

# The face module emerges from domain-general visual experience: a deprivation study on deep convolutional neural network

1 Shan Xu<sup>#\*</sup>, Yiyuan Zhang<sup>#</sup>, Zonglei Zhen, & Jia Liu<sup>\*</sup>

2 **#: equal contribution**

3 **\*: correspondence author**

4 Beijing Key Laboratory of Applied Experimental Psychology, Faculty of Psychology, Beijing  
5 Normal University, Beijing 100875, China.

6

7 **\* Correspondence:**

8 Jia Liu, Ph.D. or Shan Xu, Ph.D. or; Room 1415, New Main Building, 19 Xijiekouwai St, Haidian  
9 District, Beijing 100875, China. Tel.: +86-10-58806154; fax: +86-10-58806154. E-mail:

10 liujia@bnu.edu.cn or shan.xu@bnu.edu.cn.

11

12 **Abstract**

13 Can faces be accurately recognized with zero experience on faces? The answer to this question is  
14 critical because it examines the role of experiences in the formation of domain-specific modules in  
15 the brain. However, thorough investigation with human and non-human animals on this issue cannot  
16 easily dissociate the effect of the visual experience from that of genetic inheritance, i.e., the  
17 hardwired domain-specificity. The present study addressed this problem by building a model of  
18 selective deprivation of the experience on faces with a representative deep convolutional neural  
19 network (DCNN), AlexNet. We trained a new AlexNet with the same image dataset, except that all  
20 images containing faces of human and nonhuman primates were removed. We found that the  
21 experience-deprived AlexNet (d-AlexNet) did not show significant deficits in face categorization and  
22 discrimination, and face-selective modules also automatically emerged. However, the deprivation  
23 made the d-AlexNet to process faces in a more parts-based fashion, similar to the way of processing  
24 objects. In addition, the face representation of the face-selective module in the d-AlexNet was more  
25 distributed and the empirical receptive field was larger, resulting in less degree of selectivity of the  
26 module. In sum, our study provides undisputable evidence on the role of nature versus nurture in  
27 developing the domain-specific modules that domain-specificity may evolve from non-specific  
28 stimuli and processes without genetic predisposition, which is further fine-tuned by domain-specific  
29 experience.

30 **Keywords:** face perception, face domain, deep convolutional neural network, visual deprivation,  
31 experience

## 32 1 Introduction

33 A fundamental question in cognitive neuroscience is how nature and nurture form our cognitive  
34 modules. In the center of the debate is the origin of face recognition ability. Numerous studies have  
35 revealed both behavioral and neural signatures of face-specific processing, indicating a face module  
36 in the brain (for reviews, see Freiwald, Duchaine, & Yovel, 2016; Kanwisher & Yovel, 2006).  
37 Further studies from behavioral genetics revealed the contribution of genetics on the development of  
38 the face-specific recognition ability in humans (Wilmer et al., 2010; Zhu et al., 2010). Collectively,  
39 these studies suggest an innate domain-specific module for face cognition. However, it is unclear  
40 whether the visual experience is also necessary for the development of the face module.

41 A direct approach to address this question is visual deprivation. Two studies on monkeys  
42 selectively deprived the visual experience of faces since birth, while leaving the rest of experiences  
43 untouched (Arcaro, Schade, Vincent, Ponce, & Livingstone, 2017; Sugita, 2008). They report that  
44 face-deprived monkeys are still capable of categorizing and discriminating faces (Sugita, 2008),  
45 though less prominent in selective looking preference to faces over non-face objects (Arcaro et al.,  
46 2017). Further examination of the brain of the experience-deprived monkeys fails to localize typical  
47 face-selective cortical regions with the standard criterion; however, in the inferior temporal cortex  
48 where face-selective regions are normally localized, face-selective activation (i.e., neural responses to  
49 faces larger than nonface objects) is observed (Arcaro et al., 2017). Taken together, without visual  
50 experiences of faces, rudimental functions to process faces may still evolve to some extent.

51 Two related but independent hypotheses may explain the emergence of the face module without  
52 face experiences. An intuitive answer is that the rudimental functions are hardwired in the brain by  
53 genetic predisposition (McKone, Crookes, Jeffery, & Dilks, 2012; Wilmer et al., 2010).  
54 Alternatively, we argue that the face module may emerge from experiences on nonface objects and  
55 related general-purpose processes, because representations for faces may be constructed by abundant  
56 features derived from nonface objects. Unfortunately, studies on humans and monkeys are unable to  
57 thoroughly decouple the effect of nature and nurture to test these two hypotheses.

58 Recent advances in deep convolutional neural network (DCNN) provide an ideal test platform to  
59 examine the role of visual experiences alone on face modules without genetic predisposition.  
60 Previous studies have shown that DCNNs are similar to human visual cortex both structurally and  
61 functionally (Kriegeskorte, 2015), but free of any predisposition on functional modules. Therefore,

62 with DCNNs we can manipulate experiences without considering interactions from genetic  
63 predisposition. In this study, we asked whether DCNNs can achieve face-specific recognition ability  
64 when visual experiences on faces were selectively deprived.

65 To do this, we trained a representative DCNN, AlexNet (Krizhevsky, 2014; Krizhevsky,  
66 Sutskever, & Hinton, 2012), to categorize nonface objects with face images carefully removed from  
67 the training dataset. Once this face-deprived DCNN (d-AlexNet) was trained, we compared its  
68 behavioral performance to that of a normal AlexNet of the same architecture but with faces present  
69 during training in both face categorization (i.e., differentiating faces from nonface objects) and  
70 discrimination (i.e., discriminating faces among different individuals) tasks. We predicted that the d-  
71 AlexNet, though without predisposition and experiences of faces, may still develop face selectivity  
72 through its visual experiences of nonface objects.

73

## 74 **2 Materials and methods**

### 75 **2.1 Stimuli**

76 **Deprivation dataset** The deprivation dataset was constructed to train the d-AlexNet. It was based on  
77 the ImageNet Large-Scale Visual Recognition Challenge (ILSVRC) 2012 dataset (Deng et al., 2009),  
78 which contains 1,281,167 images for training and 50,000 images for validation, in 1000 categories.  
79 These images were first subjected to automated screening with an in-house face-detection toolbox  
80 based on VGG-Face (Parkhi, Vedaldi, & Zisserman, 2015), and then further screened by two human  
81 raters, who separately judged whether a given image contains faces of humans or non-human  
82 primates regardless of the orientation and intactness of the face, or anthropopathic artwork, cartoons,  
83 and artifacts. We removed images judged by either rater as containing any above-mentioned contents.  
84 Finally, we removed categories whose remaining images were less than 640 images (approximately  
85 half of the original number of images in a category). The resultant dataset consists of 736 categories,  
86 with 662,619 images for training and 33,897 for testing the performance.

87 **Classification dataset** To train a classifier that can classify faces, we constructed a classification  
88 dataset consisting of 204 categories of non-face objects and one face category, each of 80 exemplars.  
89 For the non-face categories, we manually screened Caltech-256 (Griffin, Holub, & Perona, 2007) to  
90 remove images containing human, primate, or cartoon faces, and then removed categories whose

91 remaining images were less than 80. In each of the 204 remaining non-face categories, we randomly  
92 chose 70 images for training and another 10 for calculating classification accuracy. The face category  
93 was constructed by randomly selecting 1000 faces images from Faces in the Wild (FITW) dataset  
94 (Berg, Berg, Edwards, & Forsyth, 2005). Among them, 70 were used as training data and another 10  
95 for classification accuracy. In addition, to characterize DCNN's ability in differentiating faces from  
96 object categories, we compiled a second dataset consisting of all images in the face category except  
97 those used in training.

98 **Discrimination dataset** To train a classifier that can discriminate faces at individual level, we  
99 constructed a discrimination dataset consisting of face images of 133 individuals, 300 images each,  
100 selected from the Casia-WebFace database (Yi, Lei, Liao, & Li, 2014). For each individual in the  
101 dataset, 250 were randomly chosen for training and another 50 for calculating discrimination  
102 accuracy.

103 **Representation dataset** To examine representational similarity of faces and non-face images  
104 between the d-AlexNet and the normal one, we constructed a representation dataset with two  
105 categories, faces and bowling pins as an 'unseen' non-face object category that was not presented to  
106 the DCNNs during training. Each category consisted of 80 images. The face images were a random  
107 subset of FITW, and images of bowling pins were randomly chosen from the corresponding category  
108 in Caltech-256.

109 **Movies clips for DCNN-brain correspondence analysis** We examined the correspondence between  
110 the face-selective response of the DCNNs and brain activity using a set of 18 clips of 8-min natural  
111 color videos from the Internet that are diverse yet representative of real-life visual experiences (Wen  
112 et al., 2017).

## 113 **2.2 The deep convolutional neural network**

114 Our model of selective deprivation, the d-AlexNet, was built with the architecture of the well-known  
115 DCNN 'AlexNet' (Krizhevsky et al., 2012, see Figure 1a for illustration). AlexNet is a feed-forward  
116 hierarchical convolutional neural network consisting of five convolutional layers (denoted as Conv1  
117 – Conv5, respectively) and three fully connected layers denoted as FC1 – FC3. Each convolutional  
118 layer consists of a convolutional sublayer, followed by a ReLU sublayer, and Conv1, 2, and 5 are  
119 further followed by a pooling sublayer. Each convolutional sublayer consists of a set of distinct  
120 channels. Each channel convolves the input with a distinct linear filter (kernel) which extracts filtered

121 outputs from all locations within the input with a particular stride size. FC1 to FC3 are fully  
122 connected layers. FC3 is followed by a sublayer using a softmax function to output a vector that  
123 represents the probability of the visual input containing the corresponding object category  
124 (Krizhevsky et al., 2012).

125 The d-AlexNet used the architecture of AlexNet but changed the number of units in FC3 to 736,  
126 so was the following softmax function, to match the number of categories in the deprivation dataset.  
127 Same to the pre-training AlexNet in pytorch 1.2.0 (<https://pytorch.org/>, Paszke et al., 2017), the d-  
128 AlexNet was initialized with values drawn from a uniform distribution, and was then trained on the  
129 deprivation dataset following the approach specified in Krizhevsky et al., (2014). We used the pre-  
130 trained AlexNet from pytorch 1.2.0 as the normal DCNN, referred to as the AlexNet in this paper for  
131 brevity.

132 The present study referred to channels in the convolutional sublayers by the layer they belong to  
133 and a channel index, following the convention of pytorch 1.2.0. For instance, Layer 5-Ch256 refers to  
134 the 256<sup>th</sup> convolutional channel of Layer 5.

### 135 **2.3 Transfer learning for classification and discrimination**

136 To examine to what extent our manipulation of the visual experience affected the categorical  
137 processing of faces, we replaced the fully-connected layers of each DCNN with a two-layer face-  
138 classification classifier. The first layer was a fully connected layer with 43,264 units as inputs and  
139 4,096 units as outputs with sigmoid activation function, and the second was a fully connected layer  
140 with 4,096 units as inputs and 205 units as outputs, each of which corresponded to one category of  
141 the classification dataset. This classifier, therefore, classified each image into one category of the  
142 classification dataset. The face-classification classifier was trained for each DCNN with the training  
143 images in the classification dataset for 90 epochs.

144 To examine to what extent our manipulation of the visual experience affected face  
145 discrimination, we similarly replaced the fully connected layers of each DCNN with a discrimination  
146 classifier. The discrimination classifier differed from the classification classifier only in its second  
147 layer, which had 133 units instead as outputs, each corresponding to one individual in the  
148 discrimination dataset. The face-discrimination classifier was trained for each DCNN with the  
149 training images in the discrimination dataset for 90 epochs.

## 150 **2.4 The face selective channels in DCNNs**

151 To identify the channels selectively responsive to faces, we submitted images in the classification  
152 dataset to each DCNN, recorded the average activation in each channel of Conv5 after ReLU in  
153 response to each image, and then averaged the channel-wise activation within each category. We  
154 selected channels where the face category evoked the highest activation, and used the Mann-Whitney  
155 U test to examine the activation difference between faces and objects that had the second-highest  
156 activation in these channels ( $p < .05$ , Bonferroni corrected). The selectivity of each face channel thus  
157 identified was indexed by the selective ratio. The selective ratio was calculated by dividing the face  
158 activation by the second-highest activation. In addition, we measured the lifetime sparseness of each  
159 face-selective channel as an index for selectivity of faces among all non-face objects. We first  
160 normalized the mean activations of a face channel in Layer5 to all the categories to the range of 0-1,  
161 and then calculated lifetime sparseness with the formula:

$$162 \quad S = (\sum_{i=1,n} r_i / n)^2 / \sum_{i=1,n} (r_i^2 / n)$$

163 where  $r_i$  is the normalized activations to the  $i$ th object category. The smaller this value is, the higher  
164 the selectivity is.

## 165 **2.5 DCNN-Brain Correspondence**

166 We submitted the movie clips to the DCNNs. Following Wen (2017)'s approach, we extracted and  
167 log-transformed the channel-wise output (the average activation after ReLU) of each face-selective  
168 channel using DNNBrain, an in-house toolbox (Chen et al., 2020), and then convolved it with a  
169 canonical hemodynamic response function (HRF) with a positive peak at 4s. The HRF convolved  
170 channel-wise activity was then down-sampled to match the sampling rate of functional magnetic  
171 resonance imaging (fMRI) and the resultant timeseries was standardized before further analysis.

172 Neural activation in the brain was derived from the preprocessed data in Wen (2017). The  
173 fMRI data were recorded while human participants viewed each movie clips twice. We averaged the  
174 standardized time series across repetition and across subjects for each clip. Then, for each DCNN, we  
175 conducted multiple regression for each clip, with the activation time series of each brain vertex as the  
176 dependent variable and that of face-selective channels in this network as independent variables. For  
177 the d-AlexNet, all face-selective channels were included. For the AlexNet, we included the same  
178 number of face-selective channels with the highest face selectivity to match the complexity of the

179 regression model. We used the  $R^2$  of each vertex as the index of the overall Goodness of fit of the  
180 regression in that vertex. The  $R^2$  values were then averaged across clips. The larger the  $R^2$  value, the  
181 higher correspondence between the DCNN and the brain in response to movie clips.

182 To determine whether cortical regions with large  $R^2$  values were traditional face-selective  
183 regions, we delineated the bilateral fusiform face areas (FFA) and the occipital face area (OFA) with  
184 the maximum-probability atlas of face-selective regions (Zhen et al., 2015). Two hundred of vertexes  
185 of the highest probability of the left FFA and 200 of the right FFA were included in the ROI of FFA,  
186 and the ROI of OFA was delineated in the same way. The correspondence with brain activation in  
187 each ROI and the impact of the visual experience was examined by submitting the vertex-wise  $R^2$   
188 into a two-way ANOVA with visual experience (d-AlexNet vs. AlexNet) as within-subject factor and  
189 regional correspondence (OFA and FFA) as between-subject factor.

## 190 **2.6 Face inversion effect in DCNNs**

191 The average activation amplitude of the top 2 face-selective channels of each DCNN in response to  
192 upright and inverted version of 20 faces from the Reconstructing Faces dataset (VanRullen & Reddy,  
193 2019) was measured. The inverted faces were generated by vertically flipping the upright ones. The  
194 face inversion effect in the d-AlexNet was measured with paired sample t-tests (two-tailed) and the  
195 impact of the experience on the face inversion effect was examined by two-way ANOVAs with  
196 visual experience (d-AlexNet vs. AlexNet) and inversion (upright vs inverted) as within-subject  
197 factors.

## 198 **2.7 Representational similarity analysis**

199 To examine whether faces in the d-AlexNet were processed in an object-like fashion, we compared  
200 the within-category representational similarity of faces to that of bowling pins, an ‘unseen’ non-face  
201 object category never exposed to either DCNN. Specifically, for each image in the representation  
202 dataset, we arranged the average activations of each channel of Conv5 after ReLU into vectors, and  
203 then for each pair of images we calculated and then Fisher-z transformed the correlation between  
204 their vectors, which served as an index of pairwise representational similarity. Within-category  
205 similarity between pairs of face images and that between pairs of object images were calculated  
206 separately. A  $2 \times 2$  ANOVA was conducted with visual experience (d-AlexNet vs AlexNet) and  
207 category (face vs object) as independent factors. In addition, cross-category similarity between faces



208 and bowling pins was also calculated for each DCNN, and a paired sample t-test (two-tailed) on two  
209 DCNNs was conducted.

## 210 **2.8 Sparse coding and empirical receptive field**

211 To quantify the degree of sparseness of the face-selective channels in representing faces, we  
212 submitted the same set of 20 natural images containing faces from FITW to each DCNN, and  
213 measured the number of activated units (i.e., the units showing above-zero activation) in the face-  
214 selective channels. The more non-zero units of the face-selective channels, the less sparse of the  
215 representation for faces. The coding sparseness of the two DCNNs was compared with a paired-  
216 sample t-test.

217 We also calculated the size of the empirical receptive field of the face-selective channels.  
218 Specifically, we obtained activation maps of 1000 images randomly chosen from FITW. Using an in-  
219 house toolbox DNNBrain (Chen et al., 2020), we up-sampled each activation maps to the same size  
220 of the input. For each image, we averaged the up-sampled activation within the theoretical receptive  
221 field of each unit (the part of the image covered by the convolution of this unit and the preceding  
222 computation, decided by the network architecture), and selected the unit with the highest average  
223 activation. We then cropped the up-sampled activation map by the theoretical receptive field of this  
224 unit, to locate the image part that activated this channel most across all the units. Then, we averaged  
225 corresponding cropped activation maps across all the face images, and the resultant map denotes the  
226 empirical receptive field of this channel, delineating the part of the theoretical receptive field that  
227 causes this channel to respond strongly in viewing its preferred stimuli.

228

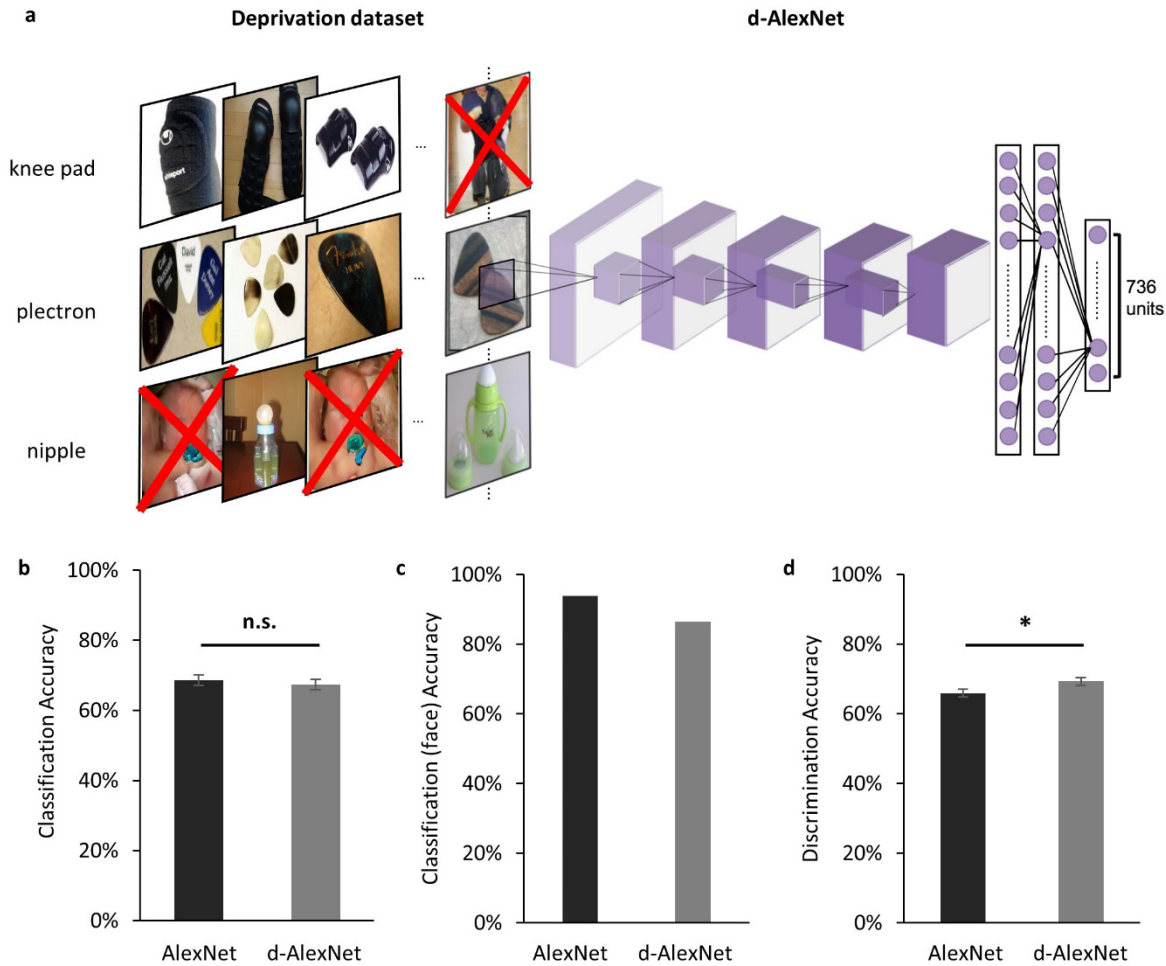
## 229 **3 Results**

230 The d-AlexNet was trained with a dataset of 662,619 non-face images consisting of 736 non-face  
231 categories, generated by removing images containing faces from the ILSVRC 2012 dataset (Figure  
232 1a). The d-AlexNet was initialized and trained in the same way as the AlexNet, and the resultant top-  
233 1 accuracy (57.29%) and the top-5 accuracy (80.11%) were comparable with the pre-trained  
234 AlexNet.

235 We first examined the performance of the d-AlexNet in two representative tasks of face  
236 processing, face categorization (i.e., differentiating faces from non-face objects) and face

237 discrimination (i.e., identifying different individuals). The output of Conv5 after ReLU of the d-  
238 AlexNet was used to classify objects in the classification dataset. The averaged categorization  
239 accuracy of the d-AlexNet (67.40%) was well above the chance level (0.49%), and comparable to  
240 that in the AlexNet (68.60%,  $t(204) = 1.26$ ,  $p = 0.209$ , Cohen's  $d = 0.007$ , Figure 1b). Critically, the  
241 d-AlexNet, although with no experience on faces, succeeded in the face categorization task, with an  
242 accuracy of 86.50% in categorizing faces from non-face objects. Note that the accuracy was  
243 numerically smaller than the AlexNet's accuracy in categorizing faces (93.90%) though (Figure 1c).

244 A similar pattern was observed in the face discrimination task. In this task, the output of Conv5  
245 after ReLU of each DCNN was used to identify 33,250 face images into 133 identities in the  
246 discrimination dataset. As expected, the AlexNet was capable of face discrimination (65.9%), well  
247 above the chance level (0.75%), consistent with previous studies (AbdAlmageed et al., 2016;  
248 Grundstrom, Chen, Ljungqvist, & Astrom, 2016). Critically, the d-AlexNet also showed the  
249 capability of discriminating faces, with an accuracy of 69.30% that was even significantly higher  
250 than that of the AlexNet,  $t(132) = 3.16$ ,  $p = .002$ , Cohen's  $d = 0.20$ , (Figure 1d). Taken together,  
251 visual experiences on faces seemed not necessary for developing basic functions of processing faces.



252

253 Figure 1. (a) An illustration of the screening to remove images containing faces for the d-AlexNet.

254 The 'faces' shown in the figure were AI-generated for illustration purpose only, and therefore have

255 no relation to real person. In the experiment, face images were from the ImageNet, with real persons'

256 faces. (b) The classification performance across categories of the two DCNNs was comparable. (c)

257 Both DCNNs achieved high accuracy in categorizing faces from other images. (d) Both DCNNs'

258 performance in discriminating faces was above the chance level, and the d-AlexNet's accuracy was

259 significantly higher than that of the AlexNet. The error bar denotes standard error. The asterisk

260 denotes statistical significance ( $\alpha = .05$ ). n.s. denotes no significance.

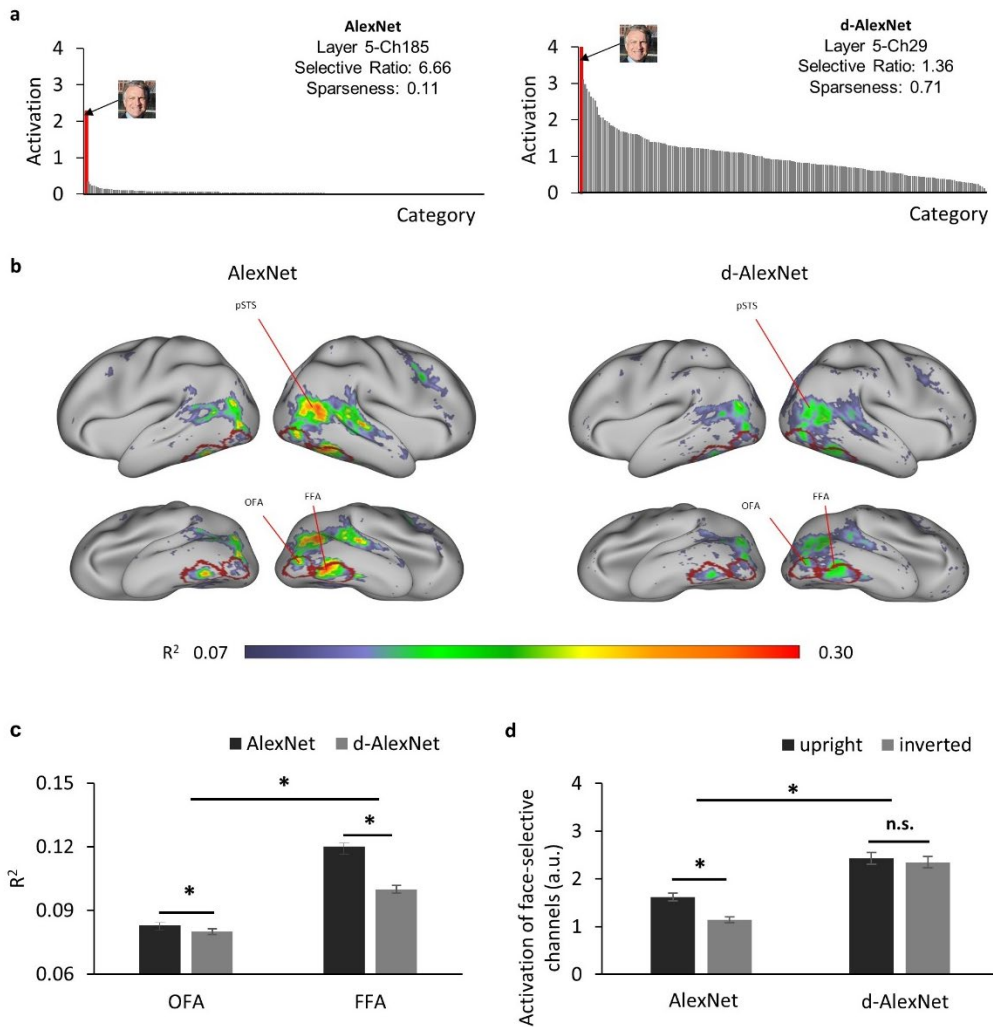
261

262 Was a face module formed in the d-AlexNet to support these functions? To answer this  
263 question, we searched all the channels in Conv5 of the d-AlexNet, where face-selective channels  
264 have been previously identified in the AlexNet (Baek, Song, Jang, Kim, & Paik, 2019). To do this,  
265 we calculated the activation of each channel in Conv5 after ReLU in response to each category of the  
266 classification dataset, and then identified channels that showed significantly higher response to faces  
267 than non-face images with Mann-Whitney U test ( $p < .05$ , Bonferroni corrected). Two face-selective  
268 channels (Ch29 and Ch50) met this criterion in the d-AlexNet (for an example channel, see Figure  
269 2a, right), whereas four face-selective channels (Ch195, Ch125, Ch60, and Ch187) were identified in  
270 the AlexNet (for an example channel, see Figure 2a, left). The face-selective channels in two DCNNs  
271 differed in selectivity. The averaged selective ratio, the ratio of the activation magnitude to faces by  
272 that to the most activated non-face object category, was 1.29 (range: 1.22 - 1.36) in the d-AlexNet,  
273 much lower than that in the AlexNet (average ratio: 3.63, range: 1.43 - 6.66). The lifetime sparseness,  
274 which measures the breadth of tuning of a channel in response to a set of categories, also showed a  
275 similar result. The average lifetime sparseness index of the face channels in the AlexNet (mean =  
276 0.25, range: 0.11 - 0.51) was smaller than that in the d-AlexNet (mean = 0.71, range: 0.70 - 0.71),  
277 indicating higher face selectivity in the AlexNet than that in the d-AlexNet. Taken together, this  
278 finding suggested that the face-selective channels already emerged in the d-AlexNet, though the face  
279 selectivity was weaker.

280 How did the face-selective channels correspond to face-selective cortical regions in humans,  
281 such as the FFA and OFA? To answer this question, we calculated the coefficient of determination  
282 ( $R^2$ ) of the multiple regression with the output of the face-selective channels as regressors and the  
283 fMRI signals from human visual cortex in response to movies on natural vision as the regressand. As  
284 shown in Figure 2b (right), the face-selective channels identified in the d-AlexNet corresponded to  
285 the bilateral FFA, OFA, and the posterior superior temporal sulcus face area (pSTS-FA). Similar  
286 correspondence was also found with the top two face-selective channels in the AlexNet (Figure 2b,  
287 left). Direct visual inspection revealed that the deprivation weakened the correspondence between the  
288 face-selective channels and face-selective regions in human brain. This observation was confirmed  
289 by the main effect of visual experiences ( $F(1,798) = 161.97, p < .001, \text{partial } \eta^2 = 0.17$ ) in a two-way  
290 ANOVA of visual experiences (d-AlexNet vs. AlexNet) by regional correspondence (the OFA versus  
291 the FFA). In addition, the main effect of the regional correspondence showed that the response  
292 profile of the face-selective channels in the DCNNs fitted better with the activation of the FFA than

293 that of the OFA ( $F(1,798) = 98.69, p = .001$ , partial  $\eta^2 = 0.11$ ), suggesting that the face-selective  
294 channels in DCNNs may in general prefer to process faces as a whole than face parts. Critically, the  
295 two-way interaction was significant ( $F(1,798) = 84.9, p < .001$ , partial  $\eta^2 = 0.10$ ), indicating that the  
296 experience affected the correspondence to the FFA and OFA disproportionately. A simple effect  
297 analysis revealed that the correspondence to the FFA (MD = 0.023,  $p < .001$ ) was increased by face-  
298 specific experiences to a significantly larger extent than that to the OFA (MD = 0.004,  $p = .013$ ,  
299 Figure 2c). Since the FFA is more involved in holistic processing of faces and the OFA is more  
300 dedicated to the part-based analysis, the disproportional decrease in correspondence between the  
301 face-selective channels in the d-AlexNet and the FFA implied that the role of the experience on faces  
302 was to facilitate the processing of faces as a whole.

303 To test this conjecture, we examined how the d-AlexNet responded to inverted faces, a  
304 behavioral signature of face-specific processing. As expected, there was a face inversion effect in the  
305 AlexNet's face-selective channels, with the magnitude of the activation to upright faces significantly  
306 larger than that to inverted faces ( $t(19) = 6.45, p < .001$ , Cohen's  $d = 1.44$ ) (Figure 2d). However, no  
307 inversion effect was observed in the d-AlexNet, as the magnitude of the activation to upright faces  
308 was not significantly larger than that to inverted faces ( $t(19) = 0.86, p = .40$ ). The lack of the  
309 inversion effect in the d-AlexNet was further supported by a two-way interaction of visual experience  
310 by orientation of faces,  $F(1, 19) = 7.79, p = .012$ , partial  $\eta^2 = 0.29$ . That is, unlike the AlexNet, the  
311 d-AlexNet processed upright faces in the same fashion as inverted faces.



312

313 Figure 2. (a) The category-wise activation profiles of example face-selective channels of the AlexNet  
 314 (left) and the d-AlexNet (right). The ‘faces’ shown here were AI-generated for illustration purpose  
 315 only. (b) The  $R^2$  maps of the regression with the activation of the d-AlexNet’s (right) or the  
 316 AlexNet’s face-selective channels (left) as the independent variables. The higher  $R^2$  in multiple  
 317 regression, the better correspondence between the face channels in the DCNNs and the face-selective  
 318 regions in the human brain. The crimson lines delineate the ROIs of the OFA and the FFA. (c) The  
 319 face-channels of both DCNNs corresponded better with the FFA than the OFA, and the difference  
 320 between the AlexNet and the d-AlexNet was larger in the FFA. (d) Face inversion effect. The  
 321 average activation amplitude of the top two face-selective channels differed in response to upright  
 322 and inverted faces in the AlexNet but not the d-AlexNet. The error bar denotes standard error. The  
 323 asterisk denotes statistical significance ( $\alpha = .05$ ). n.s. denotes no significance.

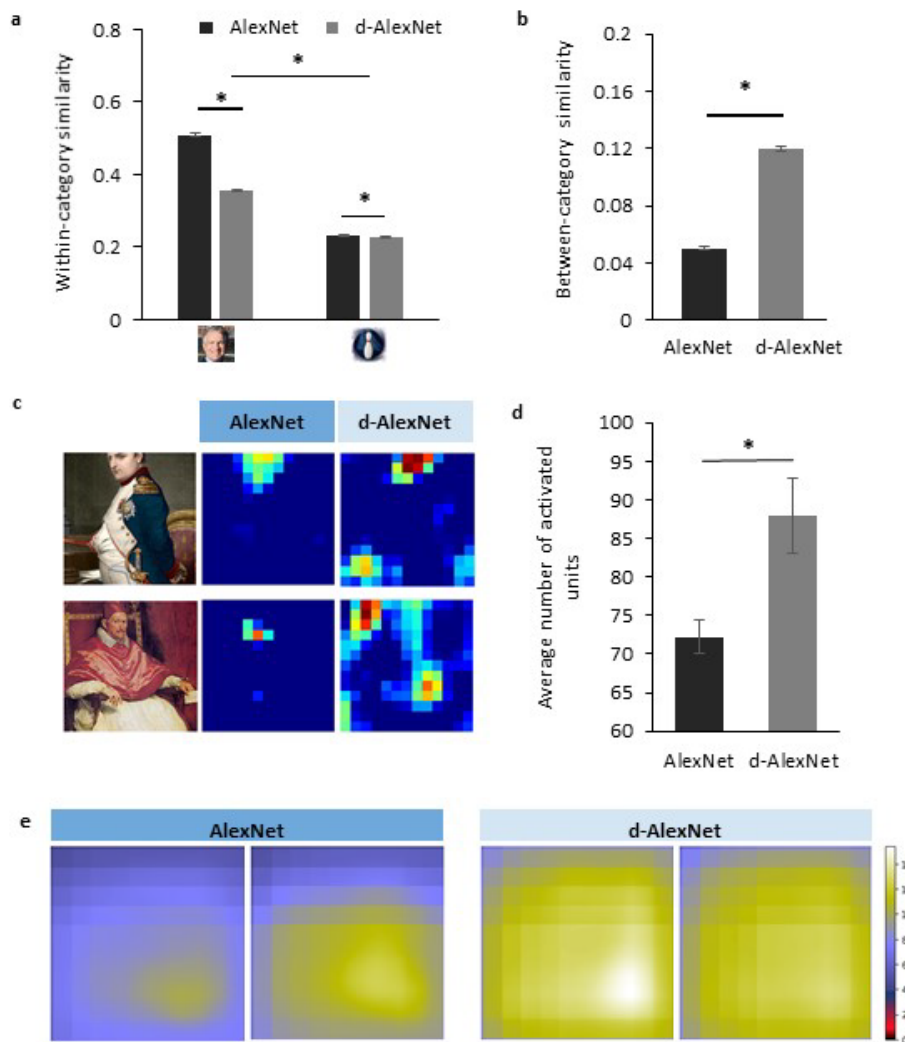
324

325 Previous studies on human suggested that inverted faces are processed in an object-like fashion.  
326 That is, it relies more on the parts-based analysis than the holistic processing. Therefore, we  
327 speculated that in the d-AlexNet faces were also represented more like non-face objects. To test this  
328 speculation, we first compared the representational similarity among responses in Conv5 to faces and  
329 bowling-pins, a novel object category that was not exposed to either DCNNs. As expected, the two-  
330 way interaction of experience (AlexNet versus d-AlexNet) by category (faces versus bowling-pins)  
331 was significant ( $F(1, 6,318) = 4,110.88, p < .001, \text{partial } \eta^2 = 0.39$ ), and the simple effect analysis  
332 suggested that the representation for faces in the AlexNet was more similar between each other than  
333 in the d-AlexNet ( $MD = 0.16, p < .001$ ), whereas the within-category representation similarity for  
334 bowling-pins showed the same but numerically smaller between-DCNN difference ( $MD = 0.005, p$   
335  $= .002$ ) (Figure 3a).

336 A more critical test was to examine how face-specific experiences made faces being processed  
337 differently from objects. Here we calculated between-category similarities between faces and  
338 bowling-pins. We found that the between-category similarity between faces and bowling-pins was  
339 significantly higher in the d-AlexNet than that in the AlexNet ( $t(3,159) = 42.42, MD = 0.07, p$   
340  $< .001, \text{Cohen's } d = 0.76$ ) (Figure 3b), suggesting that faces in the d-AlexNet were indeed  
341 represented more like objects. In short, although d-AlexNet was able to perform face tasks similar to  
342 the one with face-specific experiences, it represented faces in an object-like fashion.

343 Finally, we asked how faceness was achieved in DCNNs with face-specific experiences.  
344 Neurophysiological studies on monkeys demonstrate experience-associated sharpening of neural  
345 response, with fewer neurons activated after learning. Here we performed a similar analysis by  
346 measuring the number of non-zero units (i.e., units with above-zero activation) of the face-selective  
347 channels activated by natural images containing faces. As shown in the activation map (Figure 3c), a  
348 smaller number of units were activated by faces in the AlexNet than that in the d-AlexNet ( $t(19) =$   
349  $3.317, MD = 15.78, \text{Cohen's } d = 0.74$ ) (Figure 3d), suggesting that the experience on faces made the  
350 representation to faces sparser, and thus more effective. Another effect of visual experiences  
351 observed in neurophysiological studies is that experiences reduce the size of neurons' receptive field.  
352 Here we also mapped the empirical receptive field of the face-selective channels. Similarly, we found  
353 that the empirical receptive field of the AlexNet was smaller than that of the d-AlexNet. That is,

354 within the theoretical receptive field, the empirical receptive field of the face-selective channels in  
355 the AlexNet was tuned to focus on a smaller region by face-specific experiences (Figure 3e).



356

357 Figure 3. (a) The within-category similarity in the face category and an unseen non-face category  
358 (bowling pins) in the DCNNs. (b) The between-category similarity between faces and bowling pins.  
359 (c) The activation maps of a typical face-selective channel of each DCNN in responses to natural  
360 images containing faces. Each pixel denotes activation in one unit. The images shown here were  
361 historical portrait paintings downloaded from the Internet for illustration purpose only, and are  
362 different from the images used in this study. (d) The extent of activation of the face-selective  
363 channels of each DCNN in responses to natural images containing faces. (e) The empirical receptive  
364 fields of the face-selective channels of each DCNN. The error bar denotes standard error. The  
365 asterisk denotes statistical significance ( $\alpha = .05$ ).



## 366 4 Discussion

367 This study presented a DCNN model of selective visual deprivation of faces. We found that without  
368 genetic predisposition and face-specific visual experiences, DCNNs were still capable of face  
369 perception. In addition, face-selective channels were also present in the d-AlexNet, which  
370 corresponded to human face-selective regions. That is, the visual experience of faces was not  
371 necessary for an intelligent system to develop a face-selective module. On the other hand, besides the  
372 slightly compromised selectivity of the module, the deprivation led the d-AlexNet to process faces in  
373 a more parts-based fashion, similar to the way of processing objects. Indeed, face-inversion effect  
374 was absent in the d-AlexNet, and the representation of faces was more similar to objects as compared  
375 to the AlexNet. Finally, the functionality of face-specific experiences that led the AlexNet to process  
376 faces as a whole might be achieved by fine-tuning the sparse coding and the size of the receptive  
377 field of the face-selective channels. In sum, our study addressed a long-standing debate on nature  
378 versus nurture in developing the face-specific module, and illuminated the role of visual experiences  
379 in shaping the module.

380 The observation that without domain-specific visual experience, the face-selective processing  
381 and module still emerged in the d-AlexNet seems surprising; yet this finding is consistent with  
382 previous studies on non-human primates and new-born human infants (Bushneil, Sai, & Mullin,  
383 1989; Goren, Sarty, & Wu, 1975; Morton & Johnson, 1991; Sugita, 2008; Valenza, Simion, Cassia,  
384 & Umiltà, 1996), where the face-specific experience is found not necessary for face detection and  
385 recognition. However, the experience-independent face processing is largely attributed to either  
386 innate face-specific mechanisms (McKone et al., 2012; Morton & Johnson, 1991) or domain-general  
387 processing with predisposed biases (Cassia, Turati, & Simion, 2004; Simion & Di Giorgio, 2015;  
388 Simion, Macchi Cassia, Turati, & Valenza, 2001). Our study argued against this conjecture, because  
389 unlike any biological system, DCNNs have no predefined genetic inheritance or processing biases.  
390 Therefore, the face-specific processing observed in DCNNs had to derive from domain-specific  
391 factors.

392 We speculated that the face module in the d-AlexNet may result from a tremendous amount of  
393 features represented in the multiple layers of the network, with which face-like features were selected  
394 to construct face-specific module. In fact, previous studies on DCNNs have shown that DCNN's  
395 lower layers showed sensitivity to myriad visual features similar to primates' primary visual cortex  
396 (Krizhevsky et al., 2012), while the higher layers are tuned to complex features resembling those

397 represented in the ventral visual pathway (Güçlü & van Gerven, 2015; Khaligh-Razavi &  
398 Kriegeskorte, 2014; Pospisil, Pasupathy, & Bair, 2018; Yamins et al., 2014). With such repertoire of  
399 rich features, a representational space for faces, or for any natural object, may be constructed by  
400 selecting face-like features and features that are potentially useful in a variety of face tasks.

401 Supporting evidence for this conjecture came from the observation that the d-AlexNet processed  
402 faces in an object-like fashion. For example, the face inversion effect, a behavioral signature of face-  
403 specific processing in human (Kanwisher, Tong, & Nakayama, 1998; Rossion & Gauthier, 2002;  
404 Yin, 1969) was absent in the d-AlexNet. That is, similar to inverted faces, upright faces may also be  
405 processed like objects in the d-AlexNet. A more direct illustration of the object-like representation of  
406 faces came from the analysis on the representational similarity between faces and objects. As  
407 compared to the AlexNet, faces in the representational space of the d-AlexNet were less congregated  
408 among each other; instead they were more intermingled with non-face object categories. The finding  
409 that face representation was no longer qualitatively different from object representation may help  
410 explaining the performance of the d-AlexNet. Because faces were less segregated from objects in the  
411 representational space, the d-AlexNet's accuracy of face categorization was worse than that of the  
412 AlexNet. In contrast, within the face category, individual faces were less congregated in the  
413 representational space; therefore, the discrimination of individual faces became easier instead,  
414 suggested by the slightly higher face discrimination accuracy in the d-AlexNet than the AlexNet. In  
415 short, when the representational space of the d-AlexNet was formed exclusively based on features  
416 from non-face stimuli, faces were represented no longer qualitatively different from non-face objects,  
417 which inevitably led to 'object-like' face processing.

418 The face-specific processing is likely achieved through prior exposure to faces. At first glance,  
419 the effect of face-specific experiences seemed quantitative, as in the AlexNet, both the selectivity to  
420 faces and the number of the face-selective channels were increased, and the correspondence between  
421 the face-selective channels and the face-selective regions in human brain was tighter. However,  
422 careful scrutiny of the difference between the two DCNNs revealed that the changes led by the  
423 experience may be qualitative. For example, the deprivation of visual experiences disproportionately  
424 weakened the DCNN-brain correspondence in the FFA as comparing to the OFA, and the FFA is  
425 engaged more in the configural processing and the OFA in parts-based analysis (Liu, Harris, &  
426 Kanwisher, 2010; Nichols, Betts, & Wilson, 2010; Zhao et al., 2014). Therefore, the 'face-like' face  
427 processing may come from the fact that face-specific experiences led the representation of faces more

428 congregated within face category and more separable from the representation of non-face objects  
429 stimuli (see also Gomez, Barnett, & Grill-Spector, 2019). In this way, a relative encapsulated  
430 representation may help developing a unique way of processing faces, qualitatively different from  
431 non-face objects.

432         The advantage of the computational transparency of DCNNs may shed light on the  
433 development of domain specificity of the face module. First, we found that face-specific experiences  
434 increased the sparseness of face representation, as fewer units of the face channels were activated by  
435 faces in the AlexNet. The experience-dependent sparse coding has been widely discovered in the  
436 visual cortex, such as the V4, MT, and IT (for reviews, see Desimone, 1996; Grill-Spector, Henson,  
437 & Martin, 2006; Wiggs & Martin, 1998). The experience-induced increase of sparseness is thought to  
438 reflect a preference-narrowing process that tunes neurons to a smaller range of stimuli (Kohn &  
439 Movshon, 2004); therefore, with sparse coding faces are less likely to be intermingled with non-face  
440 objects, which may lead to more congregated representations in the representational space in the  
441 AlexNet, as compared to the d-AlexNet. Second, we found that the empirical receptive field of the  
442 face channel in the AlexNet was smaller than that in the d-AlexNet, suggesting that the visual  
443 experience on faces decreased the size of the receptive field of the face channels. This finding fits  
444 perfectly with neurophysiological studies that the size of receptive fields of visual neurons is reduced  
445 after eye-opening (Braastad & Heggelund, 1985; Cantrell, Cang, Troy, & Liu, 2010; Koehler,  
446 Akimov, & Renteria, 2011; Tavazoie & Reid, 2000). Importantly, along with the refined receptive  
447 fields, the selectivity of neurons increases (Spilmann, 2014), possibly because neurons can avoid  
448 distracting information by focusing on a more restricted part of stimuli, which may further allowed  
449 finer representation of the selected regions. This is especially important for processing faces because  
450 faces are highly homogeneous, and some information is identical across faces, such as parts  
451 composition (eyes, noses, and mouth) and their configural arrangements. Therefore, the reduced  
452 receptive field of the face channels may facilitate selective analyses of discriminative face features  
453 while avoiding irrelevant information. Further, the sharpening of the receptive field and the fine-  
454 tuned selectivity may result in superior discrimination ability on faces, and allow faces to be  
455 processed at the sub-ordinate level (i.e., identification), whereas the rest of objects are largely  
456 processed at the basic level (i.e., categorization).

457         It has long been assumed that domain-specific visual experiences and inheritance are the pre-  
458 requisites in the development of the face module. In our study with DCNNs as a model, we

459 completely decoupled the genetic predisposition and face-specific visual experiences, and found that  
460 the representation for faces can be constructed with features from non-face objects to realize basic  
461 functions for face recognition. Therefore, in many situations, the difference between faces and  
462 objects is ‘quantitative’ rather than ‘qualitative’, as they are represented in a continuum of the  
463 representational space. In addition, we also found that face-specific experiences likely fine-tuned the  
464 face representation, and thus transformed the ‘object-like’ face processing into ‘face-specific’  
465 processing. However, we shall be cautious that our finding may not be applicable for the  
466 development of face module in human, as in the biological brain experience-induced changes are  
467 partly attributed to the inhibition from lateral connections (Grill-Spector et al., 2006; Norman &  
468 O'Reilly, 2003), whereas there is no lateral or feedback connection in DCNNs. However, despite  
469 structural differences, recent studies have shown similar representation for faces between DCNNs  
470 and humans (Song, Qu, Xu, & Liu, 2020), suggesting that a common mechanism may be shared by  
471 both artificial and biological intelligent systems. Future studies are needed to examine the  
472 applicability of our finding to humans. On the other hand, our study illustrated the advantages of  
473 using DCNNs as a model to understand human mind because of its computational transparency and  
474 its dissociation of factors in nature and nurture. Thus, our study invites future studies with DCNNs to  
475 understand the development of domain specificity in particular and a broad range of cognitive  
476 modules in general.

477

## 478 **5 Conflict of Interest**

479 The authors declare no competing interests.

## 480 **6 Author Contributions**

481 J. L. conceived and designed the study. Y.Z. analyzed the data with input from all authors. S.X. wrote the  
482 manuscript with input from J. L., Y.Z. and Z. Z.

## 483 **7 Funding**

484 This study was funded by the National Natural Science Foundation of China (31861143039, 31771251,  
485 and 31600925), the Fundamental Research Funds for the Central Universities, and the National Basic  
486 Research Program of China (2018YFC0810602).

## 487 **8 Data Availability Statement**

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

## 490 **9 Reference**

- 491 AbdAlmageed, W., Wu, Y., Rawls, S., Harel, S., Hassner, T., Masi, I., et al. (2016). Face  
492 Recognition Using Deep Multi-Pose Representations. In *2016 Ieee Winter Conference on*  
493 *Applications of Computer Vision*.
- 494 Arcaro, M. J., Schade, P. F., Vincent, J. L., Ponce, C. R., & Livingstone, M. S. (2017). Seeing faces  
495 is necessary for face-domain formation. *Nature Neuroscience*, *20*(10), 1404-+.
- 496 Baek, S., Song, M., Jang, J., Kim, G., & Paik, S.-B. (2019). Spontaneous generation of face  
497 recognition in untrained deep neural networks. *bioRxiv*, 857466.
- 498 Berg, T. L., Berg, A. C., Edwards, J., & Forsyth, D. A. (2005). *Who's in the picture*. Paper presented  
499 at the Advances in neural information processing systems.
- 500 Braastad, B. O., & Heggelund, P. (1985). Development of spatial receptive-field organization and  
501 orientation selectivity in kitten striate cortex. *Journal of Neurophysiology*, *53*(5), 1158-1178.
- 502 Bushneil, I., Sai, F., & Mullin, J. (1989). Neonatal recognition of the mother's face. *British Journal of*  
503 *Developmental Psychology*, *7*(1), 3-15.
- 504 Cantrell, D. R., Cang, J., Troy, J. B., & Liu, X. (2010). Non-Centered Spike-Triggered Covariance  
505 Analysis Reveals Neurotrophin-3 as a Developmental Regulator of Receptive Field Properties  
506 of ON-OFF Retinal Ganglion Cells. *Plos Computational Biology*, *6*(10).
- 507 Cassia, V. M., Turati, C., & Simion, F. (2004). Can a nonspecific bias toward top-heavy patterns  
508 explain newborns' face preference? *Psychological Science*, *15*(6), 379-383.
- 509 Chen, X., Zhou, M., Gong, Z., Xu, W., Liu, X., Huang, T., et al. (2020). DNNBrain: a unifying  
510 toolbox for mapping deep neural networks and brains. *bioRxiv*, 2020.2007.2005.188847.
- 511 Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., & Fei-Fei, L. (2009). *Imagenet: A large-scale*  
512 *hierarchical image database*. Paper presented at the 2009 IEEE conference on computer  
513 vision and pattern recognition.
- 514 Desimone, R. (1996). Neural mechanisms for visual memory and their role in attention. *Proceedings*  
515 *of the National Academy of Sciences of the United States of America*, *93*(24), 13494-13499.
- 516 Freiwald, W., Duchaine, B., & Yovel, G. (2016). Face Processing Systems: From Neurons to Real-  
517 World Social Perception. In S. E. Hyman (Ed.), *Annual Review of Neuroscience*, *Vol 39* (Vol.  
518 39, pp. 325-346).
- 519 Gomez, J., Barnett, M., & Grill-Spector, K. (2019). Extensive childhood experience with Pokemon  
520 suggests eccentricity drives organization of visual cortex. *Nature Human Behaviour*, *3*(6),  
521 611-624.
- 522 Goren, C. C., Sarty, M., & Wu, P. Y. (1975). Visual following and pattern discrimination of face-like  
523 stimuli by newborn infants. *Pediatrics*, *56*(4), 544-549.

- 524 Griffin, G., Holub, A., & Perona, P. (2007). Caltech-256 object category dataset.
- 525 Grill-Spector, K., Henson, R., & Martin, A. (2006). Repetition and the brain: neural models of  
526 stimulus-specific effects. *Trends in Cognitive Sciences*, *10*(1), 14-23.
- 527 Grundstrom, J., Chen, J., Ljungqvist, M. G., & Astrom, K. (2016). Transferring and Compressing  
528 Convolutional Neural Networks for Face Representations. In A. Campilho & F. Karray  
529 (Eds.), *Image Analysis and Recognition* (Vol. 9730, pp. 20-29).
- 530 Güçlü, U., & van Gerven, M. A. (2015). Deep neural networks reveal a gradient in the complexity of  
531 neural representations across the ventral stream. *Journal of Neuroscience*, *35*(27), 10005-  
532 10014.
- 533 Kanwisher, N., Tong, F., & Nakayama, K. (1998). The effect of face inversion on the human  
534 fusiform face area. *Cognition*, *68*(1), B1-B11.
- 535 Kanwisher, N., & Yovel, G. (2006). The fusiform face area: a cortical region specialized for the  
536 perception of faces. *Philosophical Transactions of the Royal Society B-Biological Sciences*,  
537 *361*(1476), 2109-2128.
- 538 Khaligh-Razavi, S.-M., & Kriegeskorte, N. (2014). Deep supervised, but not unsupervised, models  
539 may explain IT cortical representation. *PLoS computational biology*, *10*(11), e1003915.
- 540 Koehler, C. L., Akimov, N. P., & Renteria, R. C. (2011). Receptive field center size decreases and  
541 firing properties mature in ON and OFF retinal ganglion cells after eye opening in the mouse.  
542 *Journal of Neurophysiology*, *106*(2), 895-904.
- 543 Kohn, A., & Movshon, J. A. (2004). Adaptation changes the direction tuning of macaque MT  
544 neurons. *Nature Neuroscience*, *7*(7), 764-772.
- 545 Kriegeskorte, N. (2015). Deep neural networks: a new framework for modeling biological vision and  
546 brain information processing. *Annual review of vision science*, *1*, 417-446.
- 547 Krizhevsky, A. (2014). One weird trick for parallelizing convolutional neural networks. *arXiv*  
548 *preprint arXiv:1404.5997*.
- 549 Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). *Imagenet classification with deep*  
550 *convolutional neural networks*. Paper presented at the Advances in neural information  
551 processing systems.
- 552 Liu, J., Harris, A., & Kanwisher, N. (2010). Perception of Face Parts and Face Configurations: An  
553 fMRI Study. *Journal of Cognitive Neuroscience*, *22*(1), 203-211.
- 554 McKone, E., Crookes, K., Jeffery, L., & Dilks, D. D. (2012). A critical review of the development of  
555 face recognition: Experience is less important than previously believed. *Cognitive*  
556 *Neuropsychology*, *29*(1-2), 174-212.
- 557 Morton, J., & Johnson, M. H. (1991). CONSPEC and CONLERN: a two-process theory of infant  
558 face recognition. *Psychological Review*, *98*(2), 164.
- 559 Nichols, D. F., Betts, L. R., & Wilson, H. R. (2010). Decoding of faces and face components in face-  
560 sensitive human visual cortex. *Frontiers in Psychology*, *1*.
- 561 Norman, K. A., & O'Reilly, R. C. (2003). Modeling hippocampal and neocortical contributions to  
562 recognition memory: A complementary-learning-systems approach. *Psychological Review*,  
563 *110*(4), 611-646.

- 564 Parkhi, O. M., Vedaldi, A., & Zisserman, A. (2015). *Deep face recognition*. Paper presented at the  
565 bmvc.
- 566 Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., et al. (2017). Automatic  
567 differentiation in pytorch.
- 568 Pospisil, D. A., Pasupathy, A., & Bair, W. (2018). 'Artiphysiology' reveals V4-like shape tuning in a  
569 deep network trained for image classification. *Elife*, 7, e38242.
- 570 Rossion, B., & Gauthier, I. (2002). How does the brain process upright and inverted faces?  
571 *Behavioral and cognitive neuroscience reviews*, 1(1), 63-75.
- 572 Simion, F., & Di Giorgio, E. (2015). Face perception and processing in early infancy: inborn  
573 predispositions and developmental changes. *Frontiers in Psychology*, 6.
- 574 Simion, F., Macchi Cassia, V., Turati, C., & Valenza, E. (2001). The origins of face perception:  
575 specific versus non - specific mechanisms. *Infant and Child Development: An International  
576 Journal of Research and Practice*, 10(1 - 2), 59-65.
- 577 Song, Y., Qu, Y., Xu, S., & Liu, J. (2020). Implementation-independent representation for deep  
578 convolutional neural networks and humans in processing faces. *bioRxiv*.
- 579 Spilmann, L. (2014). Receptive fields of visual neurons: The early years. *Perception*, 43(11), 1145-  
580 1176.
- 581 Sugita, Y. (2008). Face perception in monkeys reared with no exposure to faces. *Proceedings of the  
582 National Academy of Sciences of the United States of America*, 105(1), 394-398.
- 583 Tavazoie, S. F., & Reid, R. C. (2000). Diverse receptive fields in the lateral geniculate nucleus during  
584 thalamocortical development. *Nature Neuroscience*, 3(6), 608-616.
- 585 Valenza, E., Simion, F., Cassia, V. M., & Umiltà, C. (1996). Face preference at birth. *Journal of  
586 experimental psychology: Human Perception and Performance*, 22(4), 892.
- 587 VanRullen, R., & Reddy, L. (2019). Reconstructing faces from fMRI patterns using deep generative  
588 neural networks. *Communications biology*, 2(1), 193-193.
- 589 Wen, H., Shi, J., Zhang, Y., Lu, K.-H., Cao, J., & Liu, Z. (2017). Neural encoding and decoding with  
590 deep learning for dynamic natural vision. *Cerebral Cortex*, 28(12), 4136-4160.
- 591 Wiggs, C. L., & Martin, A. (1998). Properties and mechanisms of perceptual priming. *Current  
592 Opinion in Neurobiology*, 8(2), 227-233.
- 593 Wilmer, J. B., Germine, L., Chabris, C. F., Chatterjee, G., Williams, M., Loken, E., et al. (2010).  
594 Human face recognition ability is specific and highly heritable. *Proceedings of the National  
595 Academy of Sciences of the United States of America*, 107(11), 5238-5241.
- 596 Yamins, D. L., Hong, H., Cadieu, C. F., Solomon, E. A., Seibert, D., & DiCarlo, J. J. (2014).  
597 Performance-optimized hierarchical models predict neural responses in higher visual cortex.  
598 *Proceedings of the National Academy of Sciences*, 111(23), 8619-8624.
- 599 Yi, D., Lei, Z., Liao, S., & Li, S. Z. (2014). Learning face representation from scratch. *arXiv preprint  
600 arXiv:1411.7923*.
- 601 Yin, R. K. (1969). Looking at upside-down faces. *Journal of Experimental Psychology*, 81(1), 141-  
602 &.

- 603 Zhao, M., Cheung, S.-h., Wong, A. C. N., Rhodes, G., Chan, E. K. S., Chan, W. W. L., et al. (2014).  
604 Processing of configural and componential information in face-selective cortical areas.  
605 *Cognitive Neuroscience*, 5(3-4), 160-167.
- 606 Zhen, Z., Yang, Z., Huang, L., Kong, X.-z., Wang, X., Dang, X., et al. (2015). Quantifying  
607 interindividual variability and asymmetry of face-selective regions: a probabilistic functional  
608 atlas. *Neuroimage*, 113, 13-25.
- 609 Zhu, Q., Song, Y., Hu, S., Li, X., Tian, M., Zhen, Z., et al. (2010). Heritability of the Specific  
610 Cognitive Ability of Face Perception. *Current Biology*, 20(2), 137-142.

611

612