained above in Behavioural data).
263 Looking across all networks, in the first layer of all CNNs, shape is already represented above
264 the significance threshold in most cases, whereas category is not. Shape correlations at the first
265 layer of CNNs are lower and closer to the significance threshold for Set A (design $0.12 < \rho <$
266 0.22 , behavioural $0.12 < \rho < 0.24$) than Set B (design $0.26 < \rho < 0.44$, behavioural $0.24 < \rho <$
267 0.36). CaffeNet shows the highest correlation in shape information at the first layer with both

CATEGORY VS SHAPE IN DEEP NETWORKS & VISION

268 behavioural and design models for both stimulus sets. In CaffeNet, there is a single rise and fall
269 in shape information, except in the Set A behavioural model. In all other networks, shape
270 correlations fluctuate along the layers, with peaks at different layers before decreasing at the
271 final layer in all cases except for Set A GoogLeNet and ResNet50. For Set A, shape correlations
272 remain relatively high at the final layer (design $0.34 < \rho < 0.51$, behavioural $0.29 < \rho < 0.59$). In
273 contrast, for Set B, shape correlation levels increase in the networks before falling in the final
274 layers of all networks, to below their first layer levels for the design model correlations ($0.11 < \rho$
275 < 0.14), or to roughly their initial values for the behavioural model correlations ($0.32 < \rho < 0.36$).
276 For all networks, category information remains low across the majority of layers, hovering at or
277 below the significance level until the final few layers, where it increases above the significance
278 threshold to peak at the final layer. At the final layer, for Set A, category correlations reach
279 between $0.31 < \rho < 0.42$ for design models and between $0.34 < \rho < 0.42$ for behavioural. For Set
280 B, category correlations reach between $0.11 < \rho < 0.21$ for design models and between $0.24 < \rho$
281 < 0.37 for design models at the final layer.

282 To investigate the interaction between shape and category and CNN layers, we tested
283 correlation values in a 2 X 2 ANOVA with Layer (modelled linearly with intercept and slope)
284 and Condition (Shape or Category). *Table 1* summarises the statistical results of the main effects
285 (layer, condition) and their interaction in CNNs and models. For Set A, for both types of models
286 across all networks, layer has a highly significant main effect and condition is also significant
287 (*Table 1*) which suggests that correlation values can be predicted given the CNN layer and the
288 condition of interest (shape or category information). Their interaction is significant in
289 GoogLeNet and VGG-19, but not in CaffeNet and ResNet50, suggesting that as category
290 increases, shape decreases significantly in two out of the four networks tested. For Set B, across

CATEGORY VS SHAPE IN DEEP NETWORKS & VISION

291 all networks, condition is highly significant, and layer has a significant main effect in
292 behavioural model correlations, however regarding design model correlations, layer is only
293 significant in one CNN (ResNet50). This suggests that it is possible to make significant
294 predictions of behavioural shape and category judgements given CNN layer information,
295 however this prediction does not extend to design models of shape and category. Condition is
296 highly significant across all networks, and the interaction between layer and condition is
297 significant for both models and CaffeNet, and the design model and GoogleNet.

298 In summary, across both Sets A and B, we can see that shape information gradually
299 increases and/or wavers as the network is traversed, before falling in the final layers. The peak
300 value in shape information remains roughly the same regardless of network depth. Peak category
301 correlations also remain roughly the same regardless of network depth. Across both Sets A and B,
302 category information is at or below the significance threshold in the initial layer before reaching
303 the maximum value at the final layer, showing the opposite trend with shape correlations.
304 Interestingly, the maximum levels of shape and category correlations do not depend on network
305 depth, nor on architectural design differences, such as the use of inception modules. Figure 3
306 contains multidimensional scaling plots of peak design shape and category information for Sets
307 A and B.

308

309 **Shape Versus Category information in Visual Ventral Stream Regions**

310 Figure 4 summarises the representational similarity in three regions of interest (ROIs)
311 along the visual ventral pathway, from low-level area V1 through to posterior and anterior VTC,
312 compared with design and behavioural models of shape and category. For Set A, shape
313 information reduces slightly along the ventral stream, from 22% to 19% in design models, and

CATEGORY VS SHAPE IN DEEP NETWORKS & VISION

314 18% to 10% in behavioural models. Category information increases along the ventral pathway,
315 from -3% to 41% in design models, and -6% to 40% in behavioural models. We tested RSA
316 results using a two-factor ANOVA, with ROI (V1, VTC ant, VTC post) and Condition (category,
317 shape) as within-subject factors. For Set A, results reveal a significant main effect for ROI ($F_{2, 15}$
318 = 26.34, $p < 0.001$ for the design model; $F_{2, 15} = 35.81$, $p < 0.001$ for behavioural), whereas the
319 main effect of Condition (shape vs category) is not significant ($F_{1, 15} = 0.56$, for design; $F_{1, 15} =$
320 1.02, for behavioural). There is a significant interaction between ROI and Condition ($F_{2, 15} =$
321 68.14, $p < 0.001$ for design, $F_{2, 15} = 73.34$, $p < 0.001$ for behavioural), indicating that as category
322 information increases from V1 to VTC ant, shape information decreases. Post hoc pairwise t-
323 tests further confirmed the dissociation between shape and category along the visual ventral
324 stream: category divisions were able to significantly better explain the neural pattern in later
325 ventral areas (VTC ant) relative to shape ($t_{(15)} = 8.57$, $p < 0.0001$ for design models, $t_{(15)} = 5.67$, p
326 < 0.0001 for behavioural models); whereas the opposite was true in early visual area V1, where
327 shape was significantly more related to the neural data compared to category divisions ($t_{(15)} =$
328 6.34, $p < 0.0001$ for design models, $t_{(15)} = 8.16$, $p < 0.0001$ for behavioural models).

329 For Set B, we see a qualitatively similar trend of decreasing shape information from V1
330 to VTC anterior (from 10% to 0% in the design models, and from 18% to 4% in the behavioural
331 models) and increasing category information (from -1% to 6% in the design models, and from
332 1% to 6% in the behavioural models). The two-factor ANOVA, with ROI (V1, VTC ant, VTC
333 post) and Condition (category, shape), revealed that when correlating ROI representations with
334 the design models for Set B, ROI has no significant effect ($F_{2, 14} = 0.57$, ns), the effect of
335 Condition is significant ($F_{1, 14} = 11.39$, $p < 0.01$) and there is a highly significant interaction
336 effect between area and condition ($F_{2, 14} = 36.71$, $p < 0.001$). Analysing correlations with the

CATEGORY VS SHAPE IN DEEP NETWORKS & VISION

337 behavioural models for Set B, the effect of area is significant ($F_{2, 14} = 3.79, p = 0.027$), as is
338 condition ($F_{1, 14} = 33.84, p < 0.001$) and there is a highly significant interaction effect between
339 area and condition ($F_{2, 14} = 13.33, p < 0.001$). Again, pairwise t-tests further confirmed the
340 dissociation between shape and category in visual ventral brain regions, with shape being
341 significantly more related to neural data in early visual area V1 than category ($t_{(14)} = 7.56, p <$
342 0.0001 for design models, $t_{(14)} = 5.28, p = 0.0001$ for behavioural models); and category able to
343 explain neural patterns more in VTC ant than shape (significantly for design models $t_{(14)} = 3.89,$
344 $p = 0.0007$, but not significantly for behavioural models: $t_{(14)} = 1.20, p = 0.24$). Thus, there is a
345 two-way interaction between shape and category across the visual ventral stream that is
346 significant for both stimulus sets and both model types, illustrating a decrease in shape combined
347 with an increase in category going from V1 to VTC anterior.

348

349 **RSA for fMRI Brain Data and all CNN layers**

350 Neural fMRI responses for each participant, and ROI, for Set A and Set B were
351 correlated with the RDMs of every layer for each CNN. Results are shown in Figure 5. For each
352 stimulus set and network, correlation values were tested in a 2 X 3 ANOVA with Layer
353 (modelled linearly with intercept and slope) and ROI as within subject factors. In CaffeNet, V1
354 and VTC posterior correlations peaked at the third convolutional layer, and VTC anterior peaks
355 at the final layer for both stimulus sets. For both stimulus sets, the 2 X 3 ANOVA results reveal a
356 significant main effect of ROI (Set A: $F_{2, 15} = 88.73, p < 0.001$; Set B: $F_{2, 14} = 57.00, p < 0.001$)
357 and Layer (Set A: $F_{1, 15} = 41.06, p < 0.001$; $F_{1, 14} = 48.38, p < 0.001$) and their interaction (Set A:
358 $F_{2, 15} = 133.72, p < 0.001$; Set B: $F_{2, 14} = 44.88, p < 0.001$). In VGG-19, both stimulus sets show
359 similar peaks in correlations, with V1 reaching a maximum at layer 13, VTC posterior at layer 15,

CATEGORY VS SHAPE IN DEEP NETWORKS & VISION

360 and VTC anterior at the final 19th layer. For both sets, there is a significant main effect of ROI
361 (Set A: $F_{2,15} = 59.12, p < 0.001$; Set B: $F_{2,14} = 26.98, p < 0.001$) and Layer (Set A: $F_{1,15} = 294.14,$
362 $p < 0.001$; $F_{1,14} = 40.30, p < 0.001$). The ROI x Layer interaction is significant in Set A ($F_{2,15} =$
363 $55.49, p < 0.001$), but does not reach significance in Set B ($F_{2,14} = 2.76, p = 0.06$). GoogLeNet
364 has multiple peaks for correlations with V1 and VTC posterior, and there is a clear peak in VTC
365 anterior in the final layer for both stimulus sets. For both Sets, ROI (Set A: $F_{2,15} = 73.76, p <$
366 0.001 ; Set B: $F_{2,14} = 37.07, p < 0.001$), Layer (Set A: $F_{1,15} = 152.19, p < 0.001$; Set B: $F_{1,14} =$
367 $18.08, p < 0.001$) and their interaction (Set A: $F_{2,15} = 130.85, p < 0.001$; Set B: $F_{2,14} = 12.46, p <$
368 0.001) are all highly significant. Finally, in ResNet50, V1 peaks at layers 44 to 47, VTC
369 posterior peaks at layers 47 to 49, and VTC anterior peaks at the final layer. For both Sets, ROI
370 (Set A: $F_{2,15} = 31.20, p < 0.001$; Set B: $F_{2,14} = 20.26, p < 0.001$) and Layer (Set A: $F_{1,15} =$
371 $1431.40, p < 0.001$; Set B: $F_{1,14} = 895.32, p < 0.001$) are highly significant, and their interaction
372 is significant (Set A: $F_{2,15} = 5.97, p = 0.003$; Set B: $F_{2,14} = 52.54, p < 0.001$). Together these
373 results show that across all deep neural networks, there is a cascade in correlation peaks from V1
374 to VTC posterior to VTC anterior along the layers of each network, matching with the flow of
375 activation along the human visual ventral pathway. For all networks, and both stimulus sets, the
376 highest correlation of VTC anterior occurs at the final layer.

CATEGORY VS SHAPE IN DEEP NETWORKS & VISION

377 Discussion

378 In this study, we investigated orthogonal shape and category representations in biological
379 and artificial networks by making comparisons between: (i) CNNs and models of shape and
380 category; (ii) models and the brain; and (iii) CNNs and the brain. First, comparing artificial
381 networks and models, we found that CNNs represent category information as well as shape, and
382 that category information peaks at the final layer for all tested CNNs, regardless of network
383 depth. Peak correlation levels for shape and category do not increase with network depth, and
384 remain roughly at the same level regardless of architectural design differences, including the use
385 of inception modules or residual networks. Second, comparing models and the brain, there is a
386 two-way interaction between shape and category in the human visual ventral pathway, where
387 shape is best represented earlier in V1, and category emerges later in anterior VTC. This
388 interaction between shape and category is significant across both stimulus sets and for both
389 design and behavioural models. Third, comparing artificial networks and the brain, V1 correlates
390 highest with early to mid-level layers of deep networks, and anterior VTC correlates best with
391 the final layer of CNNs. Across both stimulus sets and for all networks, peak correlations with
392 V1 always occur in earlier network layers than peak correlations with anterior VTC,
393 demonstrating that CNNs reflect a similar order of computational stages as the human ventral
394 pathway when processing these object images.

395 Our results allow for a greater understanding of how shape and category are represented
396 in deep networks and in the visual ventral pathway, in particular: (i) how differing shape and
397 category definitions between the two stimulus sets reveal differences between low-level and
398 high-level shape representations in CNNs and the brain; (ii) how shape and category processing

CATEGORY VS SHAPE IN DEEP NETWORKS & VISION

399 along deep network layers maps onto brain regions; and iii) how careful stimulus design allows
400 us to make better inferences about category semantics in the brain and in CNNs.

401 One major advantage of this study is that we consider two stimulus sets that carefully
402 control shape and category to draw conclusions about their interaction and interplay, rather than
403 broadly extrapolating results based on a single set of images. These two well-controlled stimulus
404 sets are similar in design but differ slightly in how shape and category are defined, allowing us to
405 extract a finer interpretation of results. Looking at the differences in shape definitions between
406 these stimulus sets, in Set A, shape is defined with a low to high aspect ratio (described as “bar-
407 like” or “blob-like”), while it is characterized retinotopically in Set B. Comparing CNNs and
408 models, both low-level (Set B) and high-level (Set A) shape information is preserved until the
409 very last layer of all networks, however there is a visible reduction in low-level compared to
410 high-level shape information in the final layers. Comparing models and the brain, we see that the
411 high-level (Set A) shape information remains quite high in VTC ant, compared to low-level (Set
412 B) shape information, which reduces to correlation levels that are at or near zero. The plausible
413 explanation for why shape information drops off in Set B but not in A, is that higher level
414 regions represent a more abstract form of shape, which is factored into the design of Set A, but
415 not B. Indeed, previous studies showed that perceived shape similarity strongly overlaps with
416 higher-level brain representations in humans⁴⁰, and in monkeys^{12, 41}. Kalfas *et al.*¹² found that the
417 deepest layers of networks, rather than IT responses, correlated best with human shape similarity
418 judgements. We also found that CNNs correlated much higher with behavioural shape
419 judgements than fMRI. This finding suggests that there is at least some correspondence between
420 how humans and models use shape, even though there are very likely also differences (see e.g.
421 Baker *et al.*¹⁹).

CATEGORY VS SHAPE IN DEEP NETWORKS & VISION

422 Considering the differences in category definitions between the stimulus sets, Set A has
423 only two category clusters defined by the animate-inanimate division, whereas Set B has six
424 object clusters. The number of groups clearly affects the size difference in correlation levels
425 between category models and CNNs as well as the brain, where fewer groupings boost the signal.
426 In the final layer of all CNNs, we see that category, as defined by animacy in Set A, reaches
427 correlation levels up to three times the magnitude of Set B. Considering brain data, category as
428 defined by animacy in Set A reaches six times the magnitude in VTC ant compared to Set B.
429 This is consistent with existing studies that show a strong animacy division in higher-level
430 regions of visual cortex²⁴. We find that in all four networks, human similarity judgements of
431 category are best explained by the final layer of CNNs, more so than fMRI representations in late
432 ventral areas.

433 Our use of multiple CNNs allows us to observe the influence of network depth on peak
434 correlations with brain regions. Hong *et al.*⁹ compared their brain data to a CNN consisting of 6
435 parallelised convolutional layers, finding that the model's top hidden layer was most predictive
436 of IT response patterns and that lower layers had higher resemblance to V1-like Gabor patterns.
437 Consistent with their findings, we also found that the final layer of CNNs had maximum
438 correspondence with later ventral stream areas, and that earlier layers showed higher correlation
439 with V1. Cichy *et al.*¹⁴ found peak V1 correlations in the second layer of an 8-layer CNN trained
440 for object recognition. Similarly in our experiments, we found that peak V1 correlations occurred
441 at layer 3 in an 8-layer network (CaffeNet) for both stimulus sets. As network depth increases,
442 peak correlations with V1 shift from earlier tiers in the network to later layers. Interestingly,
443 some of the highest V1 correlations occur immediately prior to fully connected layers, as is the
444 case in ResNet50 and VGG-19. Figure 5 illustrates peak V1 correlations occurring as late as the

CATEGORY VS SHAPE IN DEEP NETWORKS & VISION

445 45th layer in ResNet50, bringing into question the explanatory value of additional processing
446 stages in deeper networks, especially when an 8-layer network achieves similar magnitudes of
447 correlation with V1 by the third layer. Nevertheless, while the maximum correlation values of
448 brain regions shift to later layers in larger networks, the rank-order of correlation peaks with
449 brain regions still matches the order of information processing along the ventral pathway. That is,
450 correlations with V1 always peak before VTC ant, regardless of network depth. We extend upon
451 the findings of Cichy *et al.*¹⁴ on the order of visual information processing from a single 8 layer
452 network to multiple networks, including a 50 layer network.

453 Recently, there has been some effort directed towards investigating the role of semantic
454 representations in deep visual networks, and where category semantics may be represented in the
455 ventral pathways¹³. Deriving high-level semantic meaning from low-level feature descriptions is
456 commonly referred to as the “semantic gap” in computer vision literature⁴². In order to fully
457 establish the level at which CNNs are able bridge the semantic gap, and extract meaningful
458 information from images, it is necessary to remove all possible reliance on low-level features,
459 which could be exploited to improve performance, and test network performance on carefully
460 designed images that minimise potential dependencies between category and influencing features.
461 Devereux *et al.*¹³ do not properly control for the influence of shape, as we have, and include
462 many low-level visual features labelled misleadingly as “semantic” descriptors, such as “is
463 circular/round” or “is “green”, which we would argue do not allow for a dissociation between
464 vision and semantics¹⁵. Our study explicitly defines category semantics as falling within the
465 animacy division in Set A, or in multiple object categories (animals, minerals, fruit/vegetables,
466 music, sports equipment and tools) in Set B. Our stimulus sets do not confound category
467 semantics with shape information, allowing us to draw firmer conclusions.

CATEGORY VS SHAPE IN DEEP NETWORKS & VISION

468 In conclusion, despite shape and category often being confounded in natural images, and
469 the possibility for artificial neural networks to exploit this correlation when performing
470 classification tasks, we find that deep convolutional neural networks are able to represent
471 category information independently from low-level shape in a manner similar to higher level
472 visual cortex in humans.
473

CATEGORY VS SHAPE IN DEEP NETWORKS & VISION

474

References

- 475 1. Krizhevsky, A., Sutskever, I., & Hinton, G. E. ImageNet Classification with Deep
476 Convolutional Neural Networks. *Advances in Neural Information Processing Systems 25*
477 (*NIPS 2012*), pp. 1097-1105. Lake Tahoe: Curran Associates, Inc. (2012).
- 478 2. Szegedy, C., *et al.* Going deeper with convolutions. *2015 IEEE Conference on Computer*
479 *Vision and Pattern Recognition (CVPR)*, pp. 1-9. Boston, MA (2015).
- 480 3. Simonyan, K., & Zisserman, A. Very Deep Convolutional Networks for Large-Scale Image
481 Recognition. *ICLR*, Preprint at: <https://arxiv.org/abs/1409.1556> (2015).
- 482 4. He, K., Zhang, X., Ren, S., & Sun, J. Delving Deep into Rectifiers: Surpassing Human-Level
483 Performance on ImageNet Classification. *2015 IEEE International Conference on Computer*
484 *Vision (ICCV)*, pp. 1026-1034. Santiago (2015).
- 485 5. Kriegeskorte, N. Deep Neural Networks: A New Framework for Modeling Biological Vision
486 and Brain Information Processing. *Annual Review of Vision Science*, **1**, 417-446.
487 doi:10.1146/annurev-vision-082114-035447 (2015).
- 488 6. Kietzmann, T. C., McClure, P., & Kriegeskorte, N. Deep Neural Networks In Computational
489 Neuroscience. *bioRxiv*. Preprint at: <https://doi.org/10.1101/133504> (2017).
- 490 7. Khaligh-Razavi, S.-M., & Kriegeskorte, N. Deep Supervised, but Not Unsupervised, Models
491 May Explain IT Cortical Representation. *PLoS Computational Biology*, **10**(11), e1003915.
492 doi:10.1371/journal.pcbi.1003915 (2014).
- 493 8. Cadieu, C. F., *et al.* Deep Neural Networks Rival the Representation of Primate IT Cortex for
494 Core Visual Object Recognition. *PLoS Computational Biology*, **10**(12), e1003963.
495 doi:10.1371/journal.pcbi.1003963 (2014).

CATEGORY VS SHAPE IN DEEP NETWORKS & VISION

- 496 9. Yamins, D. L., *et al.* Performance-optimized hierarchical models predict neural responses in
497 higher visual cortex. (T. J. Sejnowski, Ed.) *PNAS*, **111**(23), 8619-8624.
498 doi:10.1073/pnas.1403112111 (2014).
- 499 10. Hong, H., Yamins, D. L., Majaj, N. J., & DiCarlo, J. J. Explicit information for category-
500 orthogonal object properties increases along the ventral stream. *Nature Neuroscience*, **19**(4),
501 613–622. doi:10.1038/nn.4247 (2016).
- 502 11. Güçlü, U., & van Gerven, M. A. Deep Neural Networks Reveal a Gradient in the Complexity
503 of Neural Representations across the Ventral Stream. *The Journal of Neuroscience*, **35**(27),
504 10005-10014 (2015).
- 505 12. Kalfas, I., Vinken, K., & Vogels, R. Representations of regular and irregular shapes by deep
506 Convolutional Neural Networks, monkey inferotemporal neurons and human judgments.
507 *PLoS Computational Biology*, **14**(10), e1006557. doi:10.1371/journal.pcbi.1006557 (2018).
- 508 13. Devereaux, B. J., Clarke, A., & Tyler, L. K. Integrated deep visual and semantic attractor
509 neural networks predict fMRI pattern-information along the ventral object processing
510 pathway. *Scientific Reports*, **8**,10636. doi:10.1038/s41598-018-28865-1 (2018).
- 511 14. Cichy, R. M., Khosla, A., Pantazis, D., Torralba, A., & Oliva, A. Comparison of deep neural
512 networks to spatio-temporal cortical dynamics of human visual object recognition reveals
513 hierarchical correspondence. *Scientific Reports*, **6**, 27755. doi:10.1038/srep27755 (2016).
- 514 15. Bracci, S., Ritchie, J. B., & Op de Beeck, H. On the partnership between neural
515 representations of object categories and visual features in the ventral visual pathway.
516 *Neuropsychologia*, **105**, 153-164 (2017).
- 517 16. Bracci, S., & Op de Beeck, H. Dissociations and Associations between Shape and Category.
518 *The Journal of Neuroscience*, **36**(2), 432-444 (2016).

CATEGORY VS SHAPE IN DEEP NETWORKS & VISION

- 519 17. Belongie, S., Malik, J., & Puzicha, J. Shape Matching and Object Recognition Using Shape
520 Contexts. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, **24**(4), 509–522
521 (2002).
- 522 18. Kubilius, J., Bracci, S., & Op de Beeck, H. P. Deep Neural Networks as a Computational
523 Model for Human Shape Sensitivity. *PLoS Computational Biology*, **12**(4), e1004896.
524 doi:10.1371/journal.pcbi.1004896 (2016).
- 525 19. Baker, N., Lu, H., Erlikhman, G., & Kellman, P. Deep convolutional networks do not
526 classify based on global object shape. *PLoS Computational Biology*, **14**(12), e1006613.
527 doi:10.1371/journal.pcbi.1006613 (2018).
- 528 20. Ritter, S., Barrett, D. G., Santoro, A., & Botvinick, M. M. Cognitive Psychology for Deep
529 Neural Networks: A Shape Bias Case Study. *Proceedings of the 34th International
530 Conference on Machine Learning, PMLR 70*, pp. 2940-2949. Sydney, Australia (2017).
- 531 21. Kaiser, D., Azzalini, D. C., & Peelen, M. V. Shape-independent object category responses
532 revealed by MEG and fMRI decoding. *Journal of Neurophysiology*, **115**, 2246-2250.
533 doi:10.1152/jn.01074.2015 (2016).
- 534 22. Proklova, D., Kaiser, D., & Peelen, M. V. Disentangling Representations of Object Shape
535 and Object Category in Human Visual Cortex: The Animate-Inanimate Distinction. *Journal
536 of Cognitive Neuroscience*, **28**(5), 680-692. (2016).
- 537 23. Grill-Spector, K., & Weiner, K. S. The functional architecture of the ventral temporal cortex
538 and its role in categorization. *Nature Reviews Neuroscience*, **15**(8), 536-548.
539 doi:10.1038/nrn3747 (2014).

CATEGORY VS SHAPE IN DEEP NETWORKS & VISION

- 540 24. Kriegeskorte, N., *et al.* Matching Categorical Object Representations in Inferior Temporal
541 Cortex of Man and Monkey. *Neuron*, **60**(6), 1126-41. doi:10.1016/j.neuron.2008.10.043
542 (2008).
- 543 25. Kiani, R., Esteky, H., Mirpour, K., & Tanaka, K. Object category structure in response
544 patterns of neuronal population in monkey inferior temporal cortex. *Journal of*
545 *Neurophysiology*, **97**, 4296-4309. (2007).
- 546 26. Rice, G. E., Watson, D. M., Hartley, T., & Andrews, T. J. Low-Level Image Properties of
547 Visual Objects Predict Patterns of Neural Response across Category-Selective Regions of the
548 Ventral Visual Pathway. *Journal of Neuroscience*, **34**(26), 8837-8844.
549 doi:10.1523/JNEUROSCI.5265-13.2014 (2014).
- 550 27. Andrews, T. J., Watson, D. M., Rice, G. E., & Hartley, T. Low-level properties of natural
551 images predict topographic patterns of neural response in the ventral visual pathway. *Journal*
552 *of Vision*, **15**(7), 1-12. doi:10.1167/15.7.3 (2015).
- 553 28. Baldassi, C., Alemi-Neissi, A., Pagan, M., DiCarlo, J., Zecchina, R., & Zoccolan, D. Shape
554 Similarity, Better than Semantic Membership, Accounts for the Structure of Visual Object
555 Representations in a Population of Monkey Inferotemporal Neurons. *PLoS Computational*
556 *Biology*, **9**(8), e1003167. doi:10.1371/journal.pcbi.1003167 (2013).
- 557 29. Ritchie, J. B., & Op de Beeck, H. Using neural distance to predict reaction time for
558 categorizing the animacy, shape, and abstract properties of objects. *bioRxiv*. Preprint at:
559 <https://doi.org/10.1101/496539> (2018).
- 560 30. Oliva, A., & Torralba, A. Modeling the shape of the scene: a holistic representation of the
561 spatial envelope. *International Journal of Computer Vision*, **42**(3), 145-175. (2001).

CATEGORY VS SHAPE IN DEEP NETWORKS & VISION

- 562 31. Kriegeskorte, N., & Mur, M. Inverse MDS: inferring dissimilarity structure from multiple
563 item arrangements. *Frontiers in Psychology*, **3**:245. doi:10.3389/fpsyg.2012.00245 (2012).
- 564 32. Op de Beeck, H. P. Against hyperacuity in brain reading: spatial smoothing does not hurt
565 multivariate fMRI analyses? *Neuroimage*, **49**, 1943–1948. (2010).
- 566 33. Bracci, S., Kalfas, I., & Op de Beeck, H. The ventral visual pathway represents animal
567 appearance over animacy, unlike human behavior and deep neural networks. *bioRxiv*.
568 Preprint at: <http://dx.doi.org/10.1101/228932> (2017)
- 569 34. Russakovsky, O., *et al.* ImageNet Large Scale Visual Recognition Challenge. *International*
570 *Journal of Computer Vision (IJCV)*, **115**(3), 211–252. Preprint at:
571 <https://arxiv.org/abs/1409.0575> doi:10.1007/s11263-015-0816-y (2015).
- 572 35. Jia, Y., *et al.* Caffe: Convolutional Architecture for Fast Feature Embedding. Preprint at
573 <https://arxiv.org/abs/1408.5093> (2014).
- 574 36. He, K., Zhang, X., Ren, S., & Sun, J. Deep Residual Learning for Image Recognition.
575 *ArXiv:1512.03385 [Cs]*. Preprint at:<http://arxiv.org/abs/1512.03385> (2015).
- 576 37. Greff, K., Srivastava, R. K., & Schmidhuber, J. Highway and Residual Networks learn
577 Unrolled Iterative Estimation. *International Conference on Learning Representations (ICLR)*.
578 Preprint at: <https://arxiv.org/abs/1612.07771> (2017).
- 579 38. Kriegeskorte, N., Mur, M., & Bandettini, P. Representational similarity analysis – connecting
580 the branches of systems neuroscience. *Frontiers in Systems Neuroscience*, **2**(4).
581 doi:10.3389/neuro.06.004.2008 (2008).
- 582 39. Shepard, R. N., & Chipman, S. Second-order isomorphism of internal representations:
583 Shapes of states. *Cognitive Psychology*, **1**(1), 1-17. (1970).

CATEGORY VS SHAPE IN DEEP NETWORKS & VISION

- 584 40. Op de Beeck, H. P., Torfs, K., & Wagemans, J. Perceived shape similarity among unfamiliar
585 objects and the organization of the human object vision pathway. *J. Neurosci.*, **28**(40)
586 10111–10123. doi:10.1523/JNEUROSCI.2511-08.2008 (2008).
- 587 41. Op de Beeck, H., Wagemans, J., & Vogels, R. Inferotemporal neurons represent low-
588 dimensional configurations of parameterized shapes. *Nature Neuroscience*, **4**(12), 1244-1252.
589 (2001).
- 590 42. Markowska-Kaczmar U., Kwaśnicka H. (2018) Deep Learning—A New Era in Bridging the
591 Semantic Gap. In: *Bridging the Semantic Gap in Image and Video Analysis*. (eds. Kwaśnicka
592 H., Jain L.) Intelligent Systems Reference Library, vol 145. doi:[https://doi.org/10.1007/978-](https://doi.org/10.1007/978-3-319-73891-8_7)
593 [3-319-73891-8_7](https://doi.org/10.1007/978-3-319-73891-8_7) (Springer, 2018)
- 594

CATEGORY VS SHAPE IN DEEP NETWORKS & VISION

595 Acknowledgements

596 A.A.Z. and H.O.d.B were funded by grant C14/16/031 of the KULeuven Research
597 Council. J.B.R. received funding from the FWO and European Union's Horizon 2020 research
598 and innovation programme under the Marie Skłodowska-Curie grant agreement No 665501, via
599 a FWO [PEGASUS]² Marie Skłodowska-Curie fellowship (12T9217N). S.B. is funded by FWO
600 (Fonds Wetenschappelijk Onderzoek) postdoctoral fellowship 516 (12S1317N). Neuroimaging
601 was funded by the Flemish Government Hercules Grant ZW11_10. We would like to thank Tim
602 Leers for help with data analysis.

603 Data Availability

604 The datasets generated during and/or analysed during the current study are available from
605 the corresponding author on reasonable request.

606

607 Author Information

608 **Affiliations**

609 Laboratory of Biological Psychology – KU Leuven, Leuven, Belgium

610 Astrid A. Zeman*, J. Brendan Ritchie, Stefania Bracci, Hans Op de Beeck

611

612 **Contributions**

613 All authors contributed to the study design. SB, JBR and HOdB provided pre-processed
614 neuroimaging data and collected behavioural data. AAZ ran network simulations, analysed the
615 data and wrote the manuscript with input from all authors. All authors interpreted the data, edited
616 the manuscript and approved the final version.

617

CATEGORY VS SHAPE IN DEEP NETWORKS & VISION

618 **Competing Interests**

619 The authors declare no competing interests.

620

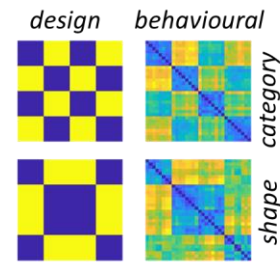
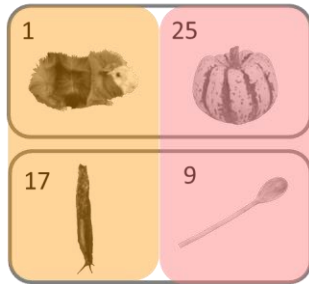
621 **Corresponding Author**

622 Correspondence to Astrid Zeman at astrid.zeman@kuleuven.be.

623

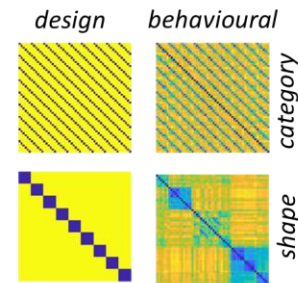
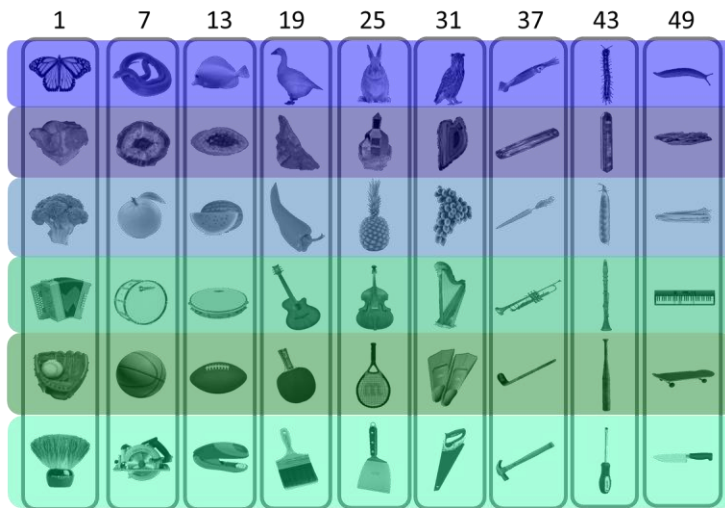
CATEGORY VS SHAPE IN DEEP NETWORKS & VISION

A



624

B



625

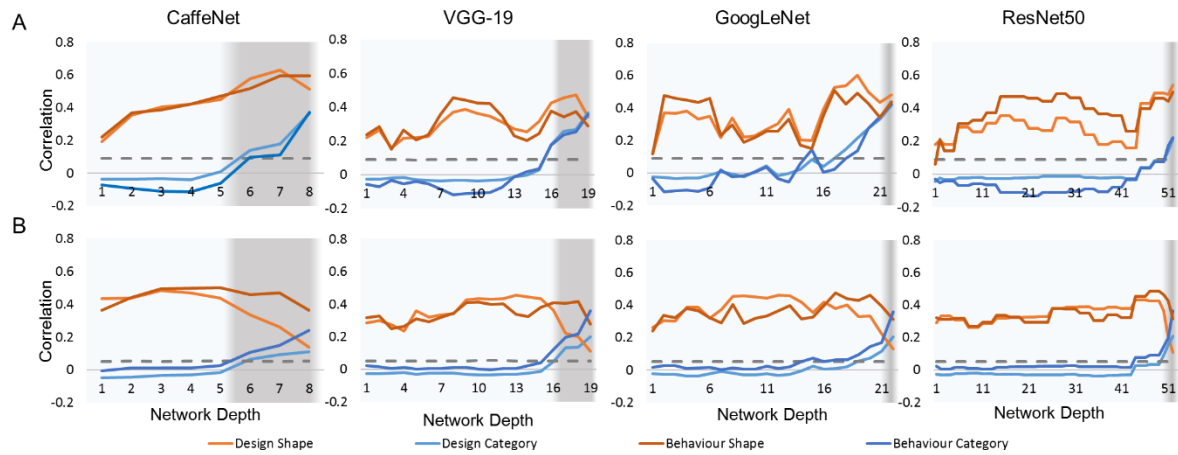
626

627 *Figure 1 (A) 32 stimuli in 2 categories (animal and non-animal), (B) 54 stimuli in 6 categories*
628 *(animals, minerals, fruit/vegetables, music, sports equipment, tools). Left: Each category*
629 *division is highlighted by a distinct colour. Common shape information is circled in grey.*
630 *Numbers indicate indexing for RDMs. Due to copyright restrictions, not all images are shown in*
631 *Set A and the ones displayed are representative. Set A images are published in compliance with*
632 *a CC BY-SA license (<https://creativecommons.org/licenses/by-sa/3.0/>) and their sources are:*
633 *guinea pig (<https://commons.wikimedia.org/wiki/File:AniarasKelpoKalle.jpg> by Tavu); squash*
634 *(<https://commons.wikimedia.org/wiki/File:Festival-Squash.jpg> by Evan-Amos); slug (Black Slug*
635 *at Aggregate Ponds, <https://www.flickr.com/photos/brewbooks/2606728819> by brewbooks); and*
636 *wooden spoon (https://upload.wikimedia.org/wikipedia/commons/7/7b/Wooden_Spoon.jpg by*
637 *Donovan Govan). Images have been changed to greyscale and have the background removed.*
638 *The final two images have also been rotated. Set B images are published in compliance with a*
639 *CC-BY license (<https://creativecommons.org/licenses/by/4.0/>) and are re-used from Figure 5a in*
640 *Kubilius, Bracci and Op de Beeck¹⁸. Right: Shape and category RDMs. The design models are*
641 *based on the experimental design. The behavioural models are obtained via multiple object*
642 *arrangement³¹; see methods.*

CATEGORY VS SHAPE IN DEEP NETWORKS & VISION

643

644



645

646

647 *Figure 2: Correlation between layers in CNNs and shape (orange/red) versus category (blue) in*

648 *Set A (top row) and B (bottom row). The horizontal axis indicates network depth and the vertical*

649 *axis indicates correlation (Spearman's ρ). For GoogLeNet and ResNet architectures, the*

650 *correlations shown are for 3x3 convolutional operations, while other parallel operations*

651 *(projections and convolutions of different sizes) are omitted. Dashed line indicates significance*

652 *threshold of $p < 0.05$. Grey shading indicates fully-connected layers.*

653

654

CATEGORY VS SHAPE IN DEEP NETWORKS & VISION

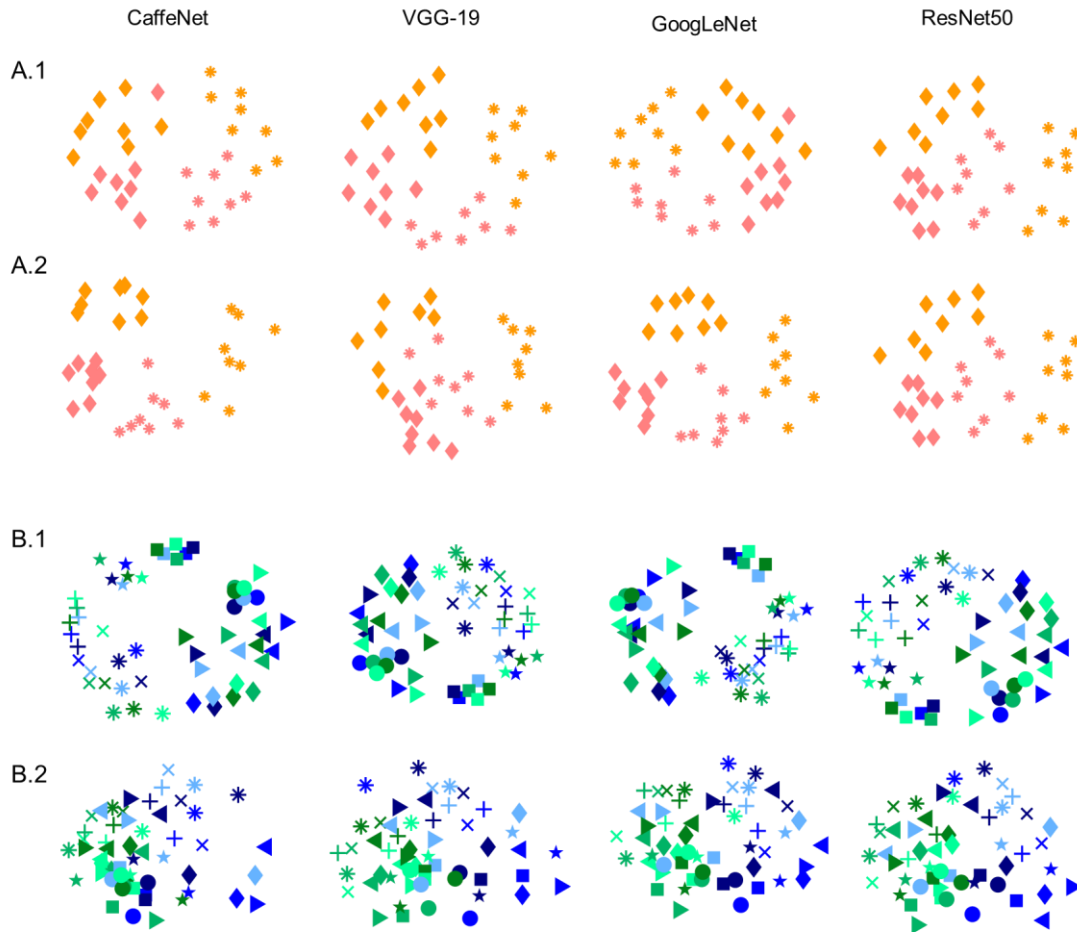
Stimulus Set	Network	Number of Layers	Model	Layer $F_{1,1}$	Layer p	Condition $F_{1,1}$	Condition p	Interaction $F_{1,1}$	Interaction p
A	CaffeNet	8	D	92.338	<0.001	14.825	0.002	4.729	0.050
			B	41.233	<0.001	126.651	<0.001	0.202	0.661
	VGG-19	19	D	53.304	<0.001	11.477	0.002	6.496	0.016
			B	17.370	<0.001	99.161	<0.001	6.252	0.017
	GoogLeNet	22	D	46.438	<0.001	8.390	0.006	4.464	0.041
			B	18.59	<0.001	87.68	<0.001	10.21	0.003
	ResNet50	52	D	28.788	<0.001	41.075	<0.001	1.662	0.200
			B	25.010	<0.001	750.551	<0.001	0.323	0.571
B	CaffeNet	8	D	1.766	0.208	173.677	<0.001	29.577	<0.001
			B	8.306	0.014	212.106	<0.001	7.774	0.016
	VGG-19	19	D	2.567	0.118	160.955	<0.001	4.023	0.053
			B	22.075	<0.001	207.91	<0.001	3.536	0.069
	GoogLeNet	22	D	2.026	0.162	312.186	<0.001	8.551	0.006
			B	27.727	<0.001	329.938	<0.001	1.833	0.183
	ResNet50	52	D	26.517	<0.001	1377.504	<0.001	0.012	0.913
			B	61.007	<0.001	1108.272	<0.001	0.311	0.578

655
656

657 *Table 1: 2 X 2 ANOVA results of Layer (modelled linearly with slope and intercept) and*
 658 *Condition (shape or category) and their interaction in CNNs and models (D = design, B =*
 659 *behavioural).*

660

CATEGORY VS SHAPE IN DEEP NETWORKS & VISION



661

662 *Figure 3: Multidimensional scaling plots of 1) Peak design shape correlations with common*
663 *shape represented by common symbols, and 2) peak category correlations, with common*
664 *category represented by shared colour, for each network and Set A (top 2 rows) and B (bottom 2*
665 *rows). Colour coding corresponds to Figure 1.*

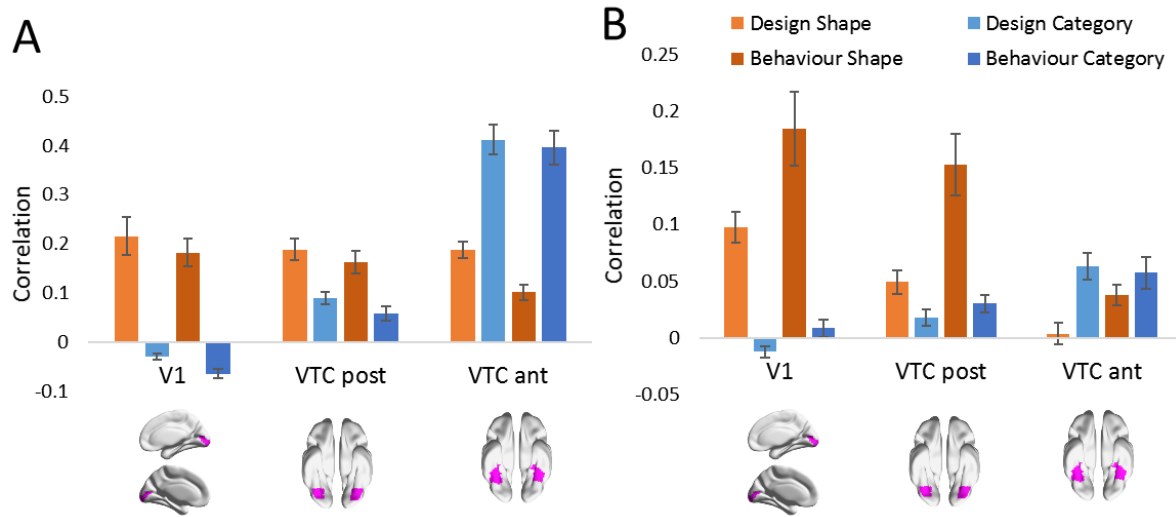
666

667

668

669

CATEGORY VS SHAPE IN DEEP NETWORKS & VISION



670

671

672 *Figure 4 RSA results for shape and category models for Set A (left) and B (right) in ROIs. Three*

673 *regions along the ventral visual pathway are analysed: V1, VTC post and VTC ant. Error bars*

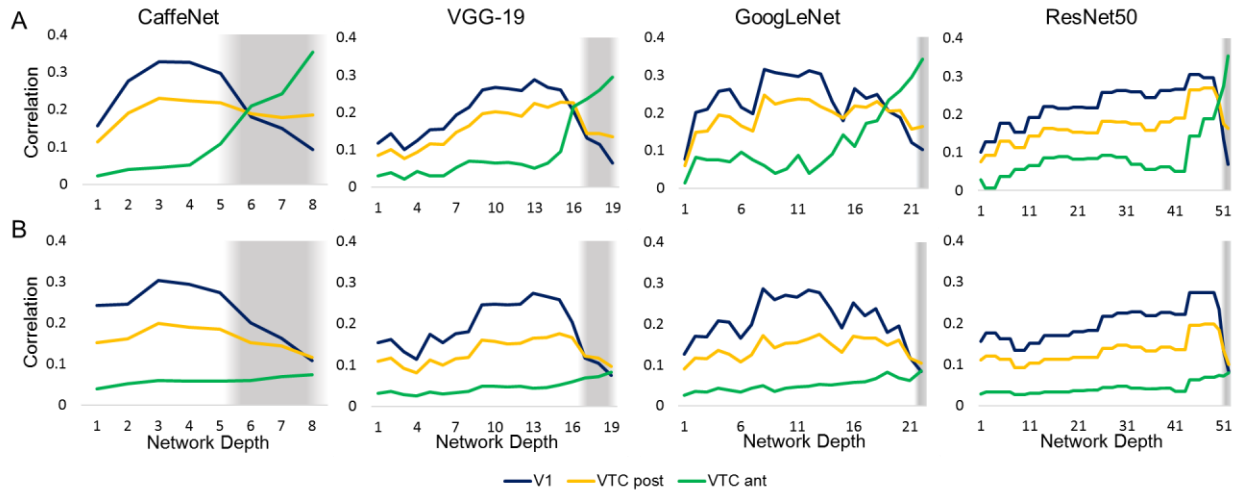
674 *represent standard error. ROI visualisations are re-used from Fig 4A in (Bracci, Kalfas, & Op*

675 *de Beeck³³, p. 8). Note the difference in scale between A and B.*

676

CATEGORY VS SHAPE IN DEEP NETWORKS & VISION

677



678

679 *Figure 5: RSA comparing models (CaffeNet, VGG-19, GoogLeNet and ResNet50) and fMRI*

680 *activation in V1 (navy), VTC post (yellow) and VTC ant (green) ROIs for Sets A (top row) and B*

681 *(bottom row). Grey shading indicates fully-connected layers.*

682