# The Face Module Emerged in a Deep Convolutional Neural Network Selectively Deprived of Face Experience

Shan Xu[1]*[†], Yiyuan Zhang[1†], Zonglei Zhen[1] and Jia Liu[2]*

[1] Beijing Key Laboratory of Applied Experimental Psychology, Faculty of Psychology, Beijing Normal University, Beijing, China,
[2] Department of Psychology & Tsinghua Laboratory of Brain and Intelligence, Tsinghua University, Beijing, China

Can we recognize faces with zero experience on faces? This question is critical because it examines the role of experiences in the formation of domain-specific modules in the brain. Investigation with humans and non-human animals on this issue cannot easily dissociate the effect of the visual experience from that of the hardwired domain-specificity. Therefore, the present study built a model of selective deprivation of the experience on faces with a representative deep convolutional neural network, AlexNet, by removing all images containing faces from its training stimuli. This model did not show significant deficits in face categorization and discrimination, and face-selective modules automatically emerged. However, the deprivation reduced the domain-specificity of the face module. In sum, our study provides empirical evidence on the role of nature vs. nurture in developing the domain-specific modules that domain-specificity may evolve from non-specific experience without genetic predisposition, and is further fine-tuned by domain-specific experience.

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

## INTRODUCTION

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

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

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

Recent advances in deep convolutional neural network (DCNN) provide an ideal test platform to examine the impact of visual experiences on face modules without genetic predisposition. DCNNs are found similar to human visual cortex both structurally and functionally (Kriegeskorte, 2015), but free of any predisposition on functional modules. Therefore, with DCNNs we can manipulate experiences without considering interactions from genetic predisposition. In this study, we asked whether DCNNs can achieve face-specific recognition ability when visual experiences on faces were selectively deprived.

To do this, we trained a representative DCNN, AlexNet (Krizhevsky et al., 2012), to categorize non-face objects with face images carefully removed from the training dataset. Once this face-deprived DCNN (d-AlexNet) was trained, we compared its behavioral performance to that of a normal AlexNet of the same architecture but with faces present during training. Specifically, we examined their performance in both face categorization (i.e., differentiating faces from non-face objects) and discrimination (i.e., discriminating faces among different individuals) tasks. We predicted that the d-AlexNet, though without predisposition and experiences of faces, may still develop face selectivity through its visual experiences of non-face objects.

## MATERIALS AND METHODS

### Stimuli

#### Deprivation Dataset

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

### Classification Dataset

To train a classifier that can classify faces, we constructed a classification dataset consisting of 204 categories of non-face objects and one face category, each of 80 exemplars. For the non-face categories, we manually screened Caltech-256 (Griffin et al., 2007) to remove images containing human, primate, or cartoon faces, and then removed categories whose remaining images were <80. In each of the 204 remaining non-face categories, we randomly chose 70 images for training and another 10 for calculating classification accuracy. The face category was constructed by randomly selecting 1,000 faces images from Faces in the Wild (FITW) dataset (Berg et al., 2005). Among them, 70 were used as training data and another 10 for classification accuracy. In addition, to characterize DCNN's ability in differentiating faces from object categories, we compiled a second dataset consisting of all images in the face category except those used in training.

### Discrimination Dataset

To train a classifier that can discriminate faces at individual level, we constructed a discrimination dataset consisting of face images of 133 individuals, 300 images each, selected from the Casia-WebFace database (Yi et al., 2014). For each individual in the dataset, 250 were randomly chosen for training and another 50 for calculating discrimination accuracy.
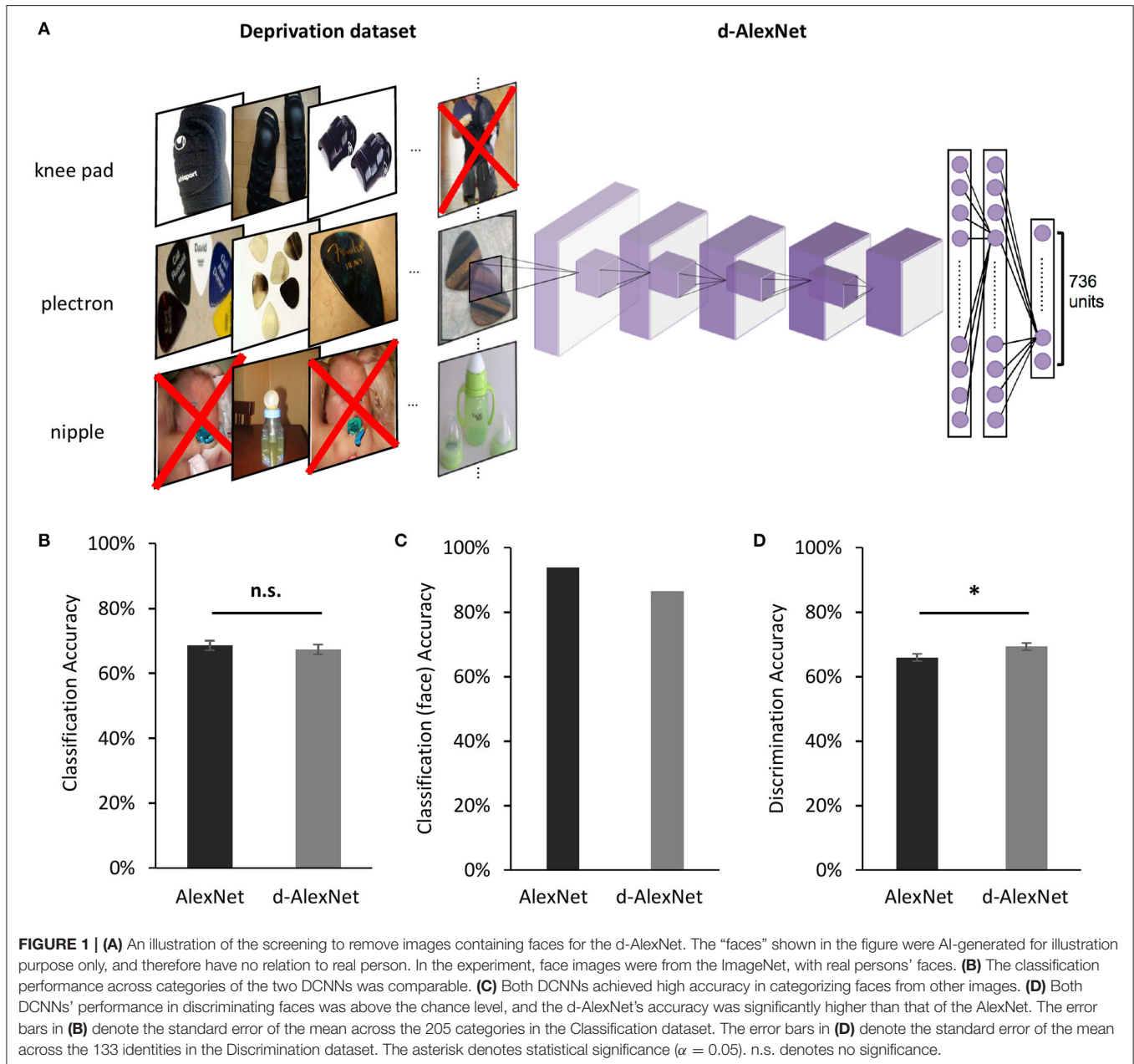
### Representation Dataset

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

### Movies Clips for DCNN-Brain Correspondence Analysis

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

## The Deep Convolutional Neural Network

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

**FIGURE 1 | (A)** An illustration of the screening to remove images containing faces for the d-AlexNet. The "faces" shown in the figure were AI-generated for illustration purpose only, and therefore have no relation to real person. In the experiment, face images were from the ImageNet, with real persons' faces. **(B)** The classification performance across categories of the two DCNNs was comparable. **(C)** Both DCNNs achieved high accuracy in categorizing faces from other images. **(D)** Both DCNNs' performance in discriminating faces was above the chance level, and the d-AlexNet's accuracy was significantly higher than that of the AlexNet. The error bars in **(B)** denote the standard error of the mean across the 205 categories in the Classification dataset. The error bars in **(D)** denote the standard error of the mean across the 133 identities in the Discrimination dataset. The asterisk denotes statistical significance ($\alpha = 0.05$). n.s. denotes no significance.

layers. FC3 is followed by a sublayer using a softmax function to output a vector that represents the probability of the visual input containing the corresponding object category (Krizhevsky et al., 2012).

The d-AlexNet used the architecture of AlexNet but changed the number of units in FC3 to 736 and changed the following softmax function accordingly to match the number of categories in the deprivation dataset. The d-AlexNet was initialized with values drawn from a uniform distribution, and was then trained on the deprivation dataset following the approach specified in Krizhevsky (2014). We used the pre-trained AlexNet from pytorch 1.2.0 as the normal DCNN, referred to as the AlexNet in this paper for brevity.

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

To test the generalizability of the main findings of the present study, we also applied the same deprivation on another well-known DCNN, "ResNet-18" (He et al., 2016). ResNet-18 introduces residual learning blocks in a DCNN to overcome the degradation problem in the training of DCNNs, and achieves even better performance than AlexNet in object categorization task with a deeper architecture. The d-ResNet used the architecture of ResNet-18 but changed the number of units in the FC layer to 736 and changed the following softmax function

accordingly to match the number of categories in the deprivation dataset. The d-ResNet was trained on the deprivation dataset following the same approach specified above. For comparison, we used the pre-trained ResNet-18 from pytorch 1.2.0 as the normal DCNN, referred to as the ResNet in this study for brevity.

## Transfer Learning for Classification and Discrimination

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

To examine to what extent our manipulation of the visual experience affected face discrimination, we similarly replaced the fully connected layers of each DCNN with a discrimination classifier. The discrimination classifier differed from the classification classifier only in its second layer, which had 133 units instead as outputs, each corresponding to one individual in the discrimination dataset. The face-discrimination classifier was trained for each DCNN with the training images in the discrimination dataset for 90 epochs. The same transfer learning was applied to the d-ResNet and the pre-trained ResNet-18.

## The Face Selective Channels in DCNNs

To identify the channels selectively responsive to faces, we submitted images in the classification dataset to each DCNN, recorded the average activation in each channel of Conv5 after ReLU in response to each image, and then averaged the channel-wise activation within each category. We selected channels where the face category evoked the highest activation, and used the Mann-Whitney U test to examine the activation difference between faces and objects that had the second-highest activation in these channels ($p < 0.05$, Bonferroni corrected). The selectivity of each face channel thus identified was indexed by the selective ratio. The selective ratio was calculated by dividing the face activation by the second-highest activation. In addition, we measured the lifetime sparseness of each face-selective channel as an index for selectivity of faces among all non-face objects. We first normalized the mean activations of a face channel in Layer5 to all the categories to the range of 0–1, and then calculated lifetime sparseness with the formula:

$$S = \frac{\left(\sum_{i=1,n} r_i/n\right)^2}{\sum_{i=1,n} \left(r_i{}^2/n\right)}$$

where $r_i$ is the normalized activations to the ith object category. The smaller this value is, the higher the selectivity is.

To confirm the face selectivity of the selected channels, we also tested their categorical selectivity with the fMRI localizer stimuli typically used to identify face-selective regions. More specifically, we recorded each channels' responses to the localizer stimuli from the face and the tool condition of the Human Connectome Project dataset (Van Essen et al., 2013), and examined the significance of face selectivity of each face channel by comparing the activation in the face condition and that of the tool condition in this channel using the Mann-Whitney U test described above.

Since we found face-selective channels in the d-AlexNet and reduced face selectivity of these channels comparing with face-selective channels in the AlexNet, we proceeded to test the robustness of these findings. Another five instances of face-deprived AlexNet were each independently trained in the same way as the d-AlexNet. In these instances, we searched for face-selective channels, computed their face selectivity, and examined the significance of their face selectivity by the Mann-Whitney U test on their responses to the classification dataset as well as on the fMRI localizer stimuli, in the same way as we did in the d-AlexNet and the AlexNet. The same procedure of channel identification was also applied to the d-ResNet and the pre-trained ResNet-18.

## DCNN-Brain Correspondence

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

Neural activation in the brain was derived from the preprocessed data in Wen et al. (2017). The fMRI data were recorded while human participants viewed each movie clips twice. We averaged the standardized time series across repetition and across subjects for each clip. Then, for each DCNN, we conducted multiple regression for each clip, with the activation time series of each brain vertex as the dependent variable and that of face-selective channels in this network as independent variables. For the d-AlexNet, all face-selective channels were included. For the AlexNet, we included the same number of face-selective channels with the highest face selectivity to match the complexity of the regression model. We used the $R^2$ of each vertex as the index of the overall Goodness of fit of the regression in that vertex. The $R^2$ values were then averaged across clips. The larger the $R^2$ value, the higher correspondence between the DCNN and the brain in response to movie clips.

To test whether the correspondence changes between networks reflected an overall increase in the correspondence between fMRI signal and the activation of the face channels of the AlexNet comparing with the d-AlexNet (in contrast to an increase selectively within the face-selective regions), we delineated the face-selective regions and the object-selective regions and compared the correspondence between the top two

face channels of each network and the face- and the object-selective regions. The face- and the object- selective regions were defined by functional localizer data of Human Connectome Project (Van Essen et al., 2013). Two hundred vertexes of the highest Z value in the tool-avg contrast were delineated as the object-selective ROIs, and two hundred vertexes of the highest Z value in the face-tool contrast were delineated as the face-selective ROIs. The channel-brain correspondence of each vertex with the ROIs was indexed by $R^2$ of the regression with the fMRI time series of this vertex as the dependent variable and the time series of the top-two face channels as the independent variables. A two-way ANOVA with visual experiences (d-AlexNet vs. AlexNet) and categorical selectivity (the object-selective regions vs. the face-selective regions) as independent variables was conducted to examine the difference between the channel-brain correspondence between the categorical-selective regions and the face-selective channels of the d-AlexNet and the AlexNet.

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

## Face Inversion Effect in DCNNs

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

## Representational Similarity Analysis

To examine whether faces in the d-AlexNet were processed in an object-like fashion, we compared the within-category representational similarity of faces to that of bowling pins, an "unseen" non-face object category never exposed to either DCNN. Specifically, for each image in the representation dataset, we arranged the average activations of each channel of Conv5 after ReLU into vectors, and then for each pair of images we calculated and then Fisher-z transformed the correlation between their vectors, which served as an index of pairwise representational similarity. Within-category similarity between pairs of face images and that between pairs of object images

were calculated separately. A $2 \times 2$ ANOVA was conducted with visual experience (d-AlexNet vs. AlexNet) and category (face vs. object) as independent factors. In addition, cross-category similarity between faces and bowling pins was also calculated for each DCNN, and a paired sample t-test (two-tailed) on two DCNNs was conducted.

## Sparse Coding and Empirical Receptive Field

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

We also calculated the size of the empirical receptive field of the face-selective channels. Specifically, we obtained the activation maps of 1,000 images randomly chosen from FITW. Using the toolbox DNNBrain (Chen et al., 2020), we up-sampled each activation map to the same size of the input. For each image, we averaged the up-sampled activation within the theoretical receptive field of each unit (the part of the image covered by the convolution of this unit and the preceding computation, decided by the network architecture), and selected the unit with the highest average activation. We then cropped the up-sampled activation map by the theoretical receptive field of this unit, to locate the image part that activated this channel most across all the units. Then, we averaged corresponding cropped activation maps across all the face images, and the resultant map denotes the empirical receptive field of this channel, delineating the part of the theoretical receptive field that causes this channel to respond strongly in viewing its preferred stimuli.

## RESULTS

The d-AlexNet was trained with a dataset of 662,619 non-face images consisting of 736 non-face categories, generated by removing images containing faces from the ILSVRC 2012 dataset (**Figure 1A**). The d-AlexNet was initialized and trained in the same way as the AlexNet. Both networks were trained following the approach specified in Krizhevsky (2014). The resultant top-1 accuracy (57.29%) and the top-5 accuracy (80.11%) were comparable with the pre-trained AlexNet.

We first examined the performance of the d-AlexNet in two representative tasks of face processing, face categorization (i.e., differentiating faces from non-face objects) and face discrimination (i.e., identifying different individuals). The output of Conv5 after ReLU of the d-AlexNet was used to classify objects in the classification dataset (see Materials and Methods). The averaged categorization accuracy of the d-AlexNet (67.40%) was well above the chance level (0.49%), and comparable to that in the AlexNet [68.60%, $t_{(204)} = 1.26$, $p = 0.209$, Cohen's $d = 0.007$, **Figure 1B**]. Critically, the d-AlexNet, although with no experience on faces, succeeded in the face categorization task,

with an accuracy of 86.50% in categorizing faces from non-face objects. Note that the accuracy was numerically smaller than the AlexNet's accuracy in categorizing faces (93.90%) though (**Figure 1C**).
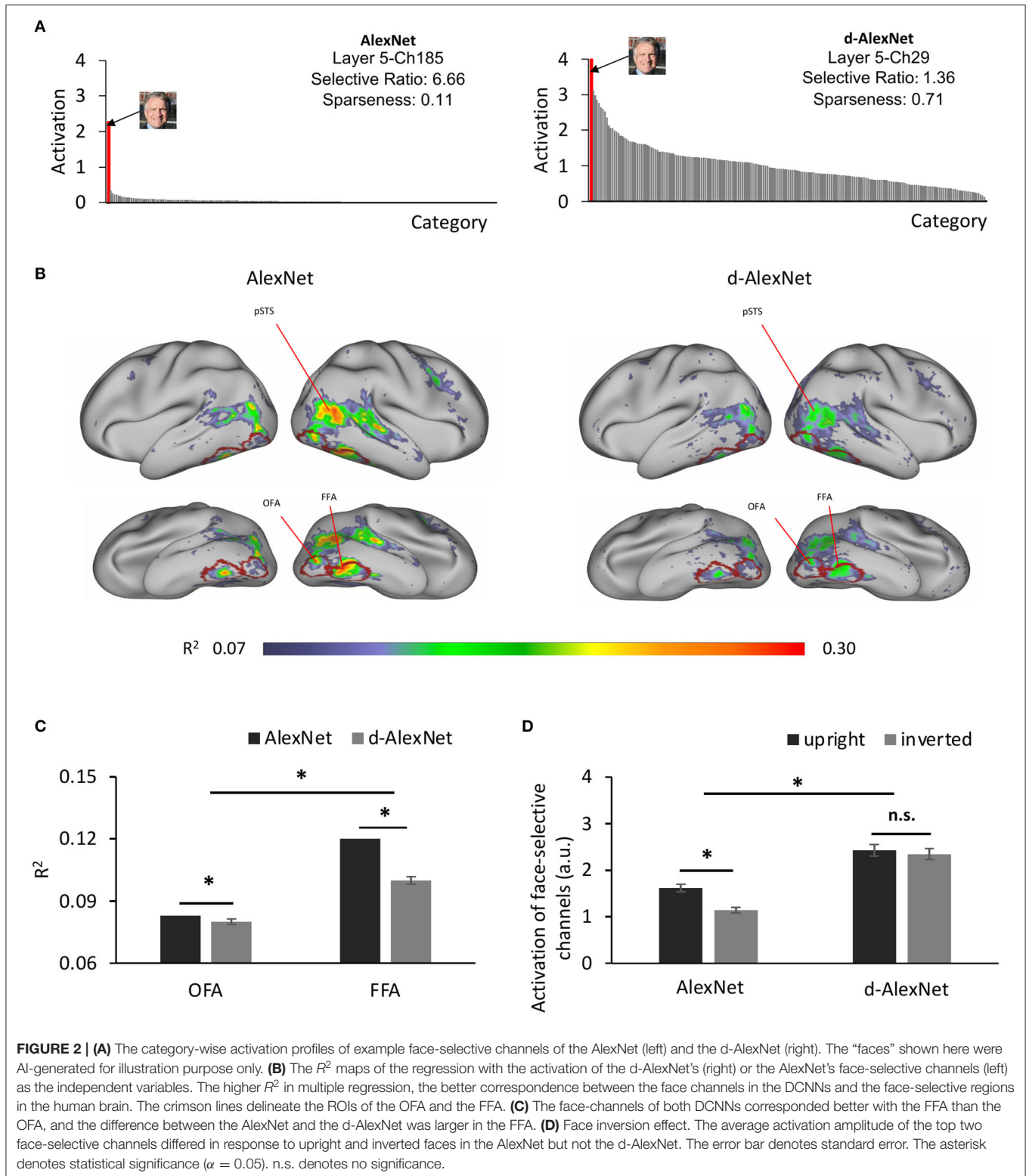
A similar pattern was observed in the face discrimination task. In this task, the output of Conv5 after ReLU of each DCNN was used to identify 33,250 face images into 133 identities in the discrimination dataset (see section Materials and Methods). As expected, the AlexNet was capable of face discrimination (65.9%), well above the chance level (0.75%), consistent with previous studies (AbdAlmageed et al., 2016; Grundstrom et al., 2016). Critically, the d-AlexNet also showed the capability of discriminating faces, with an accuracy of 69.30% that was even significantly higher than that of the AlexNet, $t_{(132)} = 3.16$, $p = 0.002$, Cohen's $d = 0.20$, (**Figure 1D**). Taken together, visual experiences on faces seemed not necessary for developing basic functions of processing faces.

Was a face module formed in the d-AlexNet to support these functions? To answer this question, we searched all the channels in Conv5 of the d-AlexNet, where face-selective channels have been previously identified in the AlexNet (Baek et al., 2019). To do this, we calculated the activation of each channel in Conv5 after ReLU in response to each category of the classification dataset, and then identified channels that showed significantly higher response to faces than non-face images with Mann-Whitney U test ($ps < 0.05$, Bonferroni corrected). Two face-selective channels (Ch29 and Ch50) met this criterion in the d-AlexNet (for an example channel, see **Figure 2A**, right), whereas four face-selective channels (Ch185, Ch125, Ch60, and Ch187) were identified in the AlexNet (for an example channel, see **Figure 2A**, left). The face-selective channels in two DCNNs differed in selectivity. The averaged selective ratio, the ratio of the activation magnitude to faces by that to the most activated non-face object category, was 1.29 (range: 1.22–1.36) in the d-AlexNet, much lower than that in the AlexNet (average ratio: 3.63, range: 1.43–6.66). The lifetime sparseness, which measures the breadth of tuning of a channel in response to a set of categories, also showed a similar result. The average lifetime sparseness index of the face channels in the AlexNet (mean = 0.25, range: 0.11–0.51) was smaller than that in the d-AlexNet (mean = 0.71, range: 0.70–0.71), indicating higher face selectivity in the AlexNet than that in the d-AlexNet. To confirm that the emergence of the face-selective channels in the d-AlexNet was not because of chance factors in network training, another five instances of face-deprived networks were independently initiated and trained respectively. One or two face-selective channels emerged in each of these face-deprived network instances, though the level of face selectivity was lower as compared to the AlexNet. In addition, we tested the face selectivity of the face channels in all face-deprived networks with the stimuli used to localize face-selective regions in fMRI studies, and found that the responses in these face-selective channels were significantly higher to the faces than those to the objects (Mann-Whitney U test, $ps < 0.05$, Bonferroni corrected). Taken together, this finding suggested that the face-selective channels indeed emerged in the d-AlexNet, though the face selectivity was weaker than the AlexNet.

To test the generalizability of these findings, we applied the same deprivation manipulation to another representative DCNN architecture, the ResNet-18, and the resultant d-ResNet reached top-1 accuracy (69.57%) and the top-5 accuracy (89.47%), comparable with those of the ResNet. Further, the face categorization accuracy of the d-ResNet (92.90%) was comparable to that of the ResNet (96.02%), and the discrimination accuracy of d-ResNet (65.34%) comparable to that of the pre-trained ResNet (59.80%). These findings were similar to those achieved with the d-AlexNet and the AlexNet.

How did the face-selective channels correspond to face-selective cortical regions in humans, such as the FFA and OFA? To answer this question, we calculated the coefficient of determination ($R^2$) of the multiple regression with the output of the face-selective channels as regressors and the fMRI signals from human visual cortex in response to movies on natural vision as the regressand (see section Materials and Methods). As shown in **Figure 2B** (right), the face-selective channels identified in the d-AlexNet corresponded to the bilateral FFA, OFA, and the posterior superior temporal sulcus face area (pSTS-FA). Similar correspondence was also found with the top two face-selective channels in the AlexNet (**Figure 2B**, left). Direct visual inspection revealed that the deprivation weakened the correspondence between the face-selective channels and face-selective regions in human brain. The increased channel-brain correspondence in the face-selective regions in the AlexNet compared with the d-AlexNet was confirmed by a two-way ANOVA of visual experience (d-AlexNet vs. AlexNet) by categorical selectivity (fMRI defined object-selective vs. face-selective regions, see section Methods). In addition to a main effect of categorical selectivity [$F_{(1, 398)} = 53.04$, $p < 0.001$, partial $\eta^2 = 0.12$], we also observed a two-way interaction [$F_{(1, 398)} = 79.99$, $p < 0.001$, partial $\eta^2 = 0.17$]. Follow-up simple effect analyses revealed that the correspondence to the face-selective regions decreased in the d-AlexNet as compared with the AlexNet in the face-selective regions (MD = −0.01, $p < 0.001$), but increased in the object-selective regions (MD = 0.013, $p < 0.001$), further indicating that the changes between the face-selective channels and human face-selective regions cannot be attributed to a global decrease in the channel-brain correspondence in the d-AlexNet comparing with the AlexNet.

We then examined whether this decrease in channel-brain correspondence affected different face-selective regions equally. A two-way ANOVA of visual experience (d-AlexNet vs. AlexNet) by regional correspondence (the OFA vs. the FFA) confirmed the decrease of channel-brain correspondence in the d-AlexNet compared with the AlexNet with a significant main effect of visual experiences [$F_{(1, 798)} = 161.97$, $p < 0.001$, partial $\eta^2 = 0.17$]. In addition, the main effect of the regional correspondence showed that the response profile of the face-selective channels in the DCNNs fitted better with the activation of the FFA than that of the OFA [$F_{(1, 798)} = 98.69$, $p = 0.001$, partial $\eta^2 = 0.11$], suggesting that the face-selective channels in DCNNs may in general tend to process faces as a whole than face parts. Critically, the two-way interaction was significant [$F_{(1, 798)} = 84.9$, $p < 0.001$, partial $\eta^2 = 0.10$], indicating that the experience affected the correspondence to the FFA and OFA disproportionally. A

**FIGURE 2 | (A)** The category-wise activation profiles of example face-selective channels of the AlexNet (left) and the d-AlexNet (right). The "faces" shown here were AI-generated for illustration purpose only. **(B)** The $R^2$ maps of the regression with the activation of the d-AlexNet's (right) or the AlexNet's face-selective channels (left) as the independent variables. The higher $R^2$ in multiple regression, the better correspondence between the face channels in the DCNNs and the face-selective regions in the human brain. The crimson lines delineate the ROIs of the OFA and the FFA. **(C)** The face-channels of both DCNNs corresponded better with the FFA than the OFA, and the difference between the AlexNet and the d-AlexNet was larger in the FFA. **(D)** Face inversion effect. The average activation amplitude of the top two face-selective channels differed in response to upright and inverted faces in the AlexNet but not the d-AlexNet. The error bar denotes standard error. The asterisk denotes statistical significance ($\alpha = 0.05$). n.s. denotes no significance.

simple effect analysis revealed that the correspondence to the FFA (MD = 0.023, $p < 0.001$) was increased by face-specific experiences to a significantly larger extent than that to the OFA

(MD = 0.004, $p = 0.013$, **Figure 2C**). Since the FFA is more involved in holistic processing of faces and the OFA is more dedicated to the part-based analysis, the disproportional decrease

in correspondence between the face-selective channels in the d-AlexNet and the FFA implied that the role of the experience on faces was to facilitate the processing of faces as a whole.

To test this conjecture, we examined whether the d-AlexNet responded stronger to upright than inverted faces, since human studies suggested that the upright faces were processed in a more holistic manner than inverted faces. As expected, there was a face inversion effect in the AlexNet's face-selective channels, with the magnitude of the activation to upright faces significantly larger than that to inverted faces [$t_{(19)} = 6.45$, $p < 0.001$, Cohen's $d = 1.44$] (**Figure 2D**). However, no inversion effect was observed in the d-AlexNet, as the magnitude of the activation to upright faces was not significantly larger than that to inverted faces [$t_{(19)} = 0.86$, $p = 0.40$]. The lack of the inversion effect in the d-AlexNet was further supported by a two-way interaction of visual experience by orientation of faces, [$F_{(1, 19)} = 7.79$, $p = 0.012$, partial $\eta^2 = 0.29$]. That is, unlike the AlexNet, the d-AlexNet processed upright faces in the same fashion as inverted faces.

Previous studies on human suggested that inverted faces are processed in an object-like fashion. That is, it relies more on the parts-based analysis than the holistic processing. Therefore, we speculated that in the d-AlexNet faces were also represented more like non-face objects. To test this speculation, we first compared the representational similarity among responses in Conv5 to faces and bowling-pins, which were not present as a category in the training dataset of either DCNNs, and therefore alien to both DCNNs. As expected, the two-way interaction of experience (AlexNet vs. d-AlexNet) by category (faces vs. bowling-pins) was significant [$F_{(1, 6,318)} = 4,110.88$, $p < 0.001$, partial $\eta^2 = 0.39$], and the simple effect analysis suggested that the representation for faces in the AlexNet was more similar between each other than in the d-AlexNet (MD = 0.16, $p < 0.001$), whereas the within-category representation similarity for bowling-pins showed the same but numerically smaller between-DCNN difference (MD = 0.005, $p = 0.002$) (**Figure 3A**).

A more critical test was to examine how face-specific experiences made faces being processed differently from objects. Here we calculated between-category similarities between faces and bowling-pins. We found that the between-category similarity between faces and bowling-pins was significantly higher in the d-AlexNet than that in the AlexNet [$t_{(3,159)} = 42.42$, MD = 0.07, $p < 0.001$, Cohen's $d = 0.76$] (**Figure 3B**), suggesting that faces in the d-AlexNet were indeed represented more like objects. In short, although d-AlexNet was able to perform face tasks similar to the one with face-specific experiences, it represented faces in an object-like fashion.

Finally, we asked how faceness was achieved in DCNNs with face-specific experiences. Neurophysiological studies on monkeys demonstrate experience-associated sharpening of neural response, with fewer neurons activated after learning. Here we performed a similar analysis by measuring the number of non-zero units (i.e., units with above-zero activation) of the face-selective channels activated by natural images containing faces. As shown in the activation map (**Figure 3C**), a smaller number of units were activated by faces in the AlexNet than that in the d-AlexNet [$t_{(19)} = 3.317$, MD = 15.78, Cohen's $d =$
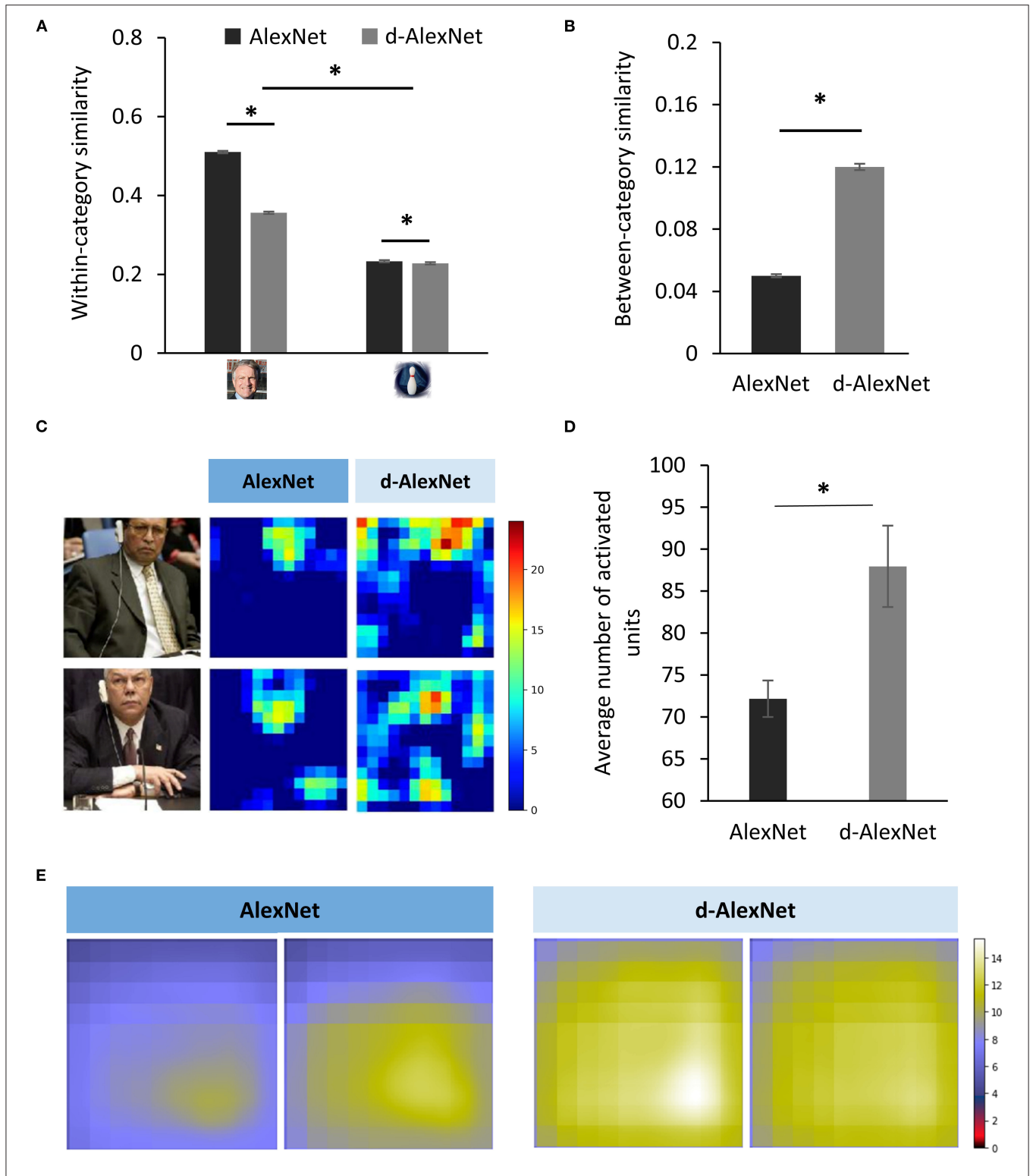
0.74] (**Figure 3D**), suggesting that the experience on faces made the representation to faces sparser, and thus allowing for more efficient coding. Another effect of visual experiences observed in neurophysiological studies is that experiences reduce the size of neurons' receptive field. Here we also mapped the empirical receptive field of the face-selective channels (see section Materials and Methods). Similarly, we found that the empirical receptive field of the AlexNet was smaller than that of the d-AlexNet. That is, within the theoretical receptive field, the empirical receptive field of the face-selective channels in the AlexNet was tuned to focus on a smaller region by face-specific experiences (**Figure 3E**).

## DISCUSSION

This study presented a DCNN model of selective visual deprivation of faces. Specifically, we chose the AlexNet as a test platform because of the functional correspondence along the hierarchy between the AlexNet and primates' ventral visual pathway (e.g., Krizhevsky et al., 2012; Cadieu et al., 2014; Wen et al., 2017; Pospisil et al., 2018; Baek et al., 2019). We found that without genetic predisposition and face-specific visual experiences, DCNNs were still capable of face perception. In addition, face-selective channels were also present in the d-AlexNet, which corresponded to human face-selective regions. That is, the visual experience of faces was not necessary for an intelligent system to develop a face-selective module. On the other hand, besides the slightly compromised selectivity of the module, the deprivation led the d-AlexNet to process faces in a fashion more similar to that of processing objects. Indeed, unlike the AlexNet, face inversion did not affect the response magnitude of the face-selective channels in the d-AlexNet, and the representation of faces was more similar to objects as compared to the AlexNet. Finally, face-specific experiences might affect face processing by fine-tuning the sparse coding and the size of the receptive field of the face-selective channels. In sum, our study addressed a long-standing debate on nature vs. nurture in developing the face-specific module, and illuminated the role of visual experiences in shaping the module.

Given the main-stream viewpoint that faces are special and therefore cannot be compensated by the presence of non-face objects, it may seem surprising that without domain-specific visual experience, the face-selective processing and modules still emerged in the d-AlexNet. These observations were further replicated with another well-known DCNN architecture, the ResNet-18, suggesting the generalizability of our findings. However, our finding is consistent with previous studies on non-human primates and new-born human infants (Bushneil et al., 1989; Valenza et al., 1996; Sugita, 2008), where the face-specific experience is found not necessary for face detection and recognition. Therefore, our study argues against the experience-independent hypothesis that face specificity is largely attributed to either innate face-specific mechanisms (Morton and Johnson, 1991; McKone et al., 2012) or domain-general processing with predisposed biases (Simion et al., 2001; Simion and Di Giorgio, 2015). Our study argues against this conjecture, because unlike

FIGURE 3 | (A) The within-category similarity in the face category and an unseen non-face category (bowling pins) in the DCNNs. (B) The between-category similarity between faces and bowling pins. (C) The activation maps of a typical face-selective channel of each DCNN in responses to natural images containing faces. Each pixel denotes activation in one unit. The color denotes the activation amplitude (a.u.). (D) The extent of activation of the face-selective channels of each DCNN in responses to natural images containing faces. (E) The empirical receptive fields of the top two face-selective channels of each DCNN. The color denotes the average activation amplitude (a.u, see section Sparse Coding and Empirical Receptive Field). The error bar denotes standard error. The asterisk denotes statistical significance ($\alpha = 0.05$). The real faces used in this figure are adapted from the FITW dataset.

any biological system, DCNNs have no domain-specific genetic inheritance or processing biases. Therefore, the face-specific processing observed in DCNNs had to derive from domain-general factors. From this sense, the present study provides one of the first direct evidence against the main-stream viewpoint and suggests that face specificity may emerge from domain-general visual experience.

We speculated that the face-selective processing and module in the d-AlexNet may result from the rich features represented in the multiple layers of the network; face-like features might be utilized when the neural network was forced to categorize faces even though these features were not learned for this purpose. In fact, previous studies on DCNNs have shown that DCNN's lower layers showed sensitivity to myriad visual features similar to primates' primary visual cortex (Krizhevsky et al., 2012), while the higher layers are tuned to complex features resembling those represented in the ventral visual pathway (Yamins et al., 2014; Güçlü and van Gerven, 2015). With such a repertoire of rich features, a representational space for faces, or any natural object, may be constructed by selecting features that are potentially useful in face tasks. With such repertoire of rich features, a representational space for faces, or for any natural object, may be constructed by selecting features that are potentially useful in face tasks.

Supporting evidence for this conjecture came from the observation that the d-AlexNet processed faces in an object-like fashion. For example, the face inversion effect, a signature of face-specific processing in human (Yin, 1969; Kanwisher et al., 1998) was absent in the d-AlexNet. Distinct from other non-face stimuli, faces are recognized better when they are upright than inverted (Yin, 1969), and the neural response to upright faces is stronger than that to inverted ones (e.g., Kanwisher et al., 1998; Rossion and Gauthier, 2002). This face inversion effect is attributed to that face processing relies particularly heavily on configural processing—processing of the relations among features instead of individual features. Since the configural information is difficult to perceive in inverted faces in a system with face specificity, inverted faces cannot engage face-specific processing as upright faces. Therefore, the finding of the lack of the face inversion effect in the DNN without face experience strengthened our argument that the lack of face experience leads to the compromise of face specificity. That is, similar to inverted faces, upright faces may also be processed like objects in the d-AlexNet. A more direct illustration of the object-like representation of faces came from the analysis of the representational similarity between faces and objects. As compared to the AlexNet, faces in the representational space of the d-AlexNet were less congregated among each other; instead they were more intermingled with non-face object categories. The finding that face representation was no longer qualitatively different from object representation may help to explain the performance of the d-AlexNet. Because faces were less segregated from objects in the representational space, the d-AlexNet's accuracy of face categorization was worse than that of the AlexNet. In contrast, within the face category, individual faces were less congregated in the representational space; therefore, the discrimination of individual faces became easier

instead, suggested by the slightly higher face discrimination accuracy in the d-AlexNet than the AlexNet. In short, when the representational space of the d-AlexNet was formed exclusively based on features from non-face stimuli, faces were represented no longer qualitatively different from non-face objects, which inevitably led to "object-like" face processing.

The face-specific processing is likely achieved through prior exposure to faces. At first glance, the effect of face-specific experiences seemed quantitative, as in the AlexNet, both the selectivity to faces and the number of the face-selective channels were increased, and the correspondence between the face-selective channels and the face-selective regions in human brain was tighter. However, careful scrutiny of the difference between the two DCNNs revealed that the changes led by the experience may be qualitative. For example, the deprivation of visual experiences disproportionally weakened the DCNN-brain correspondence in the FFA as comparing to the OFA, and the FFA is engaged more in the configural processing and the OFA in parts-based analysis (Liu et al., 2010; Nichols et al., 2010; Zhao et al., 2014). Therefore, the "face-like" face processing may come from the fact that face-specific experiences led the representation of faces more congregated within face category and more separable from the representation of non-face objects stimuli (see also Gomez et al., 2019). In this way, a relative encapsulated representation may help developing a unique way of processing faces, qualitatively different from non-face objects.

The computational transparency of DCNNs may shed light on the development of domain specificity of the face module. First, we found that face-specific experiences increased the sparseness of face representation, as fewer units of the face channels were activated by faces in the AlexNet. The experience-dependent sparse coding has been widely discovered in the visual cortex (for reviews, see Desimone, 1996; Grill-Spector et al., 2006). The experience-induced increase of sparseness is thought to reflect a preference-narrowing process that tunes neurons to a smaller range of stimuli (Kohn and Movshon, 2004); therefore, with sparse coding faces are less likely to be intermingled with non-face objects, which may lead to more congregated representations in the representational space in the AlexNet, as compared to the d-AlexNet. Second, we found that the empirical receptive field of the face channels in the AlexNet was smaller than that in the d-AlexNet, suggesting that the visual experience on faces decreased the size of the receptive field of the face channels. This finding fits perfectly with neurophysiological studies that the size of receptive fields of visual neurons is reduced after eye-opening (Braastad and Heggelund, 1985; Tavazoie and Reid, 2000; Cantrell et al., 2010). Importantly, along with the refined receptive fields, the selectivity of neurons increases (Spilmann, 2014), possibly because neurons can avoid distracting information by focusing on a more restricted part of stimuli, which may further allowed finer representation of the selected regions. This is especially important for processing faces because faces are highly homogeneous, and some information is identical across faces, such as parts composition (eyes, noses, and mouth) and their configural arrangements. Therefore, the reduced receptive field of the face channels may facilitate selective analyses of discriminative face features while avoiding irrelevant

information. Further, the sharpening of the receptive field and the fine-tuned selectivity may result in superior discrimination ability on faces, and allow faces to be processed at the subordinate level (i.e., identification), whereas the rest of objects are largely processed at the basic level (i.e., categorization).

It has long been assumed that domain-specific visual experiences and inheritance are the pre-requisites in the development of the face module in the brain. In our study with DCNN as a model, we completely decoupled the genetic predisposition and face-specific visual experiences, and found that the representation for faces can be constructed with features from non-face objects to realize basic functions for face recognition. Therefore, in many situations, the difference between faces and objects is "quantitative" rather than "qualitative," as they are represented in a continuum of the representational space. In addition, we also found that face-specific experiences likely fine-tuned the face representation, and thus transformed the "object-like" face processing into "face-specific" processing. However, we shall be cautious that our finding may not be applicable for the development of face module in human, as in the biological brain experience-induced changes are partly attributed to the inhibition from lateral connections (Norman and O'Reilly, 2003; Grill-Spector et al., 2006), whereas there is no lateral or feedback connection in DCNNs. However, despite structural differences, recent studies have shown similar representation for faces between DCNNs and humans (Song et al., 2021), suggesting that a common mechanism may be shared by both artificial and biological intelligent systems. Future studies are needed to examine the applicability of our finding to humans. In addition, higher cognitive functions such as attractiveness judgement and social-traits inference are also important components of face processing, but the present study followed the literature on face deprivation in humans and non-human primates and therefore focused on the sensory and perceptual stages of face processing. Future study may consider investigating the experiential effects on the social and affective aspects of face processing to comprehensively understand the effect of experience.

On the other hand, our study illustrated the advantages of using DCNN as a model to understand human mind because of its computational transparency and its dissociation of factors in nature and nurture. Thus, our study invites future studies with DCNNs to understand the development of domain specificity in particular and a broad range of cognitive modules in general.

## DATA AVAILABILITY STATEMENT

Publicly available datasets were analyzed in this study. This data can be found at: https://github.com/sdgds/ Deprivation_dataset; http://www.vision.caltech.edu/Image_ Datasets/Caltech256/; http://vis-www.cs.umass.edu/lfw/; https://openneuro.org/datasets/ ds001761.

## ETHICS STATEMENT

This study used human fMRI data, the acquisition of which was reviewed and approved by Institutional Review Board at Purdue University Institutional Review Board.

## AUTHOR CONTRIBUTIONS

JL conceived and designed the study. YZ analyzed the data with input from all authors. SX wrote the manuscript with input from JL, YZ, and ZZ.

## FUNDING

## REFERENCES

AbdAlmageed, W., Wu, Y., Rawls, S., Harel, S., Hassner, T., Masi, I., et al. (2016). "Face recognition using deep multi-pose representations." in *2016 IEEE Winter Conference on Applications of Computer Vision* (Lake Placid, NY).

Arcaro, M. J., Schade, P. F., Vincent, J. L., Ponce, C. R., and Livingstone, M. S. (2017). Seeing faces is necessary for face-domain formation. *Nat. Neurosci.* 20:1404. doi: 10.1038/nn.4635

Baek, S., Song, M., Jang, J., Kim, G., and Paik, S.-B. (2019). *Spontaneous generation of face recognition in untrained deep neural networks. bioRxiv, 857466.* Available online at: https://www.biorxiv.org/content/10.1101/857466v1 (accessed November 29, 2019).

Berg, T. L., Berg, A. C., Edwards, J., and Forsyth, D. A. (2005). Who's in the picture. *Adv. Neural Inf Process. Syst.* 17, 137–144. Retrieved from: https://papers.nips. cc/paper/2004/file/03fa2f7502f5f6b9169e67d17cbf51bb-Paper.pdf

Braastad, B. O., and Heggelund, P. (1985). Development of spatial receptive-field organization and orientation selectivity in kitten striate cortex. *J. Neurophysiol.* 53, 1158–1178. doi: 10.1152/jn.1985.53.5.1158

Bushneil, I., Sai, F., and Mullin, J. (1989). Neonatal recognition of the mother's face. *Br. J. Dev. Psychol.* 7, 3–15. doi: 10.1111/j.2044-835X.1989.tb0 0784.x

Cadieu, C. F., Hong, H., Yamins, D. L. K., Pinto, N., Ardila, D., Solomon, E. A., et al. (2014). Deep neural networks rival the representation of primate it cortex for core visual object recognition. *PLoS Comput. Biol.* 10:e1003963. doi: 10.1371/journal.pcbi.1003963

Cantrell, D. R., Cang, J., Troy, J. B., and Liu, X. (2010). Non-centered spike-triggered covariance analysis reveals neurotrophin-3 as a developmental regulator of receptive field properties of ON-OFF retinal ganglion cells. *PLoS Comput. Biol.* 6:e1000967. doi: 10.1371/journal.pcbi.10 00967

Chen, X., Zhou, M., Gong, Z., Xu, W., Liu, X., Huang, T., et al. (2020). DNNBrain: a unifying toolbox for mapping deep neural networks and brains. *Front. Comput. Neurosci.* 14:580632. doi: 10.3389/fncom.2020.580632

Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L. (2009). "Imagenet: a large-scale hierarchical image database," in *Paper Presented at the 2009 IEEE Conference on Computer Vision and Pattern Recognition* (Miami, FL).

Desimone, R. (1996). Neural mechanisms for visual memory and their role in attention. *Proc. Natl. Acad. Sci. U.S.A.* 93, 13494–13499. doi: 10.1073/pnas.93.24.13494

Freiwald, W., Duchaine, B., and Yovel, G. (2016). Face processing systems: from neurons to real-world social perception. *Annu. Rev. Neurosci.* 39, 325–346. doi: 10.1146/annurev-neuro-070815-013934

Gomez, J., Barnett, M., and Grill-Spector, K. (2019). Extensive childhood experience with Pokemon suggests eccentricity drives organization of visual cortex. *Nat. Hum. Behav.* 3, 611–624. doi: 10.1038/s41562-019-0592-8

Griffin, G., Holub, A., and Perona, P. (2007). *Caltech-256 Object Category Dataset.* Pasadena, CA: California University of Technology.

Grill-Spector, K., Henson, R., and Martin, A. (2006). Repetition and the brain: neural models of stimulus-specific effects. *Trends Cogn. Sci.* 10, 14–23. doi: 10.1016/j.tics.2005.11.006

Grundstrom, J., Chen, J., Ljungqvist, M. G., and Astrom, K. (2016). "Transferring and compressing convolutional neural networks for face representations," in *Image Analysis and Recognition, Vol. 9730*, eds A. Campilho and F. Karray (Cham: Springer), 20–29.

Güçlü, U., and van Gerven, M. A. (2015). Deep neural networks reveal a gradient in the complexity of neural representations across the ventral stream. *J. Neurosci.* 35, 10005–10014. doi: 10.1523/JNEUROSCI.5023-14.2015

He, K., Zhang, X., Ren, S., Sun, J., and Ieee. (2016). "Deep residual learning for image recognition," in *2016 IEEE Conference on Computer Vision and Pattern Recognition* (Las Vegas, NV), 770–778.

Kanwisher, N., Tong, F., and Nakayama, K. (1998). The effect of face inversion on the human fusiform face area. *Cognition* 68, B1–B11. doi: 10.1016/S0010-0277(98)00035-3

Kanwisher, N., and Yovel, G. (2006). The fusiform face area: a cortical region specialized for the perception of faces. *Philos. Trans. R. Soc. B Biol. Sci.* 361, 2109–2128. doi: 10.1098/rstb.2006.1934

Kohn, A., and Movshon, J. A. (2004). Adaptation changes the direction tuning of macaque MT neurons. *Nat. Neurosci.* 7, 764–772. doi: 10.1038/nn1267

Kriegeskorte, N. (2015). Deep neural networks: a new framework for modeling biological vision and brain information processing. *Annu. Rev. Vis. Sci.* 1, 417–446. doi: 10.1146/annurev-vision-082114-035447

Krizhevsky, A. (2014). One *weird trick for parallelizing convolutional neural networks. arXiv preprint arXiv:1404.5997.* Available online at: https://arxiv.org/abs/1404.5997 (accessed April 29, 2014).

Krizhevsky, A., Sutskever, I., and Hinton, G. E. (2012). Imagenet classification with deep convolutional neural networks. *Adv. Neural Inf Process. Syst.* 25, 1097–1105. doi: 10.1145/3065386

Liu, J., Harris, A., and Kanwisher, N. (2010). Perception of face parts and face configurations: an fMRI study. *J. Cogn. Neurosci.* 22, 203–211. doi: 10.1162/jocn.2009.21203

McKone, E., Crookes, K., Jeffery, L., and Dilks, D. D. (2012). A critical review of the development of face recognition: experience is less important than previously believed. *Cogn. Neuropsychol.* 29, 174–212. doi: 10.1080/02643294.2012.660138

Morton, J., and Johnson, M. H. (1991). CONSPEC and CONLERN: a two-process theory of infant face recognition. *Psychol. Rev.* 98:164. doi: 10.1037/0033-295X.98.2.164

Nichols, D. F., Betts, L. R., and Wilson, H. R. (2010). Decoding of faces and face components in face-sensitive human visual cortex. *Front. Psychol.* 1:28. doi: 10.3389/fpsyg.2010.00028

Norman, K. A., and O'Reilly, R. C. (2003). Modeling hippocampal and neocortical contributions to recognition memory: a complementary-learning-systems approach. *Psychol. Rev.* 110, 611–646. doi: 10.1037/0033-295X.110.4.611

Parkhi, O. M., Vedaldi, A., and Zisserman, A. (2015). Deep face recognition. *Proc. Br. Mach. Vis.* 1, 1–12. doi: 10.5244/C.29.41

Pospisil, D. A., Pasupathy, A., and Bair, W. (2018). "Artiphysiology" reveals V4-like shape tuning in a deep network trained for image classification. *Elife* 7:e38242. doi: 10.7554/eLife.38242

Rossion, B., and Gauthier, I. (2002). How does the brain process upright and inverted faces? *Behav. Cogn. Neurosci. Rev.* 1, 63–75. doi: 10.1177/1534582302001001004

Simion, F., and Di Giorgio, E. (2015). Face perception and processing in early infancy: inborn predispositions and developmental changes. *Front. Psychol.* 6:969. doi: 10.3389/fpsyg.2015.00969

Simion, F., Macchi Cassia, V., Turati, C., and Valenza, E. (2001). The origins of face perception: specific versus non-specific mechanisms. *Infant Child Dev. Int. J. Res. Pract.* 10, 59–65. doi: 10.1002/icd.247

Song, Y., Qu, Y., Xu, S., and Liu, J. (2021). Implementation-independent representation for deep convolutional neural networks and humans in processing faces. *Front. Comput. Neurosci.* 14:601314. doi: 10.3389/fncom.2020.601314

Spilmann, L. (2014). Receptive fields of visual neurons: the early years. *Perception* 43, 1145–1176. doi: 10.1068/p7721

Sugita, Y. (2008). Face perception in monkeys reared with no exposure to faces. *Proc. Natl. Acad. Sci. U.S.A.* 105, 394–398. doi: 10.1073/pnas.0706079105

Tavazoie, S. F., and Reid, R. C. (2000). Diverse receptive fields in the lateral geniculate nucleus during thalamocortical development. *Nat. Neurosci.* 3, 608–616. doi: 10.1038/75786

Valenza, E., Simion, F., Cassia, V. M., and Umilt,à, C. (1996). Face preference at birth. *J. Exp. Psychol. Hum. Percept. Perform.* 22:892. doi: 10.1037/0096-1523.22.4.892

Van Essen, D. C., Smith, S. M., Barch, D. M., Behrens, T. E. J., Yacoub, E., Ugurbil, K., et al. (2013). The WU-Minn human connectome project: an overview. *Neuroimage* 80, 62–79. doi: 10.1016/j.neuroimage.2013.05.041

VanRullen, R., and Reddy, L. (2019). Reconstructing faces from fMRI patterns using deep generative neuralnetworks. *Commun. Biol.* 2, 193–193. doi: 10.1038/s42003-019-0438-y

Wen, H., Shi, J., Zhang, Y., Lu, K.-H., Cao, J., and Liu, Z. (2017). Neural encoding and decoding with deep learning for dynamic natural vision. *Cereb. Cortex* 28, 4136–4160. doi: 10.1093/cercor/bhx268

Wilmer, J. B., Germine, L., Chabris, C. F., Chatterjee, G., Williams, M., Loken, E., et al. (2010). Human face recognition ability is specific and highly heritable. *Proc. Natl. Acad. Sci. U.S.A.* 107, 5238–5241. doi: 10.1073/pnas.0913053107

Yamins, D. L., Hong, H., Cadieu, C. F., Solomon, E. A., Seibert, D., and DiCarlo, J. J. (2014). Performance-optimized hierarchical models predict neural responses in higher visual cortex. *Proc. Natl. Acad. Sci.* 111, 8619–8624. doi: 10.1073/pnas.1403112111

Yi, D., Lei, Z., Liao, S., and Li, S. Z. (2014). *Learning face representation from scratch. arXiv preprint arXiv:1411.7923.* Available online at: https://arxiv.org/abs/1411.7923 (accessed November 28, 2014).

Yin, R. K. (1969). Looking at upside-down faces. *J. Exp. Psychol.* 81, 141–145. doi: 10.1037/h0027474

Zhao, M., Cheung, S.-H., Wong, A. C. N., Rhodes, G., Chan, E. K. S., Chan, W. W. L., et al. (2014). Processing of configural and componential information in face-selective cortical areas. *Cogn. Neurosci.* 5, 160–167. doi: 10.1080/17588928.2014.912207

Zhen, Z., Yang, Z., Huang, L., Kong, X.-,z., Wang, X., Dang, X., et al. (2015). Quantifying interindividual variability and asymmetry of face-selective regions: a probabilistic functional atlas. *Neuroimage* 113, 13–25. doi: 10.1016/j.neuroimage.2015.03.010

Zhu, Q., Song, Y., Hu, S., Li, X., Tian, M., Zhen, Z., et al. (2010). Heritability of the specific cognitive ability of face perception. *Curr. Biol.* 20, 137–142. doi: 10.1016/j.cub.2009.11.067