

20230127

1       **Analysis of convolutional neural networks reveals the computational properties**  
2                               **essential for subcortical processing of facial expression**

3  
4  
5  
6                   **Chanseok Lim<sup>1,2</sup>, Mikio Inagaki<sup>1,3</sup>, Takashi Shinozaki<sup>3,4</sup>, Ichiro Fujita<sup>1,3,5,\*</sup>**

7  
8       <sup>1</sup>Laboratory for Cognitive Neuroscience, and <sup>2</sup>Perceptual and Cognitive Neuroscience  
9       Laboratory, Graduate School of Frontier Biosciences, Osaka University, 1-4 Yamadaoka, Suita,  
10       Osaka 565-0871, Japan; <sup>3</sup>Center for Information and Neural Networks, National Institute of  
11       Information and Communications Technology, 1-4 Yamadaoka, Suita, Osaka 565-0871, Japan;  
12       <sup>4</sup>Computational Neuroscience Laboratory, Faculty of Informatics, Kindai University, 3-4-1  
13       Kowakae, Higashiosaka, Osaka 577-8502, Japan; <sup>5</sup>Research Organization of Science and  
14       Technology, Ritsumeikan University, 1-1-1 Noji-higashi, Kusatsu, Shiga 525-8577, Japan  
15  
16

17       **ORCID numbers:**      Chanseok Lim                   0000-0002-4611-2740  
18                               Mikio Inagaki                   0000-0002-5294-501X  
19                               Takashi Shinozaki           0000-0002-0641-6782  
20                               Ichiro Fujita                   0000-0003-3293-8610  
21  
22

23       **\*Corresponding author:**    Ichiro Fujita, PhD  
24    Laboratory for Cognitive Neuroscience, Graduate School of  
25    Frontier Biosciences, Osaka University, 1-4 Yamadaoka, Suita,  
26    Osaka 565-0871, Japan  
27    Phone:                   +81-6-6879-4439  
28    Fax:                       +81-6-6879-4439  
29    E-mail:                  fujita@fbs.osaka-u.ac.jp  
30

31       **Number of pages:**           36 pages  
32       **Number of figures:**          8 figures  
33       **Number of tables:**           1 table  
34       **Number of words in abstract:** 200 words  
35  
36

20230127

37 **Abstract**

38

39 Perception of facial expression is crucial in the social life of primates. This visual information  
40 is processed along the ventral cortical pathway and the subcortical pathway. Processing of face  
41 information in the subcortical pathway is inaccurate, but the architectural and physiological  
42 properties that are responsible remain unclear. We analyzed the performance of convolutional  
43 neural networks incorporating three prominent properties of this pathway: a shallow layer  
44 architecture, concentric receptive fields at the first processing stage, and a greater degree of  
45 spatial pooling. The neural networks designed in this way could be trained to classify seven  
46 facial expressions with a correct rate of 51% (chance level, 14%). This modest performance was  
47 gradually improved by replacing the three properties, one-by-one, two at a time, or all three  
48 simultaneously, with the corresponding features in the cortical pathway. Some processing units  
49 in the final layer were sensitive to spatial frequencies (SFs) in the retina-based coordinate,  
50 whereas others were sensitive to object-based SFs, similar to neurons in the amygdala.  
51 Replacement of any one of these properties affected the SF coordinate of units. All three  
52 properties constrain the accuracy of facial expression information in the subcortical pathway,  
53 and are essential for determining the coordinate of SF representation.

54

55

56 **Keywords**

57 Subcortical visual pathway; ventral pathway, convolutional neural network; amygdala; face  
58 recognition

59

20230127

## 60 **Introduction**

61

62 Perceiving the facial expressions of other individuals plays a critical role in the social life of  
63 primates, including humans. Two neural pathways, the ventral cortical pathway and the  
64 subcortical pathway, contribute to this perceptual ability (Fig. 1A; Pessoa and Adolphs, 2010;  
65 Tamietto and de Gelder, 2010; Petray and Bickford, 2019). The ventral cortical pathway consists  
66 of a network of areas in the occipito-temporal region of the cerebral cortex, and processes a  
67 variety of visual features of objects, people, and environments, including shape, color, texture,  
68 material properties, and binocular disparity (Ungerleider and Mishkin, 1982; Connor et al.,  
69 2007; Conway et al., 2010; Roe et al., 2012; Kravitz et al., 2013; Vaziri et al., 2014; Verhoef et  
70 al., 2016; Komatsu and Goda, 2018). Neurons that preferentially respond to images of faces or  
71 facial features are found in several clusters along this pathway (Desimone et al., 1984; Perrett  
72 et al., 1987; Fujita et al., 1992; Haxby et al., 2000; Tsao and Livingstone, 2008; Duchaine and  
73 Yovel, 2015; Freiwald et al., 2016). They constitute the neural system that analyzes facial details  
74 such as expression, identity, and direction of attention. The subcortical pathway consists of a  
75 few processing stages in phylogenetically ancient regions: the superior colliculus in the  
76 midbrain, the pulvinar nucleus in the posterior thalamus, and the amygdala in the medial limbic  
77 system. The subcortical pathway is suggested to mediate rapid behavioral and physiological  
78 (autonomic) responses to sensory signals related to possible dangers such as fearful faces  
79 (Tamietto and de Gelder, 2010; Nakano et al., 2013; for a critical review, see Pessoa and  
80 Adolphs, 2010). The ventral cortical pathway and the subcortical pathway intersect at the  
81 amygdala.

82

83 Psychological and brain imaging studies suggest that the subcortical pathway subserves the  
84 ability of some patients with lesions in the primary visual cortex (V1) to discriminate facial  
85 expressions despite lacking visual awareness (“affective blindsight”; deGelder et al., 1999;  
86 Pegna et al., 2005; Striemer et al., 2019). These patients also reflexively exhibit specific facial  
87 expressions and pupillary reactions when exposed to fearful or happy faces (Tamietto et al.,  
88 2009). Studies have also shown that the subcortical pathway supports unconscious face  
89 perception in neurologically healthy subjects (Morris et al., 1999, 2001). Furthermore,  
90 orientation bias toward faces or face-like patterns by newborn babies is suggested to be mediated  
91 by the subcortical pathway (Cassia et al., 2001; Johnson, 2005; but see Buiatti et al., 2019).  
92 Importantly, these perceptual abilities are not perfectly accurate, instead resulting in modest  
93 performance at above-chance levels. These findings suggest that information on faces conveyed  
94 by the subcortical pathway is less accurate than that carried by the ventral cortical pathway.

95

96 Electrophysiological studies have demonstrated that processing of facial expression in the  
97 subcortical pathway is indeed fast and not very accurate. Méndez-Bértolo and colleagues (2016)  
98 showed that intracranial local field potentials in the human amygdala respond differentially to  
99 fearful faces versus other faces within 74 ms after stimulus onset. A recent single-neuron  
100 recording study in the monkey revealed that a population of amygdala neurons responded to  
101 threatening faces even within 50 ms (Inagaki et al., 2022a). This early response, when combined  
102 across an ensemble of neurons, carries information that allows linear classifiers to discriminate  
103 threatening faces from neutral and affiliative faces. The rate of correctly discriminating the three  
104 expressions is around 50%; this is well above chance (33%), but significantly worse than perfect.

20230127

105 Some neurons in the superior colliculus and pulvinar of the monkey also respond to faces and  
106 face-like patterns with an even shorter latency of 30–50 ms (Nguyen et al., 2013, 2014).

107

108 What architectural and physiological properties of the subcortical pathway are responsible for  
109 its fast, crude processing? The fast processing most likely arises from the small number, or  
110 “shallowness,” of processing stages in the subcortical pathway, given that the ventral cortical  
111 pathway and its upstream area (the lateral geniculate nucleus) consist of a larger number of  
112 regions (at least six before reaching the amygdala) than the subcortical pathway (only two), and  
113 that every transition from one cortical region to the next takes at least 10 ms (Schmolesky et al.,  
114 1998). It is unclear whether the shallow processing similarly explains the low accuracy of the  
115 information transmitted by the subcortical pathway to the amygdala. This uncertainty arises  
116 from the fact that in addition to the difference in the number of processing stages, visual  
117 response properties differ markedly between the two pathways. Neurons in the superior  
118 colliculus at the first stage of the subcortical pathway show circular receptive fields with center-  
119 surround organization, which can be modeled using the difference-of-Gaussian (DoG) function  
120 (Cynader and Berman 1972; Updyke 1974; Marino et al. 2008; Churan et al., 2012). By contrast,  
121 simple cells of V1 at the first stage of the cortical pathway have elongated receptive fields with  
122 side-by-side ON and OFF subregions, which can be modeled by two-dimensional Gabor  
123 functions (Jones and Palmer, 1987). Furthermore, the receptive field is typically larger in the  
124 superior colliculus (for the foveal field, 1.5–10° in superficial layers, 10–20° in deep layers;  
125 Goldberg and Wurtz, 1972; Wallace et al., 1997) than in V1 (1.18° in simple cells, 1.3° in  
126 complex cells; Van den Bergh et al., 2010) and the extrastriate areas V2 and V4 (Freeman and  
127 Simoncelli, 2011). Thus, spatial pooling across ascending stages occurs over a wider visual field  
128 area in the subcortical pathway than in the ventral cortical pathway.

129

130 In the present study, we addressed how these properties of the subcortical pathway, i.e., the  
131 shallowness of processing stages, DoG-type receptive fields at the initial stage, and spatial  
132 pooling over a wider visual field, influence facial expression processing. To this aim, we  
133 constructed convolutional neural networks (CNNs) and analyzed their performance in facial  
134 expression discrimination. CNNs are one type of multilayer perceptron, and can be optimized  
135 (“learn”) to classify inputs by varying connection weights between processing units through  
136 supervised learning algorithms (LeCun et al., 2015). Typical CNNs have several to tens of layers  
137 (deep neural networks, DNNs). DNNs developed for classifying visual objects share  
138 architectural and representational features similar to those of the ventral cortical pathway  
139 (Yamins et al., 2014; Güçlü and van Gerven, 2015; Yamins and DiCarlo, 2016; Hassabis et al.,  
140 2017). We designed our CNNs to imitate the subcortical pathway by reducing the number of  
141 processing stages and by implementing DoG-type receptive fields and a wider extent of pooling.  
142 These CNNs, hereafter referred to as shallow neural networks (SNNs), learned to discriminate  
143 facial expressions with modest correct rates. Replacing the three properties, one-by-one, two at  
144 a time, or all three simultaneously, with the corresponding properties in the ventral cortical  
145 pathway gradually improved discrimination performance, suggesting that all three features are  
146 responsible for limiting the performance of the SNNs. We further showed that like some neurons  
147 in the amygdala, a major group of units in the final processing layer of the SNNs were sensitive  
148 to spatial frequency (SF) in the retina-based reference frame as initially detected in the first

20230127

149 processing layer, and that the three subcortical properties contribute to preserving the retina-  
150 based SF sensitivity.

151

152

## 153 **Materials and Methods**

154

### 155 *Architecture of SNNs.*

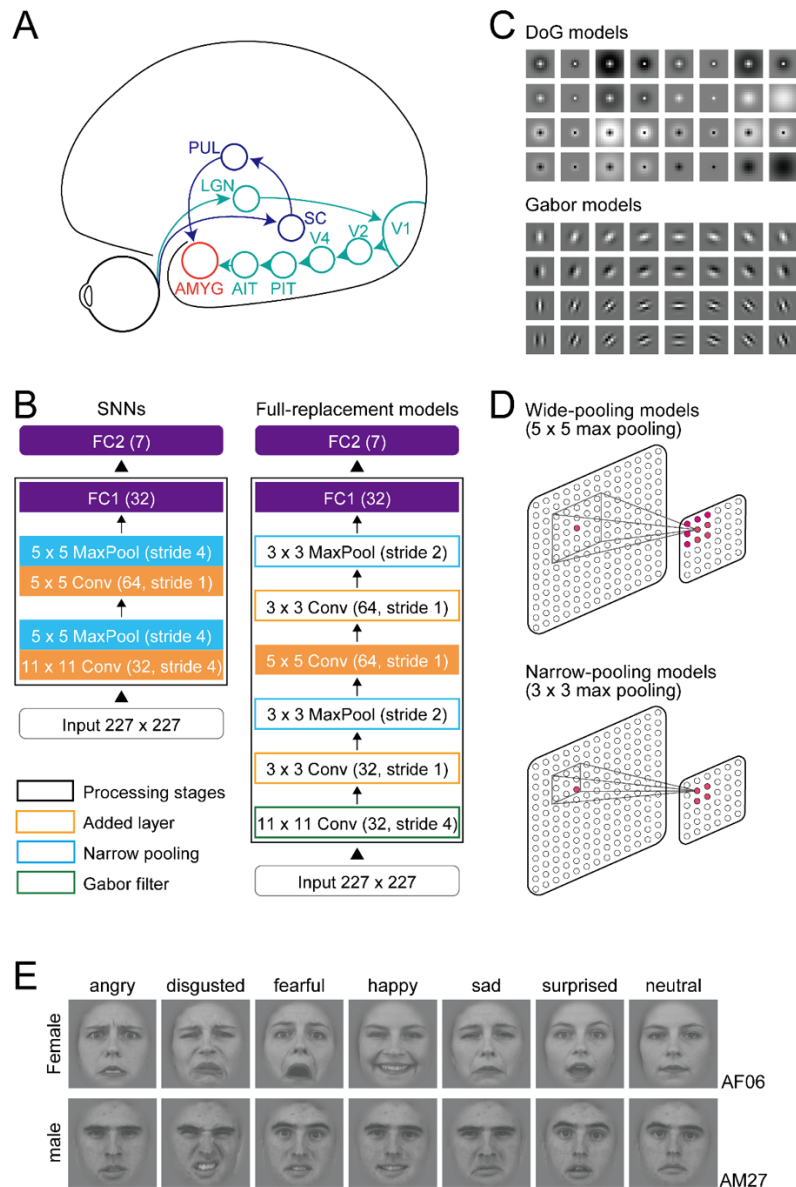
156 We constructed SNNs incorporating the distinct properties of the primate subcortical pathway  
157 (Fig. 1B–D; Table 1). Unlike typical DNNs, the SNNs consisted of only two sets of convolution  
158 and pooling layers followed by two fully connected layers (FC1, FC2), approximating the small  
159 number of processing stages of the subcortical pathway. The first convolution layer incorporated  
160 32 DoG-type filters (Fig. 1C, top) with a spatial resolution of  $11 \times 11$  pixels, whereas weights  
161 in the second convolution layer were initially random, i.e., the filters had no structure, and  
162 gradually changed through training. A rectified linear unit (ReLU) was used as the activation  
163 function of a unit in the convolution layers and FC1; the ReLU forwards the processing results  
164 directly to the next stage if they are positive, otherwise it outputs zero. A max pooling operation  
165 was performed over  $5 \times 5$  sliding regions with a stride of 4 for the outputs of convolution layers  
166 (Fig. 1D, top). Max pooling selected the largest value among the responses of units within a  
167 sliding window over the preceding convolution layer, and forwarded the value to the next layer.  
168 A local response normalization process was added after the pooling layers to aid generalization  
169 (Krizhevsky et al., 2012; we used slightly different parameters from theirs;  $k = 1$ ,  $n = 5$ ,  $\alpha = 2 \times$   
170  $10^{-5}$ ,  $\beta = 0.75$ ). Every unit in FC1 and FC2 received inputs from all units in the immediately  
171 preceding layer, i.e., each was fully connected. FC1 is the final processing layer, and FC2  
172 outputs the results of entire processing by the SNNs. These features were implemented to  
173 capture the architectural and computational properties of the subcortical pathway, i.e., fewer  
174 processing stages compared to the ventral cortical pathway (Fig. 1A), DoG-type receptive fields  
175 in the superficial layer of the superior colliculus (Churan et al., 2012), and large receptive fields  
176 of deeper superior colliculus neurons (Wallace et al., 1997). The first three processing layers  
177 were intended to represent the superior colliculus, pulvinar, and amygdala, respectively. The  
178 processing types of these layers, i.e., convolution and pooling in the first two layers and full  
179 connection in FC1, were chosen to match the retinotopic organization of the three brain regions.  
180 The convolution and pooling processes in the first two layers exploit retinotopy, as the superior  
181 colliculus and pulvinar contain retinotopic maps (Bender, 1981; Chen et al., 2019). The FC1  
182 layer loses retinotopic information because of the fully convergent connection from the earlier  
183 stage, as the amygdala does not have a retinotopic map (Morawetz et al., 2010).

184

185 The SNNs were trained to discriminate images of facial expressions representing seven basic  
186 emotions: angry, disgusted, fearful, happy, sad, surprised, and neutral (Fig. 1E; see below for  
187 details). For each input image, the seven units in FC2 yielded scores ranging from 0 to 1 for the  
188 seven expression categories, representing the probabilities of classified expressions. The  
189 expression with the highest score was taken as the output of the model.

190

20230127



191

192

193 **Figure 1.** Shallow neural networks (SNNs) and modifications. (A) Cortical and subcortical  
 194 visual pathways for processing facial expressions in the primate brain. AMYG: amygdala; AIT,  
 195 PIT: anterior and posterior parts of inferior temporal cortex; LGN: lateral geniculate nucleus;  
 196 PU: pulvinar nucleus; SC: superior colliculus; V1, V2, V4: visual areas 1, 2, 4. (B) A schematic  
 197 illustration of the SNNs and full-replacement models. In the full-replacement models,  
 198 processing layers were added, the filters in the initial layer were changed with Gabor filters, and  
 199 the range of pooling was narrowed. (C) DoG filters for the SNNs and DoG models (upper) and  
 200 Gabor filters for the Gabor models (lower). (D) The pooling range for the SNNs ( $5 \times 5$ ) and the  
 201 narrow-pooling models ( $3 \times 3$ ). (E) Examples of presented face images with seven expressions  
 202 (angry, disgusted, fearful, happy, sad, surprised, neutral) of two individuals (upper: female,  
 203 AF06; lower: male, AM27). The original images were obtained from the Karolinska Directed  
 204 Emotional Faces database (Lundqvist et al., 1998).

205

20230127

206 **Table 1.** Architecture of the Shallow Neural Network (SNN). Each row describes a layer  $i$  with  
 207 calculation operator  $F_i$ , output resolution  $H_i \times W_i$ , and the number of output channels  $C_i$ . Conv  
 208 denotes convolution layer, and Pool denotes max pooling layer.  
 209

| Input &<br>Layer<br>$i$ | Operator<br>$F_i$   | Resolution<br>$H_i \times W_i$ | #Channels<br>$C_i$ |
|-------------------------|---|--------------------------------|--------------------|
| Input                   |   | $227 \times 227$               | 3                  |
| Layer 1                 | $11 \times 11$ Conv (stride 4) & $5 \times 5$ Pool (stride 4) | $55 \times 55$                 | 32                 |
| Layer 2                 | $5 \times 5$ Conv (stride 1) & $5 \times 5$ Pool (stride 4)   | $4 \times 4$                   | 64                 |
| Layer 3                 | Fully connected   | $1 \times 1$                   | 32                 |
| Layer 4                 | Fully connected   | $1 \times 1$                   | 7                  |

210

211

212

213 The DoG-type filters of the first convolution layer were built using the following formula:

214

$$215 \text{DoG}(r) = \pm A_1 \exp\left(-\frac{r^2}{2\sigma_1^2}\right) \mp A_2 \exp\left(-\frac{r^2}{2\sigma_2^2}\right), \quad (1)$$

216

217 where  $r$  is the polar radius from the filter center,  $A_1$  and  $A_2$  are the amplitudes of exponentials of  
 218 two Gaussian functions, and  $\sigma_1$  and  $\sigma_2$  are the standard deviations. Values of  $A_1$ ,  $A_2$ ,  $\sigma_1$ , and  $\sigma_2$   
 219 were chosen empirically so that DoG curves took the shapes of Mexican hats.  $A_1$  values were  
 220 0.4, 0.67, 0.8, and 1.0.  $A_2$  values were determined based on  $A_1 - A_2 = 0.4$ . When  $A_1$  was 0.4 (i.e.,  
 221  $A_2$  is 0, and  $\sigma_2$  cannot be defined), we set the  $\sigma_1$  value at  $1/2\sqrt{2}$ ,  $1/4\sqrt{2}$ ,  $1/8\sqrt{2}$ , or  $1/16\sqrt{2}$ .  
 222 Otherwise, the  $\sigma_1$  value was  $1/2\sqrt{2}$  or  $1/4\sqrt{2}$ . The  $\sigma_2$  value was based on  $\sigma_1/\sigma_2 = 0.5$  or 0.25.  
 223 The same number of filters was generated for each  $A_1$  value.

224

225 We also constructed modified models in which the three properties of the SNNs were replaced  
 226 one-by-one, two at a time, or all three simultaneously with the corresponding properties in the  
 227 ventral cortical pathway. First, we added additional convolution layers with filters of  $3 \times 3$  pixels  
 228 after each of the first two convolution layers to increase the number of processing stages (add-  
 229 layer model). In adding the convolution layers, the stride of sliding filters was reduced to 1 to  
 230 keep the output resolutions unchanged before and after adding the new layers. Also, to keep the  
 231 number of the output channels unchanged, the new layers contained the same number of filters  
 232 as the preceding layers.

233

234 Second, we replaced the DoG-type filters with Gabor-type filters (Gabor model). Gabor-type  
 235 filtering occurs in simple cells of V1, and emerges in the first layer of DNNs after they are  
 236 trained to classify object images (Krizhevsky et al., 2012; Rai and Rivas, 2022). We constructed  
 237 the Gabor-type filters with the following formula:

238

$$239 g(x, y; f, \theta) = A \exp\left(-\frac{x^2+y^2}{2\sigma^2}\right) \exp(2\pi i f(x \cos \theta + y \sin \theta)), \quad (2)$$

240

20230127

241 where  $A$  is the amplitude of the Gaussian envelope,  $i$  is  $\sqrt{-1}$ ,  $f$  is the carrier frequency of a Gabor  
242 filter, and  $\theta$  is the orientation (Movellan, 2002).  $A$  was fixed at 0.4 to match the amplitude of  
243 the DoG filters.  $\sigma$  was fixed at 0.125 so that the half-amplitude width was half of the filter width.  
244  $f$  values were 2 or 4 cycles/image. Orientation  $\theta$  was 0, 22.5, 45, 67.5, 90, 112.5, 135, or 157.5.  
245 We built even- and odd-symmetric filters for every combination of variables. In total, we  
246 obtained 32 Gabor-type filters (Fig. 1C, bottom).

247

248 Finally, we made the pooling window (convergence field) in the max pooling layers smaller ( $3$   
249  $\times 3$ ; Fig. 1D, bottom) than that of the SNNs ( $5 \times 5$ ; Fig. 1D, top), enabling better spatial  
250 resolution of processing (narrow-pooling model) to mimic the smaller receptive fields in the  
251 visual cortices compared to the superior colliculus and pulvinar (Wallace et al., 1997; Vand den  
252 Bergh et al., 2010; Freeman and Simoncelli, 2011). The pooling range of  $3 \times 3$  is often used in  
253 DNNs (e.g., AlexNet of Krizhevsky et al., 2012; ResNet of He et al., 2016).

254

### 255 ***Face images and training of the SNNs.***

256 Face images were obtained from Karolinska Directed Emotional Faces (KDEF; Lundqvist et al.,  
257 1998) and Radboud Faces Database (RaFD; Langer et al., 2010). Images of the seven  
258 expressions of 40 individuals (half females, half males) were chosen from each database (the  
259 total number of images was  $560 = 7 \times 40 \times 2$ ). We converted the images from color into  
260 grayscale, and extracted the face region by removing hair, neck, and ears with the face-detection  
261 function of a computer vision library, OpenCV (Open Source Computer Vision Library; Bradski,  
262 2000). The isolated faces were pasted on a gray background ( $198 \times 198$  pixels; RGB values =  
263 128; Fig. 1E). We augmented the number of face images by changing size and position, and by  
264 flipping horizontally; seven sizes ( $28 \times 28$ ,  $56 \times 56$ ,  $85 \times 85$ ,  $113 \times 113$ ,  $141 \times 141$ ,  $170 \times 170$ ,  
265 and  $198 \times 198$  pixels), five positions (center, left-top, right-top, left-bottom, and right-bottom;  
266 directional displacements = 10 pixels), and two horizontally flipped images. The augmentation  
267 increased the number of images by 70 times to 39,200. At each training session, we randomly  
268 split this augmented set of face images into a training set (29,400 images), a validation set (2,450  
269 images), and a test set (7,350 images). The number of images per facial expression was identical  
270 within each of these stimulus sets. To avoid the inadvertently biased assignment of face images  
271 of a particular size, position, or horizontal flip state into a given set, all images from the same  
272 individual were assigned into the same set.

273

274 The training was performed through supervised learning, and was conducted individually 20  
275 times with randomized initial weights except for the built-in weights of the first convolution layer,  
276 i.e., 20 SNNs with different initial states were built. In training, the weights other than the first  
277 convolution layer were optimized for classification of face images into the seven categories.  
278 Stochastic gradient descent was used for weight optimization. For each iteration, 32 samples  
279 were randomly selected from the training set as a mini-batch. The averaged cross-entropy  
280 (Goodfellow et al., 2016) across the 32 images in a mini-batch was calculated as an estimate of  
281 loss value, which is a measure of the difference between a model output and a supervised signal,  
282 and is used for quantifying the training effect. The number of iterations (i.e., weight-updating  
283 processes with single mini-batches) was set at 240,000. Initial weight parameters followed a  
284 normal distribution with a mean of 0 and a standard deviation of  $\sqrt{(2/N)}$  ( $N$  is the total number  
285 of weights; He et al., 2015). Weights were updated at each iteration with a constant learning rate



20230127

286 of 0.001. This learning rate was determined empirically; a preliminary analysis based on 10  
287 constructed SNNs (different from the 20 SNNs in the main analysis) revealed no decrease in  
288 loss values (i.e., no learning) with a learning rate of 0.01, which has frequently been used for  
289 DNNs in the literature (e.g., Simonyan and Zisserman, 2014). A dropout process was added  
290 before FC2 to facilitate learning across all units. The proportion of units dropped out of each  
291 weight update was set to 0.5. The training was conducted in a Python environment (Chainer  
292 3.0.0; Tokui et al., 2015) on a graphics processing unit (GPU) machine (Intel® Core™ i7-5820K  
293 Processor. Intel, Santa Clara, CA, USA; The GeForce® GTX 1080 Ti, NVIDIA, Santa Clara,  
294 CA, USA). While the SNNs were being trained with the training set, the correct rate and loss  
295 value for the validation set were periodically checked to monitor signs of overfitting. After  
296 training was completed, the performance of the models was evaluated using the test dataset that  
297 had not been used for training. This was done to ensure that the models acquired a genuine  
298 ability to classify the facial expressions, as opposed to simply sorting the training images into  
299 the seven facial expression categories according to the instruction signals.

300

### 301 ***Test for reference frames of SF tuning of model units.***

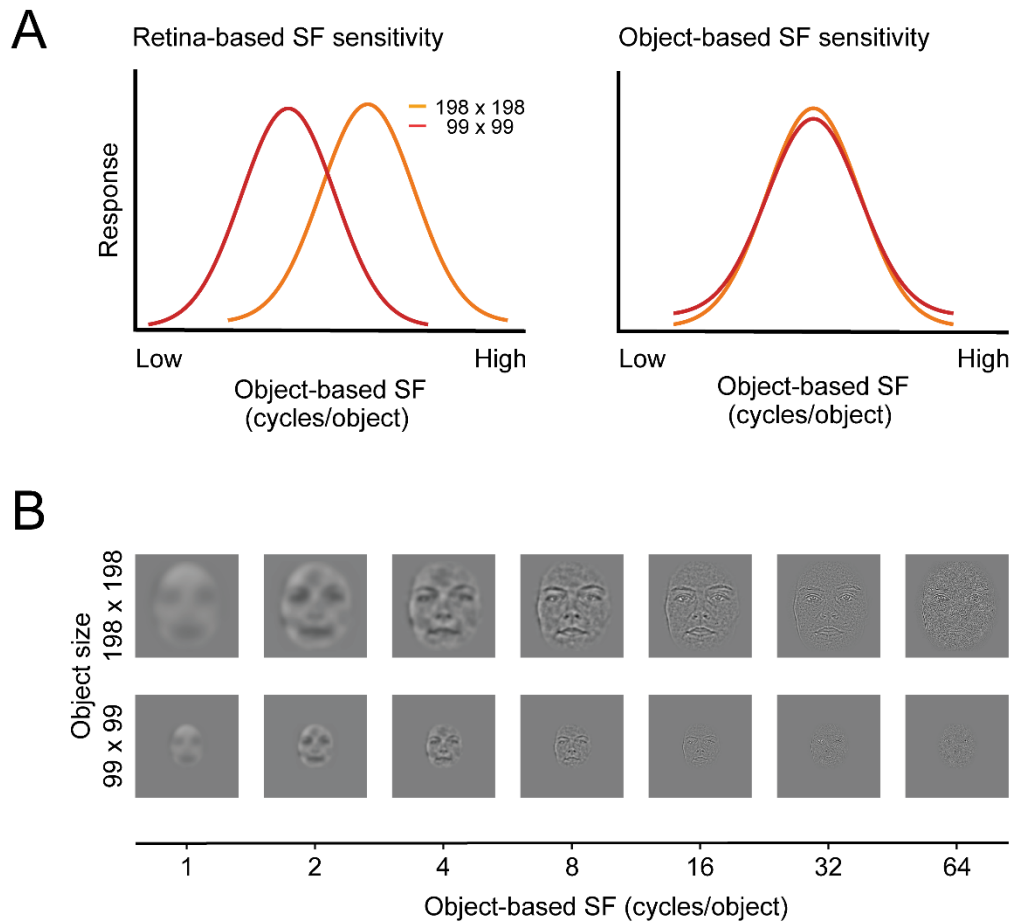
302 A difference in visual responses between the two pathways is the reference frame of neuronal  
303 tuning to SFs (Inagaki and Fujita, 2011). Neurons in the inferior temporal cortex, the final stage  
304 of the ventral cortical pathway, are tuned to object-based SFs (cycles/object) and represent face  
305 patterns in a size-invariant, hence distance-invariant, manner (Fig. 2A, right). Thus, the ventral  
306 cortical pathway converts the representation of SFs in the retina-based coordinate  
307 (cycles/degree) to that of object-based SFs. By contrast, many amygdala neurons preserve  
308 sensitivity to retina-based SFs. When the stimulus size is changed, these neurons change their  
309 preferred object-based SFs; for large stimuli, they respond to higher object-based SFs, which  
310 correspond to the same retina-based SFs (Fig. 2A, left). We analyzed the reference frame of  
311 FC1 SF tuning to evaluate how well our models captured this characteristic of subcortical  
312 processing.

313

314 Bandpass-filtered face images were used to examine the SF tunings of FC1 units (Fig. 2B).  
315 These images were created by multiplying Gaussian functions with the original face images on  
316 the polar Fourier domain. Gaussian functions had 61 different center frequencies between 1  
317 cycle/object and 64 cycles/object. The center frequencies had discrete values at steps of 0.1  
318 cycles/object on a log scale. Gaussian functions shared the same variance at 2.4 octaves,  
319 regardless of their center frequencies. The filtered images had amplitude spectra that were  
320 determined solely by the Gaussian function because their spectra were set to be flat before the  
321 multiplication. To balance the total luminance contrast among the filtered face images, the peak  
322 amplitude of the Gaussian function was set inversely proportional to the center image-based SF  
323 (Inagaki and Fujita, 2011). These bandpass-filtered images were created for the seven facial  
324 expressions at two different sizes ( $99 \times 99$  and  $198 \times 198$  pixels).

325

20230127



326

327

328

329

330

331

332

333

334

335

336

**Figure 2.** Test determining whether units are tuned for object- or retina-based spatial frequencies (SFs). (A) Hypothetical tuning curves for object-based SFs of units ideally tuned to retina-based SFs (left: peak shift = 1) or object-based SFs (right: peak shift = 0). (B) The models were fed face images with two different sizes (198 × 198 and 99 × 99 pixels) and 64 different bandpass filtering. These images were created by applying two-dimensional bandpass filters that shared the same center object-based SF across different sizes. For each unit in FC1, we obtained responses to images of different center SFs to create tuning curves for object-based SFs.

20230127

337 To characterize the reference frame of SF tuning of each unit, the peak SFs for the two stimulus  
338 sizes were estimated, defined by the SFs at which filtered face stimuli activated a unit most  
339 strongly. For a given unit, 14 peak SFs were determined (two sizes  $\times$  seven expressions). The  
340 degree to which unit responses to SFs depended on the stimulus size was quantified by  
341 calculating differences in peak SFs on a log scale between the two face sizes. A peak shift of 0  
342 indicates that a unit respond to the same cycles/object regardless of the image size, and is  
343 perfectly tuned to object-based SFs (Fig. 2A, right). A peak shift of 1 means that an SF tuning  
344 curve shifts by the amount corresponding to the change in the stimulus size, indicating that a  
345 unit is perfectly tuned to retina-based SFs (Fig. 2A, left). This analysis excluded cases in which  
346 units did not respond to face images or were not sensitive to SFs, and cases in which peak SFs  
347 were at either end of the tested range of SFs and the peak positions could not be determined.

348

#### 349 *Analysis of the effects of the max pooling operation on SF selectivity.*

350 The max pooling operation collapses positional information of edges, which is detected by  
351 convolution filters and is critical for encoding the SFs of facial images. We therefore examined  
352 the effects of the max pooling on the SF selectivity of units. We were particularly interested in  
353 the role of the max pooling in converting sensitivity to retina-based SFs to sensitivity to object-  
354 based SFs. Bandpass-filtered images of the two stimulus sizes ( $198 \times 198$  and  $99 \times 99$  pixels)  
355 were fed to the models, and it was determined how different stimulus sizes affected the  
356 responses of units to SFs in the first convolution layer (before pooling) and the first max pooling  
357 layer (after pooling). Changes in response patterns across the units associated with different  
358 stimulus sizes were quantified by calculating the dissimilarity index. The dissimilarity index  $D$   
359 ( $x, y$ ) for responses  $x$  to the large stimuli and responses  $y$  to the small stimuli was defined by the  
360 Euclidean distance between  $x$  and  $y$  as follows:

361

$$362 \quad D(x, y) = \|x - y\| / (NM) \quad (3)$$

363

364 where  $\|\cdot\|$  is the Euclidean distance, and  $N$  is the number of elements of  $x$  and  $y$ .  $M$  is the  
365 maximum value among the 6,405 Euclidean distances calculated for 61 center SFs and the seven  
366 facial expressions of 15 individuals. To probe the roles of the max pooling, the ratio of the  
367 dissimilarity index before pooling in the first convolution layer and after pooling in the first max  
368 pooling layer was then calculated. This analysis was applied to the case of wide pooling ( $5 \times 5$ )  
369 and narrow pooling ( $3 \times 3$ ), as well as to the case of the SNNs and the Gabor models. By dividing  
370  $\|x - y\|$  by  $M$ , the dissimilarity index was normalized across the layers (convolution vs. max  
371 pooling) and the models (SNNs vs. Gabor models), taking values from 0 to 1.

372

373

## 374 **Results**

375

### 376 *Performance of SNNs in facial expression classification.*

377 The SNNs were trained to classify each face image in the training set into one of the seven facial  
378 expressions. The training improved the classification performance rapidly over the initial  
379 iterations and then slowly thereafter. The correct rate across the seven facial expressions rose  
380 from the chance level (0.14), surpassed 0.6 around 50,000 iterations, further improved to around  
381 0.8 over 150,000 iterations, and reached an asymptote (Fig. 3A for an example SNN; 3B for the

20230127

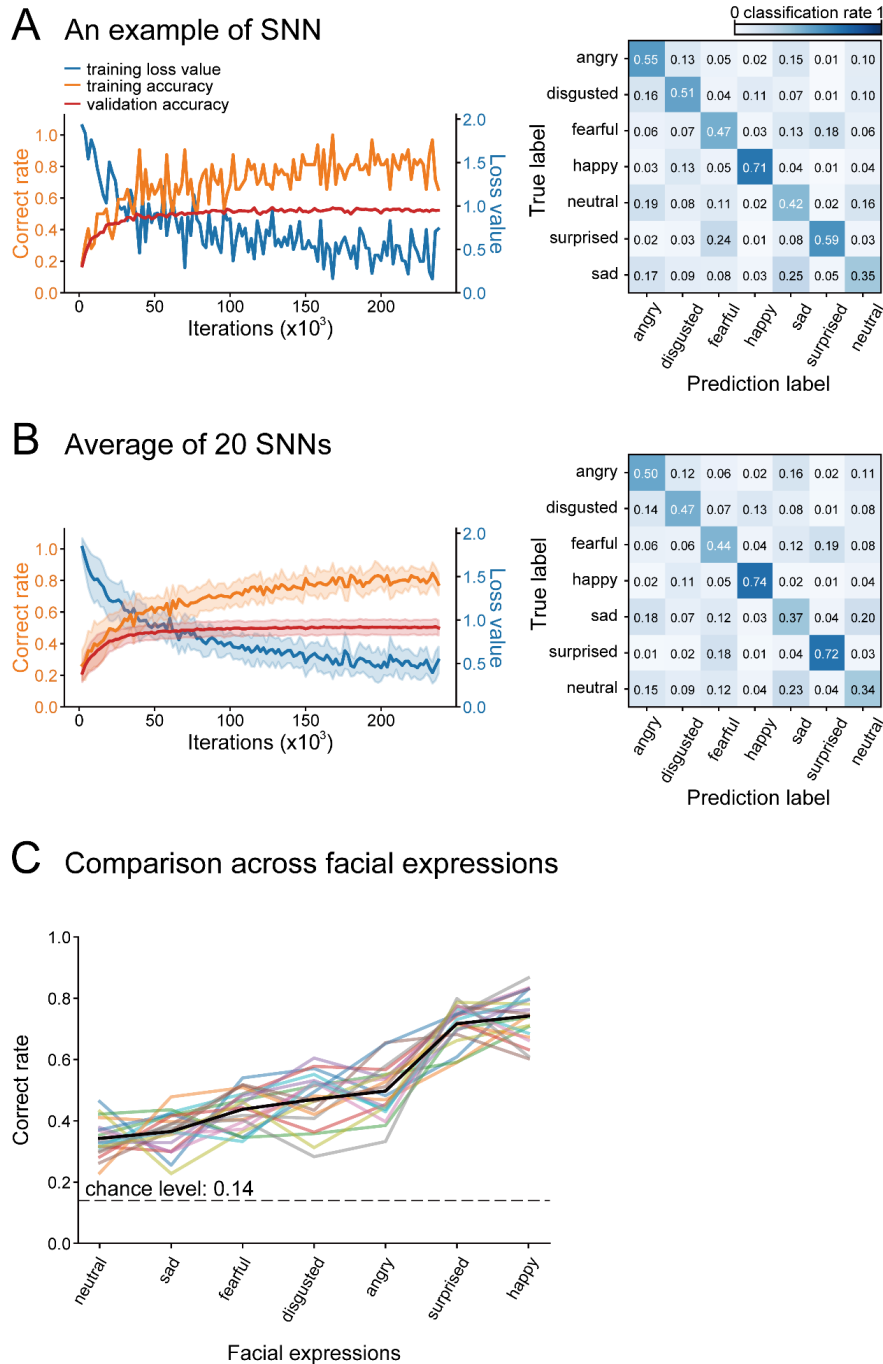
382 average of the 20 constructed SNNs; orange lines). The correct rate for the validation set  
383 saturated at around 0.5, which was substantially lower than for the training set, indicating  
384 insufficient generalization to “unseen” images. However, the validation correct rate reached a  
385 plateau in a similar way to the training correct rate. This indicates that the low correct rate was  
386 not the result of inadequate training, but represents the limited learning ability of the SNNs. It  
387 also indicates that no overfitting occurred. The loss value also quickly decreased over the initial  
388 50,000 iterations, and became gradually stable (Fig. 3A, B; cyan lines). The results indicate that  
389 within the range of adopted iterations (240,000), the SNNs were trained to classify facial  
390 expressions without overfitting.

391

392 The correct rates for the test set were higher than the chance level ( $1 / 7 = 0.14$ ) for all facial  
393 expressions, but were modest; the average correct rate across the seven expressions was 0.51  
394 for the example SNN shown in Fig. 3A. The average correct rate across the 20 constructed SNNs  
395 was  $0.51 (\pm 0.03, \text{s.d.})$ . This was not different from the average correct rate across 20 additional  
396 SNNs that were trained with 3,000,000 iterations ( $0.50 \pm 0.03; p = 0.289, t\text{-test}$ ). The training  
397 performance thus did not improve even when the SNNs underwent overly excessive training,  
398 assuring that the modest correct rate was not due to insufficient training but instead reflected the  
399 limited ability of the SNNs. Confusion matrices showed that the correct rates varied among the  
400 facial expressions (Fig. 3A, B, right panels). Based on the averaged performance, the  
401 classification performance of the SNNs was best for happy (0.74) and surprised faces (0.72),  
402 followed by angry (0.50), disgusted (0.47), and fearful (0.44) faces, and was worst for sad (0.37)  
403 and neutral (0.34) faces. Sad faces were often confused with neutral, angry, and fearful faces.  
404 Neutral faces were often confused with sad and angry faces. This expression-dependent  
405 performance was consistent across the 20 constructed SNNs (Fig. 3C;  $p < 0.001$  for expressions,  
406  $p = 0.562$  for models, two-way ANOVA).

407

20230127



408

409

410

**Figure 3.** Learning curves and confusion matrices of an example SNN (A) and the average of

20 SNNs with different initial weights (B). Left panels show changes in correct rates that

occurred during training in the training set (orange) and the validation set (red), and loss values

(cyan). Confusion matrices on the right indicate the rate of classification of each facial

expression (true label) as one of the seven expressions (prediction label). (C) The correct rates

for the seven expressions. Black line indicates the mean of the 20 SNNs, and lines with other

colors indicate the individual performance of the 20 SNNs. The order of facial expression is

based on the mean correct rate.

418

20230127

419 ***Effects of modification of SNNs on classification performance.***

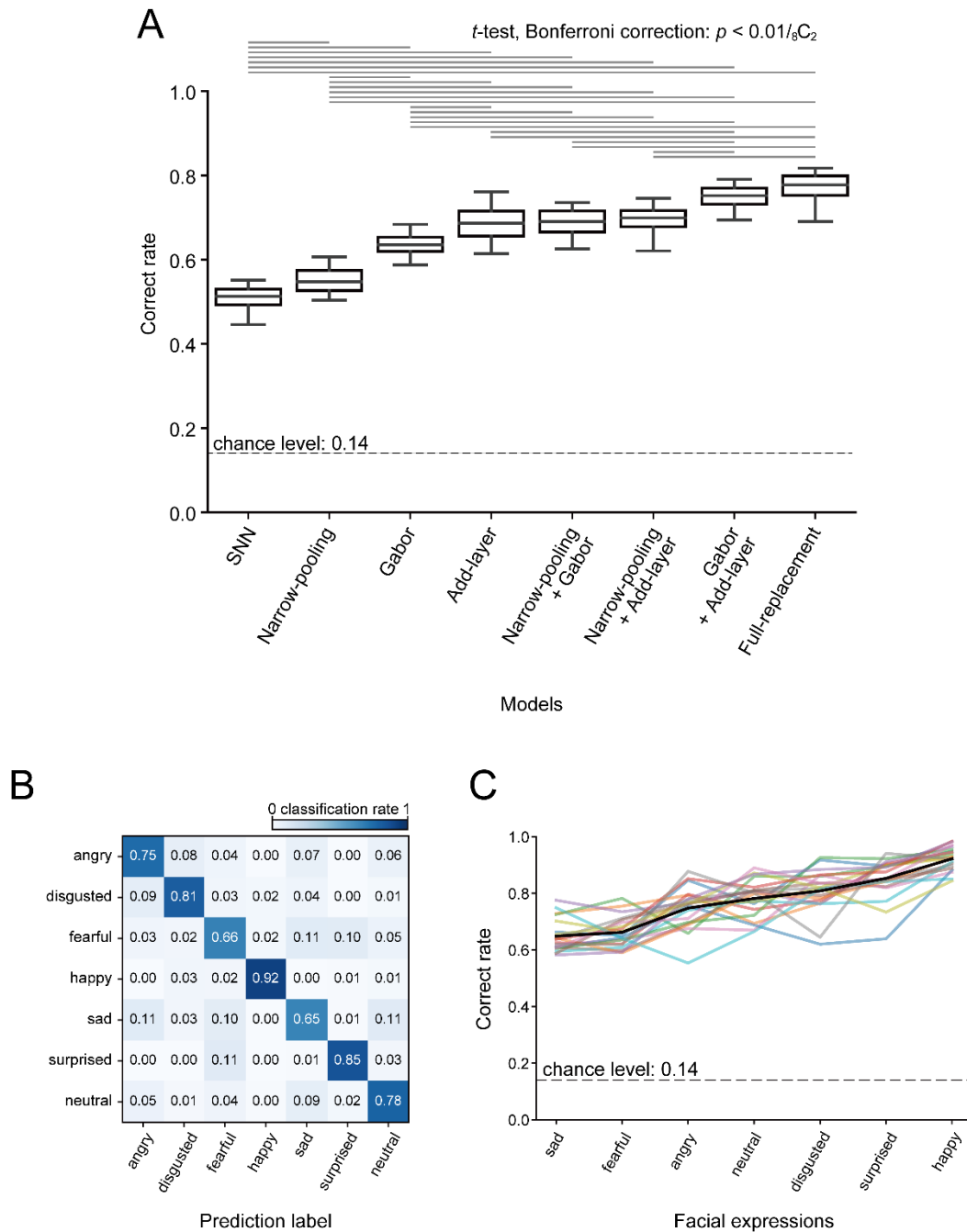
420 Replacement of one or more of the three subcortical properties, namely the shallowness of  
421 processing stages, the DoG-type receptive fields at the initial stage, and spatial pooling over a  
422 wider visual field, with the corresponding cortical properties improved the classification  
423 accuracy (Fig. 4A;  $p < 0.001$ , ANOVA). The correct rates averaged across the seven expressions  
424 and the 20 constructed models of each modification were increased from 0.51 in the SNNs to  
425 0.55 in the narrow-pooling models, 0.64 in the Gabor models, and 0.69 in the add-layer models  
426 ( $p < 0.01 / 28=8C_2$ ;  $t$ -test with Bonferroni correction). Among these models with one replaced  
427 property, the add-layer models exhibited the best performance. The results indicate that all three  
428 properties had effects on classification performance, and the layer structure was the most  
429 influential.

430

431 When two properties were replaced together, the narrow-pooling + Gabor models and the  
432 narrow-pooling + add-layer models performed better than the narrow-pooling models and the  
433 Gabor models ( $p < 0.01 / 28=8C_2$ ;  $t$ -test with Bonferroni correction), but comparably to the add-  
434 layer models (correct rate, 0.69,  $p = 0.778$  for the narrow-pooling + Gabor models; 0.69,  $p =$   
435 0.564 for the narrow-pooling + add-layer models). The Gabor + add-layer models performed  
436 better than all one-property-replacement models (correct rate, 0.75;  $p < 0.01 / 28=8C_2$ ). When  
437 all three properties were replaced together (full-replacement models), the correct rate was 0.77,  
438 which was better than all other models ( $p < 0.01 / 28=8C_2$ ) except for the Gabor + add-layer  
439 models ( $p = 0.00846 > 0.01 / 28=8C_2$ ). The performance was improved for all facial expressions  
440 (Fig. 4B; mean correct rates: happy = 0.92; surprised = 0.85; disgusted = 0.81; neutral = 0.78;  
441 angry = 0.75; fearful = 0.66; sad = 0.65). As in the SNNs, it was highest for happy and surprised  
442 faces, and lowest for fearful and sad faces. The performance was improved most for neutral  
443 faces (SNNs, 0.37; full-replacement models, 0.78). The variance of the correct rates was affected  
444 both by facial expressions and models (Fig. 4C;  $p < 0.001$  for facial expressions,  $p = 0.00285$   
445 for models, two-way ANOVA). The improved performance of the two-property-replacement  
446 and full-replacement models indicate that the effects of the three features on classification  
447 performance were partially additive, suggesting that the three features exerted their effects  
448 partially independently.

449

20230127



450

451

452

453

454

455

456

457

458

459

460

**Figure 4.** Effects on model performance of replacing subcortical properties with corresponding cortical properties. (A) Discrimination performances of the SNNs and the modified models. The discrimination performance differed across the models (ANOVA,  $p < 0.001$ ). The pairs of models with statistically significant differences in the performances are linked with horizontal lines in the upper part (*t*-test, Bonferroni correction,  $p < 0.01/28=8C_2$ ). (B) Confusion matrix for the full-replacement models (average of the 20 constructed models). (C) The correct rate for the seven expressions across the 20 full-replacement models. The black line indicates the mean, and lines with other colors indicate the data for individual full-replacement models.

20230127

461 ***Spatial frequency representation in the FC1 layer.***

462 As shown above, the SNNs exhibited modest performance in facial expression classification,  
463 and this performance was improved by changing SNN subcortical properties to corresponding  
464 cortical properties. These findings suggest that the SNNs captured aspects of processing in the  
465 subcortical pathway to the extent that they explained the suboptimal perceptual performance of  
466 V1-lesioned patients. We next looked into individual computational units to gain insights about  
467 the processing in the models. We examined SF sensitivities of FC1 units using two different  
468 sizes of input images ( $198 \times 198$  and  $99 \times 99$  pixels). This procedure allowed us to determine  
469 whether units were sensitive to retina- or object-based SFs (Fig. 2; see Materials and Methods).  
470

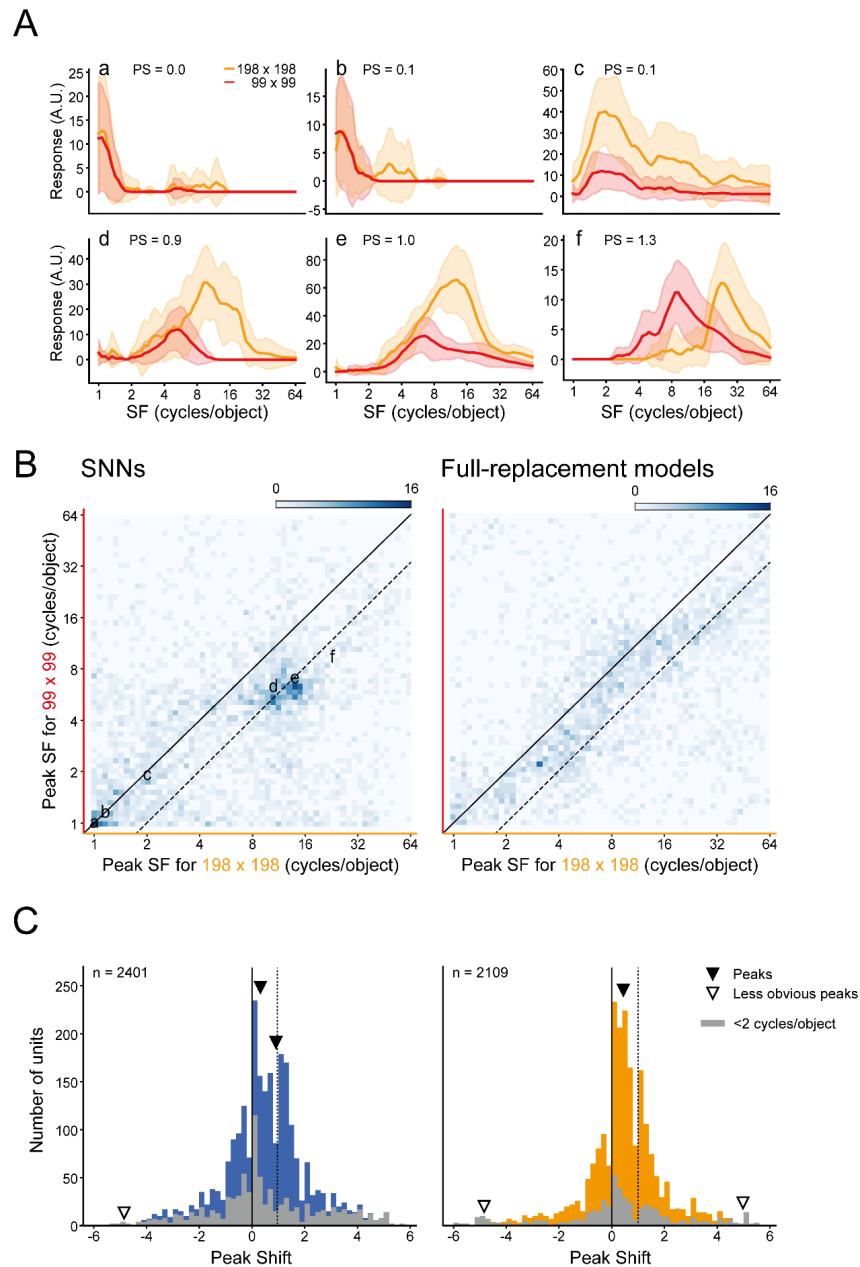
471 FC1 units of the SNNs exhibited a variety of dependencies of SF tunings on stimulus size (Fig.  
472 5A). Some units responded to the same range of object-based SFs for both large and small  
473 stimuli, and the peak positions of the SF tuning curves remained unchanged (Fig. 5Aa, Ab).  
474 Other units exhibited different preferred SFs for large and small stimuli, and in these cases the  
475 peak position shifted horizontally along the abscissa (Fig. 5Ac–f). We quantified these shifts by  
476 measuring the difference between preferred SFs on a log scale for the two stimulus sizes. A peak  
477 shift of 0 means that the unit encoded SFs in the object-based coordinate, whereas a peak shift  
478 of 1 means that the unit encoded SFs in the retina-based coordinate. The peak shifts of the  
479 example units shown in Fig. 5A were 0.0 (a), 0.1 (b), 0.1 (c), 0.9 (d), 1.0 (e), and 1.3 (f).  
480

481 We plotted the peak positions of 2,401 FC1 units of the 20 SNNs in a two-dimensional space  
482 defined by the peak SF for the large stimuli on the abscissa, and the peak SF for the small stimuli  
483 on the ordinate (Fig. 5B, left). Note that 46% of FC units were excluded from this analysis,  
484 either because they were not sensitive to SFs (23%) or because the largest responses were found  
485 at the end of the examined range of SFs and the peak SFs could not be determined (23%). The  
486 diagonal solid line in Fig. 5B represents the responses of a peak shift of 0, and the dashed line  
487 next to it represents the responses of a peak shift of 1. FC1 units of the SNNs were clustered in  
488 multiple groups in this scatter plot. One conspicuous group was selective to low SFs and was  
489 centered on the diagonal, i.e., peak shift values around 0. Another group was selective to higher  
490 SFs, and was clustered on the dashed line indicative of peak shift values around 1. The  
491 multimodality of the distribution can also be seen in the histogram (Fig. 5C, left). We applied  
492 an excess mass test for multimodality (Ameijeiras-Alonso et al., 2019, 2021) to this distribution.  
493 This test statistically determines the number of peaks in the distribution, with the null hypothesis  
494 that the true number of peaks is  $N$  ( $N = 1, 2, 3, \dots$ ). The true number of peaks is estimated as the  
495 smallest  $N$  under which the null hypothesis is not rejected. The excess mass test also estimates  
496 the locations and heights of peaks from Gaussian kernel density estimation. The test revealed  
497 that there were three peaks in the distribution of the SNNs (first  $p$ -value  $< 0.001$ , second  $p$ -value  
498  $< 0.001$ , third  $p$ -value = 0.096). Based on the probability density function derived from the  
499 histogram (Ameijeiras-Alonso et al., 2021), the peaks were estimated to be located at  $-4.85$ ,  
500  $0.144$ , and  $0.909$  (open and solid arrowheads in Fig. 5C, left). Units sensitive to low SFs below  
501 2 cycles/object were most frequent around a peak shift of 1 (gray columns). Comparing this  
502 result and the density map, the peak around 0 was mostly from the low spatial frequency group  
503 and the peak around 1 was from the high spatial frequency group. Although the third peak at the  
504 far periphery (at  $-4.85$ , open arrowhead) was statistically detected, it was much smaller in height  
505 than the other two peaks (1.1% of the peaks near 0 and 1). The results indicate that FC1  
506 contained two major groups of units, those sensitive to low SFs, encoding SFs in the object-  
507 based coordinate, and those sensitive to high SFs, encoding SFs in the retina-based coordinate.



20230127

508



509  
510

511 **Figure 5.** SF tuning reference frames of FC1 units of the SNNs and full-replacement models.  
512 Responses of FC1 units to SF-filtered face images were examined at two different sizes ( $198 \times$   
513  $198$ ,  $99 \times 99$  pixels). (A) Six example FC1 units of SNNs with a different peak shift (PS). (B)  
514 Two-dimensional histograms of peak SFs at large images versus small images for the SNNs  
515 (left) and the full-replacement models (right). Solid lines indicate peak shifts of 0, and dashed  
516 lines indicate peak shifts of 1. (C) Distribution of peak shifts of units in the 20 SNNs (left) and  
517 the 20 full-replacement models (right). Arrowheads indicate the estimated locations of multiple  
518 peaks in the distribution (solid: major peaks, open: statistically detected but less obvious peaks).  
519 Gray columns indicate units with a response peak at SFs below 2 cycles/object for large and/or  
520 small stimuli.

20230127

521 The distribution of peak shift values was drastically altered in the full-replacement models. In  
522 the two-dimensional plot shown in Fig. 5B (right), most data points were diffusely distributed  
523 in an elongated area between the diagonal and dashed lines, indicating that the SF reference  
524 frame of most units was intermediate between retina-based and object-based. An excess mass  
525 test again detected three peaks located at  $-4.84$ ,  $0.436$ , and  $4.98$  (Fig. 5C right, solid and open  
526 arrowheads; first  $p$ -value  $< 0.001$ , second  $p$ -value  $< 0.001$ , third  $p$ -value =  $0.096$ ). The second  
527 and third peaks at  $-4.84$  and  $4.98$  (open arrowheads) were smaller than the primary peak at  
528  $0.436$  (3.1% and 2.6% of the primary peak, respectively), making the distribution nearly  
529 unimodal.

530

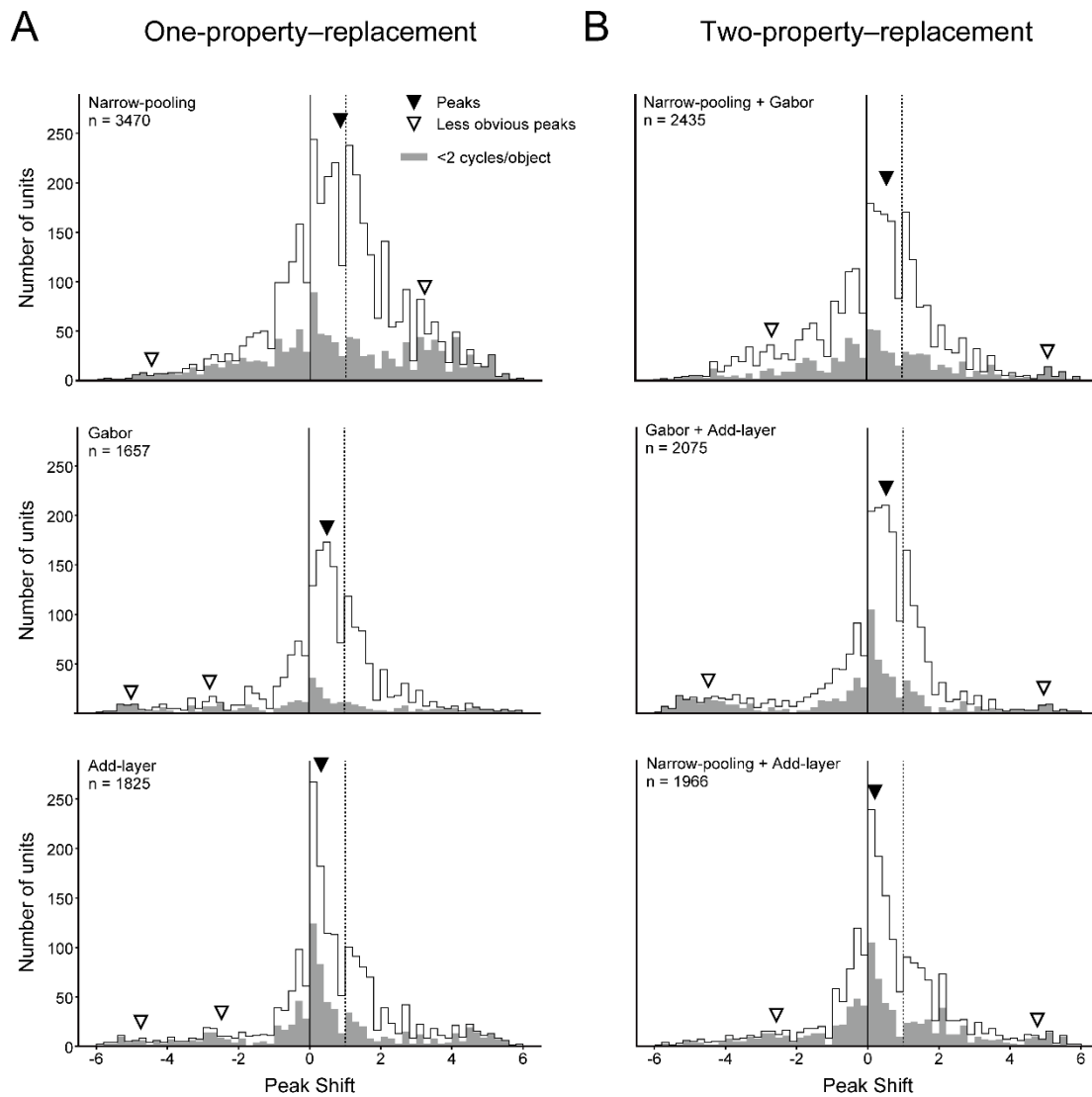
531 Given the change of the peak shift distribution in the full-replacement models, we next analyzed  
532 one- and two-property replacement models to determine which subcortical properties were  
533 essential for the multimodal distribution of the SNNs. All these modified models exhibited  
534 unimodal distributions of the major peak (solid arrowheads) at different peak positions (Fig. 6).  
535 An excess mass test for multimodality detected two other less obvious peaks (open arrowheads)  
536 in each model as in the cases of the SNNs and the full-replacement models (Fig. 5C). The heights  
537 of these smaller peaks were 2.6–26 % of those of the major peaks, and were located at the  
538 periphery of the distribution.

539

540 Each of the three one-property–replacement models showed a characteristic distribution of the  
541 peak shift values. The narrow pooling models contained units with peak shift values between 0  
542 and 1 in addition to units with peak shift values around either 0 or 1. The distribution became  
543 unimodal and broad, and was estimated to be centered at 0.85. In the Gabor models, units with  
544 peak shift values intermediate between 0 and 1 were the most abundant with a smaller number  
545 of units of peak shift values around 0 and 1. The distribution peak was estimated at 0.51. In the  
546 add-layer models, units with peak shift values around 0 were predominant, and exhibited a sharp  
547 distribution peak at 0.32. As to the two-property–replacement models, the narrow-pooling +  
548 Gabor models and the Gabor + add-layer models showed a broad distribution straddling the  
549 peak values from 0 to 1 (peak for the former, 0.56; peak for the latter, 0.52), whereas the narrow-  
550 pooling + add-layer models showed a sharp distribution peak at 0.20. As in the SNNs and the  
551 full-replacement models, units sensitive to low SFs (below 2 cycles/object) were most  
552 frequently found around the peak shift of 0 in all of the one- and two-property replacement  
553 models (gray columns). The results indicate that all of the three computational properties were  
554 responsible for the multimodal distribution of peak shift values observed in the SNNs. In  
555 particular, the smaller number of units with peak shift values around 1 in the Gabor models and  
556 the add-layer models suggests that the shallowness and the DoG-type filters were critical for  
557 preserving the unit sensitivities to retina-based SFs. The broad distribution observed for the  
558 narrow-pooling models and the narrow-pooling + Gabor models suggests that the wide pooling  
559 employed in the SNNs contributed to the two peaks at 0 and 1, by reducing units with peak shift  
560 values intermediate between 0 and 1.

561

20230127



562

563

564

565

566

567

568

569

570

571

572

**Figure 6.** Distributions of peak shifts of FC1 units of the SNNs and one-property or two-property-replacement models. (A) Data from narrow-pooling model, add-layer model, and Gabor model. (B) Data from models with two modifications: narrow-pooling + Gabor, Gabor + add-layer, narrow-pooling + add-layer. Gray columns indicate units with a response peak at SFs below two cycles/object for large or small stimuli. Arrowheads indicate the estimated locations of multiple peaks in the distribution (solid: major peaks, open: statistically detected but small peaks). Gray columns indicate units with a response peak at SFs below 2 cycles/object for large and/or small stimuli.

20230127

573 ***Effects of max pooling on SF tuning.***

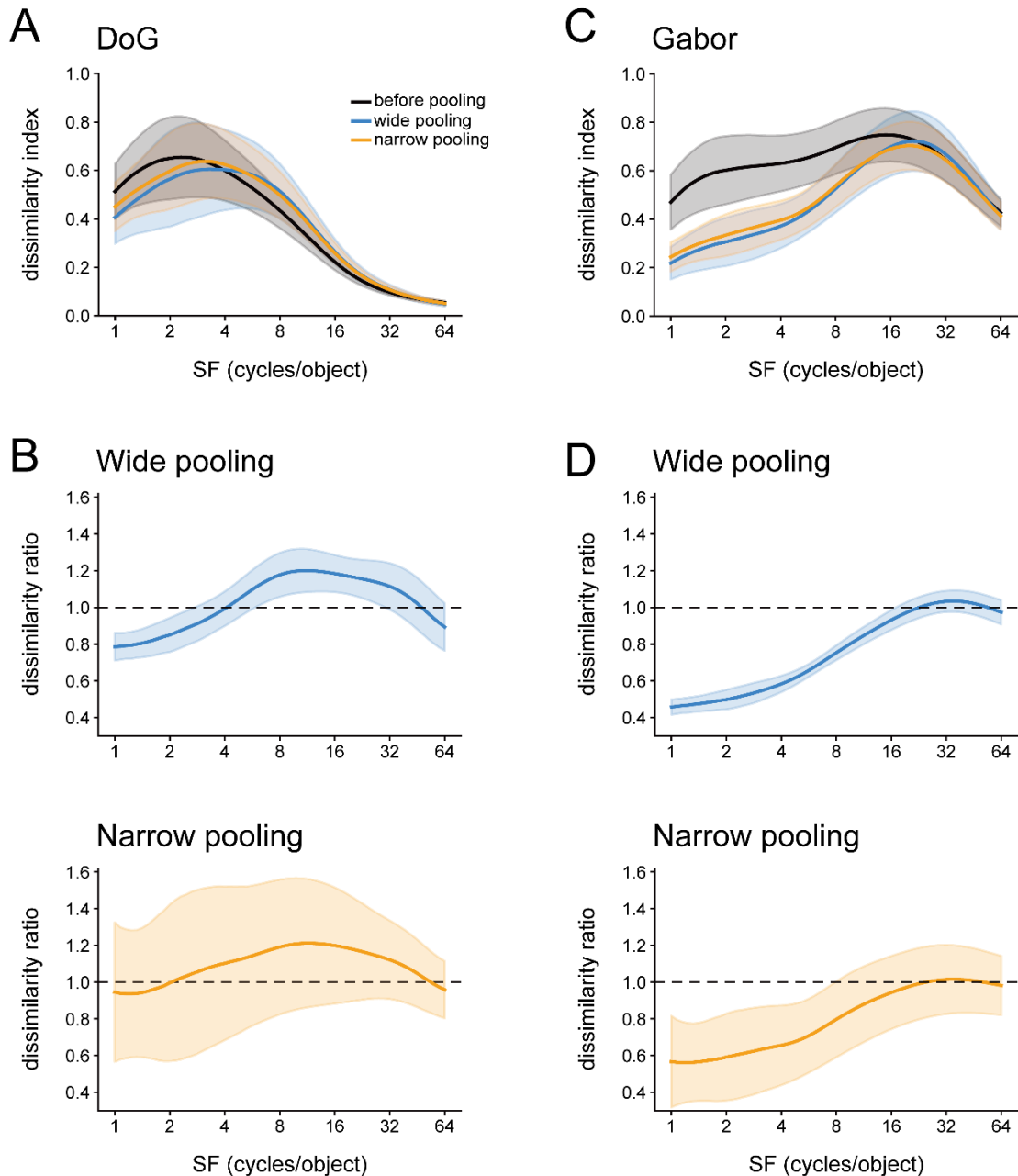
574 We showed above that FC1 units of the SNNs were roughly grouped into two populations in  
575 terms of the reference frame of SF encoding. Because the max pooling yields the same output  
576 from a population of convolution layer units in response to slightly different spatial arrangement  
577 of local features, the max pooling operation is likely to affect the encoding of global  
578 configuration of face components. This information of global configuration will be reflected in  
579 a low range of SFs. Therefore, we next compared the effect of max pooling on the representation  
580 of SFs across different SF ranges.

581  
582 We first analyzed the responses of the 96,800 units ( $32 \text{ filters} \times 55 \times 55$  resolution) in the first  
583 convolution layer. We obtained the response patterns across these units by feeding bandpass-  
584 filtered faces of two sizes ( $198 \times 198$  and  $99 \times 99$  pixels; Fig. 2B) to the models, and quantified  
585 the difference between the SF tunings obtained for the two stimulus sizes by calculating the  
586 dissimilarity index (see Materials and Methods). In the SNNs with DoG filters, the dissimilarity  
587 index was high (around 0.6) for a low SF range up to approximately four cycles/object, but  
588 gradually decreased over a higher range of SFs (Fig. 7A, black curve). In the pooling layer, the  
589 dissimilarity index of the 6,272 units ( $32 \text{ filters} \times 14 \times 14$  resolution) became lower for a low  
590 SF range of less than four cycles/object than that of the convolution layer. For a high SF range  
591 of greater than four cycles/object, by contrast, it became higher than that of the convolution  
592 layer (Fig. 7A, compare the orange and cyan curves with the black curve). Thus, max pooling  
593 resulted in the SF tuning becoming similar between the two stimulus sizes for a low SF range,  
594 consistent with the results of peak shift analysis (Fig. 5C). Although these changes were  
595 observed both for wide pooling ( $5 \times 5$ ) and narrow pooling ( $3 \times 3$ ), the effects were larger for  
596 the former than for the latter (Fig. 7A, compare the orange curve with the cyan curve). This was  
597 more evident when we plotted the ratio of dissimilarity indices between before and after pooling  
598 (Fig. 7B). Furthermore, the ratio curve for wide pooling had smaller standard deviations than  
599 that for narrow pooling (shown as shades in Fig. 7B), indicating that wide pooling exerted its  
600 effects more consistently across the 15 individual faces and the seven facial expressions than  
601 narrow pooling.

602  
603 By contrast, the convolution-layer units of the Gabor models exhibited a constantly high  
604 dissimilarity index over most of the SF range (Fig. 7C, black). However, when we applied max  
605 pooling with windows of either  $5 \times 5$  or  $3 \times 3$  in size, the dissimilarity index became small over  
606 almost the entire SF range, with the largest decrease for 1–16 cycles/object (Fig. 7C; compare  
607 the orange and blue curves with the black curve). As in the case of the SNNs, the effect was  
608 stronger for wide pooling than for narrow pooling (Fig. 7D). These results demonstrated that  
609 max pooling rendered the SF tuning more invariant to stimulus size for units sensitive to low  
610 SFs, enabling them to represent SFs in the object-based coordinate. Regardless of the filter type  
611 in the first convolution layer (i.e., DoG vs. Gabor), wide pooling was more effective than narrow  
612 pooling in creating this response property.

613

20230127



614

615

616

617

618

619

620

621

622

623

624

625

626

**Figure 7.** Effects of max pooling on the size-invariant responses to SFs. (A) Dissimilarity index curves of responses of units in the first convolution layer of the SNNs before max pooling operation (black), after wide pooling (blue), and after narrow pooling (orange). Dissimilarity indices, defined by the Euclidean distances of unit responses between different stimulus sizes (see Method and Methods), are plotted against the center SFs of input images. Solid lines indicate the means of dissimilarity indices across the seven facial expressions. Shades indicate standard deviations. Each dissimilarity index was normalized by the number of units and the maximum values. (B) Dissimilarity ratios of inputs and outputs of the max pooling operation (upper, after wide pooling; lower, after narrow pooling). (C, D) Data from units in the first convolution layer of the Gabor models. The conventions are the same as in A and B.

20230127

627 ***Effects of alternation of sliding strides on SF tuning.***

628 Finally, we examined the effect of another free parameter of our models, the stride size, on the  
629 SF sensitivity of FC1 units. We changed the stride of the two max pooling layers of the SNNs  
630 from 4 to 2. The stride size of 2 was also employed in the narrow-pooling model. This modified  
631 model with a smaller stride of 2 achieved a mean correct rate of 0.54, which was better than the  
632 SNNs (0.51) but similar to the narrow pooling models (0.55) (vs. SNN,  $p = 0.0020$ ; vs. the  
633 narrow pooling,  $p = 0.23$ ;  $t$ -test with Bonferroni correction).

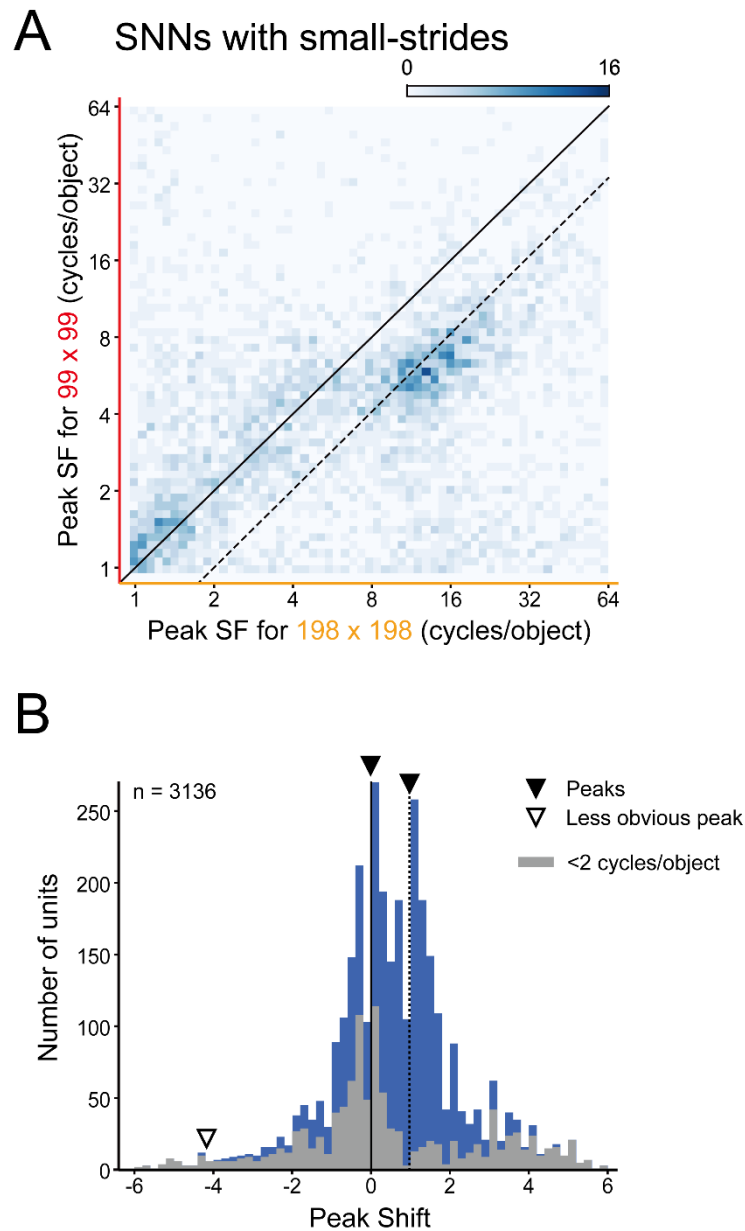
634

635 An analysis of FC1 responses to two stimulus sizes ( $198 \times 198$  and  $99 \times 99$  pixels) revealed that  
636 the distribution of the peak shift had three peaks, at  $-4.17$ ,  $-0.0107$ , and  $0.976$  (first  $p$ -value  $<$   
637  $0.001$ , second  $p$ -value  $< 0.001$ , third  $p$ -value =  $0.098$ ; excess mass test for multimodality; solid  
638 and open arrowheads in Fig. 8). Two of them were conspicuous and located near 0 or 1 (solid  
639 arrowheads), and the third one at the periphery of  $-4.17$  (open arrowhead) was small (3.2% and  
640 3.3% of the two major peaks). Comparisons with Fig. 5B, C show that the small-stride model  
641 exhibited a similar SF representation as in the SNNs, in that there were two main groups of units,  
642 one sensitive to low SFs, representing SFs in the object-based coordinate (peak shift around 0),  
643 and the other sensitive to high SFs, representing SFs in the retina-based coordinate (peak shift  
644 around 1). The change of the stride from 4 to 2 had little effects on the reference frame of SF  
645 sensitivity of FC1 units.

646

647

20230127



648

649

650 **Figure 8.** Effects of alternation of sliding strides of pooling windows on SF tuning reference

651 frames of FC1 units of the SNNs. Responses of FC1 units were obtained in the same way as in

652 Fig. 5. (A) A two-dimensional histogram of peak SFs obtained with large ( $198 \times 198$ ) versus

653 small ( $99 \times 99$ ) stimulus images. (B) Distribution of peak shifts of units. Arrowheads indicate

654 the estimated locations of multiple peaks in the distribution.

655

656

657

20230127

658 **Discussion**

659

660 We analyzed the ability of the SNNs and modified models to classify facial expressions, with  
661 the goal of determining what architectural or physiological properties underlie the modest  
662 performance of facial expression discrimination supported by the primate subcortical pathway.  
663 The SNNs were implemented with the three prominent subcortical properties, i.e., shallow  
664 processing, DoG-type filters at the first processing stage, and spatial pooling over wide areas  
665 (Fig. 1). The SNNs classified the seven basic facial expressions with modest performance (Fig.  
666 3). Replacement of any one of these properties with the corresponding cortical properties  
667 resulted in higher performances (Fig. 4). Replacement of a combination of two or three  
668 properties further improved classification performances in a partially additive manner (Fig. 4).  
669 These results suggest that all three subcortical properties of the SNNs underlie the modest  
670 performance. A major group of units in the final processing layer (FC1) of the SNNs was  
671 sensitive to SFs defined in the retina-based coordinate, whereas another group responding to  
672 low SFs encoded SFs in the object-based coordinate (Fig. 5). The number of retina-based units  
673 was reduced in most of the modified models, suggesting that the three features are also important  
674 for preserving retina-based SF information (Figs. 5, 6). Max pooling over the wide window  
675 employed in the SNNs contributed to object-based SF representation of units sensitive to low  
676 SFs (Fig. 7). These findings advance our understanding of the computational processes utilized  
677 by the subcortical pathway in facial expression recognition.

678

679 ***Modest performance of the SNNs and neural computations of the subcortical pathway.***

680

681 Based on psychological assessment and brain imaging in V1-lesioned patients, it has been  
682 proposed that affective blindsight is mediated by components of the subcortical pathway spared  
683 by the lesions, including the superior colliculus, pulvinar, and amygdala (de Gelder et al., 1999;  
684 Pegna et al., 2005; Striemer et al., 2019). One view assumes that the shortest route directly  
685 connecting the three subcortical structures conveys facial expression information from the  
686 superior colliculus via the pulvinar to the amygdala (Tamietto and de Gelder, 2010). A different  
687 view proposes that information from the pulvinar reaches the amygdala through the facial  
688 processing system in the temporal cortex, under an assumption that “the direct connections of  
689 the pulvinar with the amygdala are likely insufficient in themselves for recognizing emotional  
690 expressions” (Gerbella et al., 2019). The present study demonstrates that the SNNs with only  
691 three processing stages and the subcortical physiological properties can successfully acquire an  
692 ability to discriminate facial expressions.

693

694 The average correct rate of classifying the seven facial expressions in the present study was 0.51.  
695 This rate was well above chance ( $1/7 = 0.14$ ), but was far from perfect. The modest correct rate  
696 is in line with the performance of patients with affective blindsight. Pegna et al. (2005) reported  
697 that a patient with bilateral lesions in V1 discriminated happy faces from either angry, sad, or  
698 horrified faces at correct rates of 0.58–0.62, marginally above the chance level of 0.5. Another  
699 patient with bilateral lesions in V1 exhibited correct rates of 0.64–0.67 for happy vs. fearful or  
700 angry faces (chance level = 0.5; Striemer et al. 2019). The residual ability of facial expression  
701 classification in these patients was only moderate compared to the nearly perfect performance



20230127

702 in healthy people. This raises the question of why subcortical processing supports vision more  
703 poorly than visual functions mediated by the cortical pathway.

704

705 A traditional explanation is that neurons in the subcortical pathway respond to low SFs and are  
706 less sensitive to high SFs than the cortical pathway (e.g., Vuilleumier et al., 2003; Méndez-  
707 Bértolo et al., 2016; Burra et al., 2019). This will limit the ability of the subcortical pathway to  
708 analyze the fine details of visual images, and can itself result in the inaccurate processing of face  
709 images. However, the dependence of the subcortical response on low SFs has been disputed by  
710 other researchers (De Cesarei and Codispoti, 2013; McFadyen et al., 2017). Our results suggest  
711 that low SF sensitivity, if important, was not the only cause, because the DoG filter models  
712 combined with narrow-pooling or add-layer modifications exhibited improved performances,  
713 despite the fact that our DoG filters were tuned to low SFs, and had full width at half maximum  
714 of 0.067–1.0 cycles/degree. Note that we estimated this value on an assumption of the image  
715 size of 30.5° based on our DoG parameters, the filter resolution, and the RF size of superior  
716 colliculus neurons representing the foveal region. The range of DoG-filter width corresponds to  
717 that applied in models of the superior colliculus in a recent simulation study (Méndez et al.,  
718 2022).

719

720 Another explanation is that the small number of processing stages in the cortical pathway  
721 hampers detailed analysis of visual inputs. However, a previous study (Dailey et al., 2002)  
722 showed that CNNs that had only two processing layers, with Gabor filters at the first stage,  
723 performed highly accurate discrimination of facial expressions (the mean correct rate for  
724 classifying six facial expressions was 0.90). The performance of our SNNs incorporating the  
725 three subcortical properties was not this high. This was not due to inadequate training, because  
726 the performance reached a plateau and stayed stable over a large number of iterations in the  
727 training sessions (Fig. 3). This was further verified by showing that the correct performance for  
728 the test set remained unchanged even after overly excessive training with 3,000,000 iterations.  
729 Furthermore, replacing not only the small number of processing layers but also the filter type at  
730 the first processing layer and the width of the pooling window with the corresponding cortical  
731 properties improved the performance of the SNNs (Fig. 4). The three properties at least partially  
732 underlie the less accurate processing of facial images in the subcortical pathway, and hence,  
733 may be responsible for the low performance in affective blindsight.

734

### 735 ***Confusions of facial expressions in the SNNs, DNNs and patients.***

736

737 The classification accuracy of the SNNs varied across facial expressions (Fig. 3A, B). The  
738 classification performance of the SNNs was best for happy and surprised faces and worst for  
739 sad and neutral faces. The rank order of performance on the seven facial expressions was largely  
740 consistent across the 20 SNNs trained independently from random states (Fig. 3C). It also  
741 corresponds to the classification performance by previously developed AlexNet-based DNNs  
742 (Inagaki et al., 2022b). These DNNs were trained to discriminate between the seven expressions  
743 derived either from the KDEF database or the Kokoro Research Center (KRC) facial expression  
744 database (Ueda et al., 2019). Like the SNNs, the DNNs exhibited the best performance for happy  
745 and surprised faces (KDEF: 0.93 for happy, 0.84 for surprised; KRC: 0.91 for happy, 0.84 for  
746 surprised; chance level, 0.14). This coincidence may simply suggest that within each database,

20230127

747 facial features are consistent across faces with happy or surprised expressions, but are more  
748 diverse across faces with sad or neutral expressions. However, the variations across examples  
749 of facial expressions within a database are not the sole reason for the difference in the  
750 performance across facial expressions, because neutral faces were classified poorly by the SNNs  
751 (Fig. 3B; correct rate = 0.34), but the DNNs of Inagaki et al. (2022b) classified them with high  
752 correct rates (0.85 for KDEF, 0.82 for KRC). An alternative, yet-to-be-tested explanation is that  
753 the ease (or difficulty) of classification may vary across the facial expressions owing to  
754 differences in the conspicuousness of component facial actions underlying various expressions.  
755 Similarities between neural networks regarding expression-specific performance may vary  
756 according to these differences.

757

758 The relatively poor ability to distinguish between sad and neutral faces was also observed in  
759 another CNN with the first layer of DoG filtering and average pooling (Méndez et al., 2022).  
760 This CNN was constructed to simulate facial processing in the superior colliculus, and was  
761 trained to discriminate three facial expressions: happy, sad, and neutral. The CNN showed the  
762 best performance for happy faces and moderate performance for sad faces, but classified neutral  
763 faces into neutral faces with a classification rate of 0.49 and into sad faces with a rate of 0.39.  
764 The fact that this CNN and the SNNs in the present study demonstrated this confusion, whereas  
765 AlexNet-based DNNs and our add-layer models (0.52 for sad, 0.67 for neutral) did not, suggests  
766 that the convolution processes after the initial DoG filtering (in the case of add-layer models) or  
767 the convolution by the Gabor filters (in the case of AlexNet-based DNNs) may be critical for  
768 classification of sad and neutral faces.

769

770 Finally, we point out that the expression-dependent performance of the SNNs also had both  
771 similarities and dissimilarities to that observed in a V1-lesioned patient. The patient reported by  
772 de Gelder et al. (1999) classified happy and sad faces with a higher correct rate than angry and  
773 fearful faces; our SNNs and this patient classified happy faces well, whereas the performance  
774 for sad faces was poor in the SNNs but good in the patient.

775

776 ***Reference frame of coding SF information and invariance of visual responses.***

777

778 FC1 units of the SNNs consisted of two major groups, each with different properties regarding  
779 SF processing (Fig. 5). One group of units responded best to the same object-based SFs  
780 (cycles/object) regardless of the stimulus size (peak shift around 0). This size-invariant response  
781 indicates that these units represent SFs in the object-based coordinate. Most of these units were  
782 tuned to a low SF range (around one to two cycles/object). In the other group of units, the  
783 optimal object-based SFs shifted when testing was performed with different stimulus sizes. The  
784 direction of the shift was consistent with the interpretation that the units were tuned to retina-  
785 based SFs (cycles/degree) (peak shift around 1). That is, for larger stimuli, the units responded  
786 to higher object-based SFs that corresponded to the same retina-based SFs. The DoG filters at  
787 the initial stage and the shallow architecture appear to be critical for preserving the SF  
788 representation based on the retina-based coordinate, because FC1 units with retina-based SF  
789 sensitivity were reduced in number when the first convolution layer were changed to Gabor  
790 filters or when the number of processing layers was increased (Fig. 6A, middle, bottom).

791

20230127

792 A major group of the object-based SF units in the SNNs were tuned to low SFs (Fig. 5B). This  
793 curious bias of the object-based units towards low SF sensitivity likely resulted from the wide  
794 max pooling process. Lowpass-filtered facial images contain only coarse structure such as solid  
795 blobs at eye or mouth positions. Positional information of these blobs is initially detected by  
796 DoG filters, and is encoded as response patterns across units in the convolution layer. These  
797 blobs appear in different positions and scales for images of different sizes, and thus the response  
798 patterns vary between different sizes. After the max pooling operation, however, response  
799 patterns would become more similar between different sizes, because this operation renders  
800 units in the pooling layer insensitive to slight changes in spatial arrangement of local features.  
801 Indeed, the dissimilarity index for lowpass-filtered images decreased after max pooling in our  
802 data (Fig. 7). This effect might result in object-based SF tuning (i.e., preferential responses  
803 invariant of image size to a particular range of object-based SFs) for lowpass-filtered images.  
804 Wider pooling window would enhance this effect at the expense of losing fine details of inputs.  
805 When the pooling window is narrower, this effect would be incomplete, and units with  
806 intermediate peak shift values would increase, as we found in the narrow-pooling models (Fig.  
807 6A, top).

808  
809 One may wonder why FC1 units of the SNNs maintained sensitivity to retina-based SFs, i.e.,  
810 size-dependent representation of SFs, despite the demand that we imposed on the SNNs to  
811 classify facial expressions regardless of the seven different face image sizes. One possibility is  
812 that the architecture of our SNNs cannot achieve sufficient object-based representation, and  
813 remains suboptimal for the required task even after the excessive training sessions. This may be  
814 a reason for the modest classification performance of the SNNs. Indeed, replacement of the  
815 subcortical processing properties with the cortical properties resulted in the representation  
816 becoming more object-based (Fig. 5B, C, Fig. 6) and improved the classification performance  
817 (Fig. 4). However, if object-based SF encoding was the only requirement for optimal  
818 performance under our training conditions, the models that showed object-based SF encoding  
819 should have had the highest correct rate, but this was not the case. The add-layer models and the  
820 narrow-pooling + add-layer models exhibited the best object-based encoding of SFs (Fig. 6A  
821 bottom, 6B bottom), while they performed worse than the full-replacement model (Fig. 4A).  
822 The representation acquired for the classification depended not only on the task demand of size-  
823 invariant classification of facial expressions, but also on other, yet unspecified, constraints  
824 deriving probably from the architecture of the models.

825  
826 In the primate amygdala, the responses of many neurons are affected by retina-based SFs, and  
827 only a minority of neurons have perfect object-based SF sensitivity (Inagaki and Fujita, 2011).  
828 By contrast, many FC1 units tuned to low SFs of the SNNs exhibited object-based SF sensitivity.  
829 The paucity of evidence for units with object-based SF sensitivity in the amygdala may be  
830 related to the fact that the previous electrophysiological study (Inagaki and Fujita, 2011) did not  
831 present face images with very low SFs, and may have overlooked the neurons with object-based  
832 SF sensitivity in this range of SFs.

833  
834 Some inferior temporal cortex neurons exhibit invariant responses to changes in shape sizes  
835 (Rolls and Baylis, 1986; Ito et al., 1995). The max pooling operation may help achieve these  
836 invariant responses. To some degree, max pooling ignores positional changes of inputs in each  
837 region of interest. Because size changes involve alternations in edge positions without

20230127

838 modifications in topologies, if the changes are small enough to be covered by each region of  
839 interest, stimuli before and after the changes would yield similar responses. The effects of wide  
840 pooling shown in Figure 7 suggest that some aspects of the invariant responses of inferior  
841 temporal cortex neurons can simply be achieved by bypassing early cortical areas with high  
842 spatial resolutions such as V1. Such shortcut routes indeed exist, including the projection from  
843 the pulvinar to V2 and then to the posterior inferior temporal cortex and the projection from the  
844 pulvinar to V4 and then to the anterior inferior temporal cortex (Pessoa and Adolphs, 2010).

845  
846 The size invariance in low SFs is important in newborns. They have blurred visions that relies  
847 on low SFs (Atkinson et al., 1974; Dobson & Teller, 1978), but respond to faces or face-like  
848 patterns irrespective of the stimulus size or the viewing distance (Cassia et al., 2001; De Heering  
849 et al., 2008). These findings indicate that the ability of size-invariant face recognition based on  
850 low SFs is innately implemented in our visual system. Convergence of inputs from the  
851 superficial layer of the superior colliculus to the deep layer, which is already present in newborns  
852 (Wallace et al., 1997), may be part of the neural substrate supporting this aspect of size-invariant  
853 face recognition.

854

### 855 ***Concluding remarks***

856

857 The present study provides the first computational model for facial expression processing along  
858 the subcortical pathway (see Méndez et al., 2022 for a model of face processing in the superior  
859 colliculus). Despite the celebrated success of DNNs in modeling visual processing in the ventral  
860 cortical pathway, it has remained unclear whether and how the CNN architecture can be adapted  
861 to processing in the subcortical pathway. We demonstrated that the SNNs implemented with the  
862 three computational properties of the subcortical pathway, i.e., a shallow layer architecture,  
863 concentric receptive fields at the first processing stage, and a greater degree of spatial pooling,  
864 were successfully trained to discriminate facial expressions with a modest correct rate. The three  
865 properties were all essential for reproducing the modest performance seen in V1-lesioned  
866 patients, as well as the representation of SFs in the retina-based coordinate observed in a  
867 population of amygdala neurons. Research interest in the role of subcortical structures in  
868 cognitive functions has recently surged, but physiological data are still much sparser for  
869 subcortical structures than for the cerebral cortex (Janacsek et al., 2022). Computational  
870 approaches such as the one we present here are expected to partially compensate for this data  
871 scarcity and to guide future research.

20230127

872 **References**

- 873 Ameijeiras-Alonso, J., Crujeiras, R.M. and Rodríguez-Casal, A. Mode testing, critical  
874 bandwidth and excess mass. *TEST*, **28**, 900-919 (2019). doi: 10.1007/s11749-018-0611-5
- 875 Ameijeiras-Alonso, J., Crujeiras, R.M. and Rodríguez-Casal, A. Multimode: An R package for  
876 mode assessment. *J. Stat. Softw.* **97**, Issue 9 (2021). doi: 10.18637/jss.v097.i09
- 877 Atkinson, J., Braddick, O. & Braddick, F. Acuity and contrast sensitivity of infant vision. *Nature*  
878 **247**, 403-404 (1974). doi: 10.1038/247403a0
- 879 Bender, D.B. Retinotopic organization of macaque pulvinar. *J. Neurophysiol.* **46**, 672-693  
880 (1981). doi: 10.1152/jn.1981.46.3.672
- 881 Bradski, G. The OpenCV Library. *Dr. Dobb's Journal of Software Tools* (2000).  
882 <https://opencv.org/> (version, 2.4.8; this version is no longer available)
- 883 Buiatti, M., Di Giorgio, E., Piazza, M., Polloni, C., Menna, G., Taddei, F., Baldo, E. &  
884 Vallortigara, G. Cortical route for facelike pattern processing in human newborns. *Proc.*  
885 *Natl. Acad. Sci. U.S.A.* **116**, 4625–4630 (2019). doi: 0.1073/pnas.181241911  
886
- 887 Burra, N., Hervais-Adelman, A., Celeghin, A., De Gelder, B. & Pegna, A.J. Affective blindsight  
888 relies on low spatial frequencies. *Neuropsychol.* **128**, 44-49 (2019).  
889 doi: 10.1016/j.neuropsychologia.2017.10.009
- 890 Cassia, V.M, Simion, F. & Umiltà, C. Face preference at birth: the role of an orienting  
891 mechanism. *Develop. Sci.* **4**, 101-108 (2001). doi: 10.1111/1467-7687.00154
- 892 Chen, C-Y., Hoffmann, K.-P., Distler, C., & Hafed, Z. M. The foveal visual representation of  
893 the primate superior colliculus. *Curr. Biol.* **29**, 2109-2119 (2019). doi:  
894 10.1016/j.cub.2019.05.040
- 895 Chen, C-Y., Sonnenberg, L., Weller, S., Witschel, T. & Hafed, Z.M. Spatial frequency  
896 sensitivity in macaque midbrain. *Nat. Comm.*, **9**, pp.1-13 (2018).  
897 doi: 10.1038/s41467-018-05302-5  
898
- 899 Churan, J., Guitton, D. & Pack, C.C. Spatiotemporal structure of visual receptive fields in  
900 macaque superior colliculus. *J. Neurophysiol.* **108**, 2653-2667 (2012).  
901 doi: 10.1152/jn.00389.2012
- 902 Connor, C.E., Brincat, S.L. & Pasupathy, A. Transformation of shape information in the ventral  
903 pathway. *Curr. Opin. Neurobiol.* **17**, 140-147 (2007). doi: 10.1016/j.conb.2007.03.002

20230127

- 904 Conway, B.R., Chatterjee, S., Field, G.D., Horwitz, G.D., Johnson, E.N., Koida, K. & Mancuso,  
905 K. Advances in color science: from retina to behavior. *J. Neurosci.* **30**, 14955-14963 (2010).  
906 doi: 10.1523/JNEUROSCI.4348-10.2010
- 907 Cynader, M. & Berman, N. Receptive-field organization of monkey superior colliculus. *J.*  
908 *Neurophysiol.* **35**, 187-201 (1972). doi: 10.1152/jn.1972.35.2.187
- 909 Dailey, N.M., Cottrell, W.G., Padgett, C. & Adolphs, R. EMPATH: A neural network that  
910 categorizes facial expressions. *J. Cogni. Neurosci.* **14**, 1158-1173 (2002).  
911 doi: 10.1162/089892902760807177
- 912 De Cesarei, A. & Codispoti, M. Spatial frequencies and emotional perception. *Rev. Neurosci.*  
913 **24**, 89-104 (2013). doi: 10.1515/revneuro-2012-0053
- 914 de Gelder, B., Vroomen, J., Pourtois, G. & Weiskrantz, L. Non-conscious recognition of affect  
915 in the absence of striate cortex. *NeuroRep.* **10**, 3759-3763 (1999).  
916 doi: 10.1097/00001756-199912160-00007  
917
- 918 De Heering, A., Turati, C., Rossion, B., Bulf, H., Goffaux, V. & Simion, F. Newborns' face  
919 recognition is based on spatial frequencies below 0.5 cycles per degree. *Cognition*, 106(1),  
920 pp.444-454 (2008). doi: 10.1016/j.cognition.2006.12.012  
921
- 922 Desimone, R., Albright, T.D., Gross, C.G. & Bruce, C. Stimulus-selective properties of inferior  
923 temporal neurons in the macaque. *J. Neurosci.* **4**, 2051-2062 (1984).  
924 doi: 10.1523/JNEUROSCI.04-08-02051.1984  
925
- 926 Dobson, V. & Teller, D.Y. Visual acuity in human infants: a review and comparison of  
927 behavioral and electrophysiological studies. *Vis. Res.* **18**, 1469-1483 (1978).  
928 doi: 10.1016/0042-6989(78)90001-9
- 929 Duchaine, B. & Yovel, G. A revised neural framework for face processing. *Annu. Rev. Vis. Sci.*  
930 **1**, 393-416 (2015). doi: 10.1146/annurev-vision-082114-035518
- 931 Freeman, J. & Simoncelli, E.P. Metamers of the ventral stream. *Nat. Neurosci.* **14**, 1195-1201  
932 (2011). doi: 10.1038/nn.2889
- 933 Freiwald, W., Duchaine, B. & Yovel, G. Face processing systems: from neurons to real-world  
934 social perception. *Annu. Rev. Neurosci.* **39**, 325-346 (2016).  
935 doi: 10.1146/annurev-neuro-070815-013934  
936
- 937 Fujita, I., Tanaka, K., Ito, M. & Cheng, K. Columns for visual features of objects in monkey  
938 inferotemporal cortex. *Nature* **360**, 343-346 (1992).  
939 doi: 10.1038/360343a0

20230127

- 940 Gerbella, M., Caruana, F. & Rizzolatti, G. Pathways for smiling, disgust and fear recognition in  
941 blindsight patients. *Neuropsychol.* **128**, 6-13 (2019).  
942 doi: 10.1016/j.neuropsychologia.2017.08.028
- 943 Goldberg, M. E. & Wurtz, R. H. Activity of superior colliculus in behaving monkey. I. Visual  
944 receptive fields of single neurons. *J. Neurophysiol.* **35**, 542–559 (1972).  
945 doi: 10.1152/jn.1972.35.4.542
- 946 Goodfellow, I., Bengio, Y. & Courville, A. *Deep Learning*. Cambridge: MIT Press, pp. 71–73  
947 (2016). <http://www.deeplearningbook.org> (doi is not available)
- 948 Güçlü, U. & van Gerven, M.A.J. Deep neural networks reveal a gradient in the complexity of  
949 neural representations across the ventral stream. *J. Neurosci.* **35**, 10005-10014 (2015).  
950 doi.org/10.1523/JNEUROSCI.5023-14.2015
- 951 Hassabis, D., Kumaran, D., Summerfield, C. & Botvinick, M. Neuroscience-inspired artificial  
952 intelligence. *Neuron.* **95**, 245-258 (2017). doi: 10.1016/j.neuron.2017.06.011
- 953 Haxby, J.V, Hoffman, E.A. & Gobbini, M.I. The distributed human neural system for face  
954 perception. *Trends Neurosci.* **4**, 223-233 (2000). doi: 10.1016/s1364-6613(00)01482-0
- 955 He, K., Zhang, X., Ren, S. and Sun, J. Delving deep into rectifiers: Surpassing human-level  
956 performance on ImageNet classification. *Proc. IEEE Int. Conf. Comput. Vis. pp.* 1026-1034  
957 (2015). doi: 10.1109/ICCV.2015.123
- 958 He, K., Zhang, X., Ren, S. & Sun, J. Deep residual learning for image recognition. *Proc. IEEE*  
959 *Comput. Soc. Conf. Comput. Vis. Pattern Recognit.* 770-778 (2016).  
960 doi: 10.1109/CVPR.2016.90
- 961 Inagaki, M. & Fujita, I. Reference frames for spatial frequency in face representation differ in  
962 the temporal visual cortex and amygdala. *J. Neurosci.* **31**, 10371-1039 (2011).  
963 doi: 10.1523/JNEUROSCI.1114-11.2011
- 964 Inagaki, M., Inoue, K.-I., Tanabe, S., Kimura, K., Takada, M. & Fujita, I. Rapid processing of  
965 threatening faces in the amygdala of nonhuman primates: subcortical inputs and dual roles.  
966 *Cereb. Cortex.* 10.1093/cercor/bhac109 (2022a). doi: 10.1093/cercor/bhac109
- 967 Inagaki, M., Ito, T., Shinozaki, T. & Fujita, I. Convolutional neural networks reveal differences  
968 in action units of facial expressions between face image databases developed in different  
969 countries. *Front. Psychol.* doi: 10.3389/fpsyg.2022.988302 (2022b).
- 970 Ito, M., Tamura, H., Fujita, I. & Tanaka, K. Size and position invariance of neuronal responses  
971 in monkey inferotemporal cortex. *J. Neurophysiol.* **73**, 218-226 (1995).  
972 doi: 10.1152/jn.1995.73.1.218

20230127

- 973 Janacsek, K., Evans, T.M., Kiss, M. Shah, L., Blumenfeld, H. & Ullman, M.T. Subcortical  
974 cognition: the fruit below the rind. *Annu. Rev. Neurosci.* **45**, 361-386 (2022).  
975 doi: 10.1146/annurev-neuro-110920-013544
- 976 Johnson, M.H. Subcortical face processing. *Nat. Rev. Neurosci.* **6**, 766-774 (2005).  
977 doi: 10.1038/nrn1766
- 978 Jones, J.P. & Palmer, L.A. An evaluation of the two-dimensional Gabor filter model of simple  
979 receptive fields in cat striate cortex. *J. Neurophysiol.* **58**, 1233–1258 (1987).  
980 doi: 10.1152/jn.1987.58.6.1233
- 981 Komatsu, H. & Goda, N. Neural mechanisms of material perception: Quest on Shitsukan.  
982 *Neurosci.* **392**, 329-347 (2018). doi: 10.1016/j.neuroscience.2018.09.001
- 983 Kravitz, D.J., Saleem, K.S., Baker, C.I., Ungerleider, L.G. & Mishkin, M. The ventral visual  
984 pathway: an expanded neural framework for the processing of object quality. *Trends Cogn.*  
985 *Sci.* **17**, 26-49 (2013). doi: 10.1016/j.tics.2012.10.011
- 986 Krizhevsky, A., Sutskever, I. & Hinton, G.E. ImageNet classification with deep convolutional  
987 neural networks. *Adv. Neural Inf. Proces. Syst (NeurIPS)*. **25**, 1097-1105 (2012).  
988 [https://proceedings.neurips.cc/paper/2012/file/c399862d3b9d6b76c8436e924a68c45b-](https://proceedings.neurips.cc/paper/2012/file/c399862d3b9d6b76c8436e924a68c45b-Paper.pdf)  
989 [Paper.pdf](https://proceedings.neurips.cc/paper/2012/file/c399862d3b9d6b76c8436e924a68c45b-Paper.pdf) (doi is not available)
- 990 Langner, O., Dotsch, R., Bijlstra, G., Wigboldus, D.H., Hawk, S.T. & Van Knippenberg, A.D.  
991 Presentation and validation of the Radboud Faces Database. *Cogni. Emot.* **24**, 1377-1388  
992 (2010). doi: 10.1080/02699930903485076
- 993 LeCun, Y., Bengio, Y. & Hinton, G. Deep learning. *Nature.* **521**, 436-444 (2015).  
994 doi: 10.1038/nature14539
- 995 Lundqvist, D., Flykt, A. & Öhman, A. The Karolinska Directed Emotional Faces - KDEF.  
996 Stockholm: Department of Clinical Neuroscience, Psychology section, Karolinska Institute.  
997 CD-ROM (1998).
- 998 Marino, R.A., Rodgers, C.K., Levy, R. & Munoz, D.P. Spatial relationships of visuomotor  
999 transformations in the superior colliculus map. *J. Neurophysiol.* **100**, 2564-2576 (2008).  
1000 doi: 10.1152/jn.90688.2008
- 1001 Méndez, C.A., Celeghin, A., Diano, M., Orsenigo, D., Ocak, B. & Tamietto, M. A deep neural  
1002 network model of the primate superior colliculus for emotion recognition. *Phil. Trans. R.*  
1003 *Soc. B* **377**, 20210512 (2022). doi: 10.1098/rstb.2021.0512
- 1004 Méndez-Bértolo, C., Moratti, S., Toledano, R., Lopez-Sosa, F., Martinez-Alvarez, R., Mah,  
1005 Y.H., Vuilleumier, P., Gil-Nagel, A. & Strange, B.A. A fast pathway for fear in human  
1006 amygdala. *Nat. Neurosci.*, **19**, 1041-1049 (2016). doi: 10.1038/nn.4324



20230127

- 1007 Morawetz, C., Baudewig, J., Treue, S. & Dechent, P. Diverting attention suppresses human  
1008 amygdala responses to faces. *Front. Human Neurosci.* **4**, 226 (2010).  
1009 doi: 10.3389/fnhum.2010.00226
- 1010 Morris, J.S., DeGelder, B., Weiskrantz, L. & Dolan, R.J. Differential extrageniculostriate and  
1011 amygdala responses to presentation of emotional faces in a cortically blind field. *Brain.*  
1012 **126**, 1241-1252 (2001). doi: 10.1093/brain/124.6.1241
- 1013 Morris, J.S., Öhman, A. & Dolan, R.J. A subcortical pathway to the right amygdala mediating  
1014 “unseen” fear. *Proc. Natl. Acad. Sci. U.S.A.* **96**, 1680-1685 (1999).  
1015 doi: 10.1073/pnas.96.4.1680
- 1016 Movellan, J.R. Tutorial on Gabor filters. *Open Source Document.* **40**, 1–23 (2002).  
1017 <https://inc.ucsd.edu/mplab/75/media//gabor.pdf> (access, 2022/12/31; doi is not available)
- 1018 Nakano, T., Higashida, N. & Kitazawa, S. Facilitation of face recognition through the retino-  
1019 tectal pathway. *Neuropsychol.* **51**, 2043-2049 (2013).  
1020 doi: 10.1016/j.neuropsychologia.2013.06.018
- 1021 Nguyen, M.N., Hori, E., Matsumoto, J., Tran, H.T., Ono, T. & Nishijo, H. Neuronal responses  
1022 to face-like stimuli in the monkey pulvinar. *Eur J Neurosci.* **37**, 35-51, (2013).  
1023 doi: 10.1111/ejn.12020
- 1024 Nguyen, M.N., Matsumoto, J., Hori, E., Maior, R.S., Tomaz, C., Tran, A.H., Ono, T., & Nishijo,  
1025 H. Neuronal responses to face-like and facial stimuli in the monkey superior colliculus.  
1026 *Front. Behav. Neurosci.* **8**, 85 (2014). doi: 10.3389/fnbeh.2014.00085
- 1027 Pegna, A.J., Khateb, A., Lazeyras, F., & Seghier, M.L. Discriminating emotional faces without  
1028 primary visual cortices involves the right amygdala. *Nat. Neurosci.* **8**, 24-25 (2005).  
1029 doi: 10.1038/nn1364
- 1030 Perret, D.I., Mistlin, A.J., & Chitty, A.J. Visual neurones responsive to faces. *Trends Neurosci.*  
1031 **10**, 358-364 (1987). doi: 10.1016/0166-2236(87)90071-3
- 1032 Pessoa, L. & Adolphs, R. Emotion processing and the amygdala: from a “low road” to “many  
1033 roads” of evaluating biological significance. *Nat. Rev. Neurosci.* **11**, 773-782 (2010).  
1034 doi: 10.1038/nrn2920
- 1035 Petry, H.H. & Bickford, M.E. The second visual system of the tree shrew. *J. Comp. Neurol.* **527**,  
1036 679-693 (2019). doi: 10.1002/cne.24413
- 1037 Rai, M. & Rivas, P. A review of convolutional neural networks and Gabor filters in object  
1038 recognition. *2020 Int. Conf. Comput. Sci. Comput. Intelligence (CSCI)* 1560-1567 (2020).  
1039 doi: 10.1109/CSCI51800.2020.00289

20230127

- 1040 Roe, A.W., Chelazzi, L., Connor, C.E., Conway, B.R., Fujita, I., Gallant, J.L., Lu, H., &  
1041 Vanduffel, W. Toward a unified theory of visual area V4. *Neuron* **74**, 12-29 (2012).  
1042 doi: 10.1016/j.neuron.2012.03.011
- 1043 Rolls, E.T. & Baylis, G.C. Size and contrast have only small effects on the responses to faces  
1044 of neurons in the cortex of the superior temporal sulcus of the monkey. *Exp. Brain Res.*  
1045 **65**, 38-48 (1986). doi: 10.1007/BF00243828
- 1046 Schmolesky, M.T., Wang, Y., Hanes, D.P., Thompson, K.G., Leutgeb, S., Schall, J.D. &  
1047 Leventhal, A.G. Signal timing across the macaque visual system. *J. Neurophysiol.* **79**,  
1048 3272-3278 (1998). doi: 10.1152/jn.1998.79.6.3272
- 1049 Simonyan, K. & Zisserman, A. Very deep convolutional networks for large-scale image  
1050 recognition. *arXiv* 1409.1556 (2014). doi: 10.48550/arXiv.1409.1556
- 1051 Striemer, C.L., Whitwell, R.L. & Goodale, M.A. Affective blindness in the absence of input  
1052 from face processing regions in occipital-temporal cortex. *Neuropsychol.* **128**, 50-57  
1053 (2019). doi: 10.1016/j.neuropsychologia.2017.11.014
- 1054 Tamietto, H., Castelli, L., Vighetti, S., Perozzo, P., Geminiani, G., Weiskrantz, L. & de Gelder,  
1055 B. Unseen facial and bodily expressions trigger fast emotional reactions. *Proc. Natl. Acad.*  
1056 *Sci. U.S.A.* **106**, 17661-17666 (2009). doi: 10.1073/pnas.0908994106
- 1057 Tamietto, H. & de Gelder, B. Neural bases of the non-conscious perception of emotional signals.  
1058 *Nat. Rev. Neurosci.* **11**, 697-709 (2010). doi: 10.1038/nrn2889
- 1059 Tokui, S., Oono, K., Hido, S. & Clayton, J. Chainer: a next-generation open source framework  
1060 for deep learning. *Proc. Workshop on machine learning systems (LearningSys) in 29th*  
1061 *annual conference on neural information processing systems.* **5**, 1-6 (2015).  
1062 [http://learningsys.org/papers/LearningSys\\_2015\\_paper\\_33.pdf](http://learningsys.org/papers/LearningSys_2015_paper_33.pdf).  
1063 <https://github.com/chainer/chainer/releases/tag/v3.0.0> (version, 3.0.0; release, Oct 17,  
1064 2017) (doi is not available)
- 1065 Tsao, D.Y. & Livingstone, M.S. Mechanisms of face perception. *Annu. Rev. Neurosci.* **31**, 411-  
1066 437 (2008). doi: 10.1146/annurev.neuro.30.051606.094238
- 1067 Ueda, Y., Nunoi, M., & Yoshikawa, S. Development and validation of the Kokoro Research  
1068 Center (KRC) facial expression database. *Psychologia* **61**, 221-240 (2019).  
1069 doi: 10.2117/psysoc.2019-A009
- 1070 Ungerleider, L.G. & Mishkin, M. Two cortical visual systems. In Ingle, DJ, Goodale, MA &  
1071 Mansfield, RJW, editors. *Analysis of Visual Behavior*. Cambridge: MIT Press, pp. 549-  
1072 586 (1982). <https://www.cns.nyu.edu/~tony/vns/readings/ungerleider-mishkin-1982.pdf>  
1073 (doi is not available)

20230127

- 1074 Updyke, B.V. Characteristics of unit responses in superior colliculus of the Cebus monkey. *J.*  
1075 *Neurophysiol.* **37**, 896-909 (1974). doi: 10.1152/jn.1974.37.5.896
- 1076 Van den Bergh, G., Zhang, B., Arckens, L. & Chino, Y.M. Receptive-field properties of V1 and  
1077 V2 neurons in mice and macaque monkeys. *J. Comp. Neurol.* **518**, 2051-2070 (2010). doi:  
1078 10.1002/cne.22321
- 1079 Verhoef, B-E., Vogels, R. & Janssen, P. Binocular depth processing in the ventral visual  
1080 pathway. *Phil. Trans. R. Soc. B* **371**, 20150259 (2016). doi: 10.1098/rstb.2015.0259
- 1081 Vuilleumier, P., Armony, J.L., Driver, J. & Dolan, R.J. Distinct spatial frequency sensitivities  
1082 for processing faces and emotional expressions. *Nat. Neurosci.* **6**, 624–631 (2003).  
1083 doi: 10.1038/nm1057
- 1084 Wallace, M. T., McHaffie, J. G. & Stein, B. E. Visual response properties and visuotopic  
1085 representation in the newborn monkey superior colliculus. *J. Neurophysiol.* **78**, 2732–2741  
1086 (1997). doi: 10.1152/jn.1997.78.5.2732
- 1087 Yamins, D.L.K. & DiCarlo, J.J. Using goal-driven learning models to understand sensory cortex.  
1088 *Nat. Neurosci.* **19**, 356-365 (2016). doi: 10.1038/nm.4244
- 1089 Yamins, D.L.K., Hong, H., Cadieu, C.F., Solomon, E.A., Seibert, D. & DiCarlo, J.J.  
1090 Performance-optimized hierarchical models predict neural responses in higher visual  
1091 cortex. *Proc. Natl. Acad. Sci. U.S.A.* **111**, 8619-8624 (2014).  
1092 doi: 10.1073/pnas.1403112111

20230127

1093 **Acknowledgments**

1094 This work was supported by grants from the Ministry of Education, Culture, Sports, Science  
1095 and Technology of Japan (JP17H01381 and JP21H02596 to IF; JP18H04197, JP20H04578, and  
1096 20K12023 to MI); the Center for Information and Neural Networks; the Ministry of Internal  
1097 Affairs and Communications of Japan. CL was supported by the Research Fellowship for Young  
1098 Scientists from the Japan Society for the Promotion of Science.

1099 **Author contributions**

1100 CL, MI, TS, and IF designed the research; CL, MI, and TS performed the research; CL, MI, TS,  
1101 and IF wrote the paper. All authors approved the submitted version.

1102 **Competing interests**

1103 The authors declare no conflicts of interest.

1104 **Data availability**

1105 All data and analysis codes are available from the corresponding author upon request.

1106 **Ethics statement**

1107 Written informed consent was obtained for the publication of any identifiable images included  
1108 in this article.