

1 Perceptual Expertise and Attention: An
2 Exploration using Deep Neural Networks

3 Soukhin Das^{1,2}, G.R. Mangun^{1,2,4}, Mingzhou Ding³

4

5

6

7

1 Center for Mind and Brain, University of California, Davis

8

2 Department of Psychology, University of California, Davis

9

3 Pruitt Family Department of Biomedical Engineering, University of Florida

10

4 Department of Neurology, University of California, Davis

11

12

13 Abstract

14

15 Perceptual expertise and attention are two important factors that enable superior object
16 recognition and task performance. While expertise enhances knowledge and provides a holistic
17 understanding of the environment, attention allows us to selectively focus on task-related
18 information and suppress distraction. It has been suggested that attention operates differently
19 in experts and in novices, but much remains unknown. This study investigates the relationship
20 between perceptual expertise and attention using convolutional neural networks (CNNs), which
21 are shown to be good models of primate visual pathways. Two CNN models were trained to
22 become experts in either face or scene recognition, and the effect of attention on performance
23 was evaluated in tasks involving complex stimuli, such as superimposed images containing
24 superimposed faces and scenes. The goal was to explore how feature-based attention (FBA)
25 influences recognition within and outside the domain of expertise of the models. We found that
26 each model performed better in its area of expertise—and that FBA further enhanced task
27 performance, but only within the domain of expertise, increasing performance by up to 35% in
28 scene recognition, and 15% in face recognition. However, attention had reduced or negative
29 effects when applied outside the models' expertise domain. Neural unit-level analysis revealed
30 that expertise led to stronger tuning towards category-specific features and sharper tuning
31 curves, as reflected in greater representational dissimilarity between targets and distractors,
32 which, in line with the biased competition model of attention, leads to enhanced performance
33 by reducing competition. These findings highlight the critical role of neural tuning at single as
34 well as network level neural in distinguishing the effects of attention in experts and in novices

35 and demonstrate that CNNs can be used fruitfully as computational models for addressing
36 neuroscience questions not practical with the empirical methods.

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

52

53

54

55

56

57

58

59 Introduction

60

61 Convolutional neural networks (CNNs) are a class of deep neural networks that draw strong
62 structural parallels with the primate visual pathway (1-5). CNNs' functional relevance for
63 neuroscience has also be demonstrated in recent studies that compared single neuron activity in
64 monkeys and fMRI activities in humans with activity in CNNs, showing that there is close
65 correspondence between layers of CNNs and the areas within the visual hierarchy (6, 7) (8, 9).
66 Increasingly, CNNs are being used as models of primate visual processing, making possible
67 explorations that are not practical in biological systems (e.g., lesion), generating results that
68 inspire new questions and new empirical experimentation (10, 11) (12, 13) (14-17). The
69 introduction of feature-based attention (FBA) in CNNs has further deepened the integration
70 between AI-inspired neural models such as CNNs and cognitive neuroscience (1). It has been
71 shown that object recognition in challenging settings (e.g., images where scenes and faces are
72 superimposed) is enhanced with feature-based attention (FBA). Attention can operate both at
73 the level of elementary features and at the object level, which are essentially higher-order
74 collections of features. Attention may select objects based on the presence of relevant features,
75 implying a dual role for attention mechanisms. For example, when attention is directed toward
76 specific features, it enhances recognition of objects that are composed of these features. The
77 goal of this study is to leverage these developments to examine the relation between attention
78 and perceptual expertise computationally.

79

80 Perceptual expertise refers to the enhanced ability to recognize and categorize objects in a
81 specific category and can be acquired through extensive experience and practice. It changes how
82 objects in the category are perceived (18, 19) and represented in the visual cortex (20-22). For
83 example, expertise in face recognition is associated with enhanced activity in the face selective
84 area known as the fusiform face area (FFA) (23-25). Similarly, expertise in other object categories,
85 such as birds or cars, has been shown to increase neural activity in different category-specific
86 regions of the visual cortex, demonstrating the broad impact of experience and practice on neural
87 processing (26, 27). More relevant for this study, expertise also enables experts to more easily
88 attend to the salient features of objects falling within their area of expertise (28), suggesting a
89 possible relation between expertise and selective attention. For example, during car viewing it
90 has been observed that manipulating attention to the identity versus the location of cars had a
91 more pronounced impact on car novices compared to experts (27). Such effects transcend
92 category domains and have been found in such diverse domains as planes, animals, chessboards,
93 and radiography (26, 29-32). In all instances, experts demonstrate automatic holistic processing
94 during object recognition, outperforming novices who are often influenced by task constraints,
95 context, and various other factors (18, 30, 33, 34). We hypothesize that enhanced effectiveness
96 of selective attention in the domain of expertise is a key factor underlying the superior
97 performance of experts in their domains of expertise.

98

99 Depending on the datasets and the training objectives, CNNs can become experts in recognizing
100 different categories of images. A recent study showed that a CNN trained on the ImageNet
101 dataset to recognize objects became proficient at object recognition but less so on face

102 recognition, whereas a CNN trained on the VGGFace dataset to recognize faces become
103 proficient at face recognition, but less so on object recognition (35). In another study (9), a CNN
104 trained to recognize both scenes and objects evolved category-selective topographical units,
105 providing a computational account for the altered neural activity in category-selective brain
106 regions observed in experts relative to novices. What has not been addressed in these previous
107 studies is the relation between perceptual expertise and selective attention. We address this
108 question by considering two classes of CNNs: one that is trained to recognize scenes (Scene-
109 expert) and another that is trained to recognize faces (Face-expert). In our study, attention was
110 applied at both the feature and object level by biasing neural units that preferred certain
111 categories (faces or scenes). This approach allowed us to explore how attention modulates
112 recognition within and outside the domain of perceptual expertise. The aim was to assess the
113 effectiveness of attention-to-objects and attention-to-faces in each of the two CNN models. We
114 hypothesized that (1) the scene-trained CNN would be inferior at face recognition relative to
115 scene recognition, whereas the face-trained CNN would be inferior at scene recognition relative
116 to face recognition, and (2) that in challenging perceptual situations (e.g., superimposed stimuli),
117 the scene-trained CNN would benefit more from feature-based attention (FBA) for scene
118 recognition, whereas face-trained CNNs would benefit more from FBA for face recognition. The
119 geometry, representation, and activity of population-level neural activity were examined to
120 explore the underlying neural mechanisms.

121

122 Results

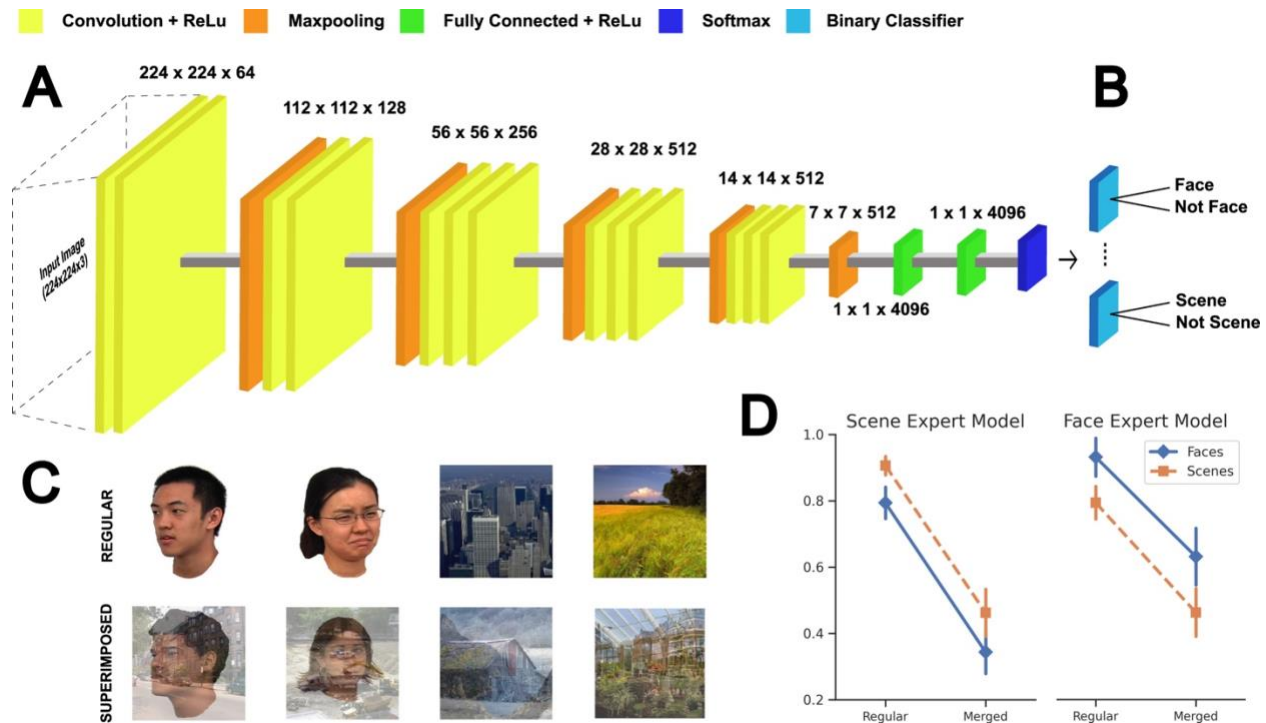
123

124 Overview

125

126 We used VGG16, a 16-layer deep convolutional neural network (DNN) (36), as a model of the
127 ventral visual stream; see Figure 1A. In this model, neurons in each layer are connected to
128 neurons in the next layer in a convolutional manner, mimicking the receptive-field based
129 feedforward retinotopic processing of visual information in the primate ventral visual system.
130 Two VGG16 models were independently trained to acquire expertise in either object recognition
131 or face recognition. Specifically, the scene network was trained on the ImageNet database (37)
132 consisting of 3.2M images where the training objective was to classify the images into 1000 object
133 categories as accurately as possible. On the other hand, the face network was trained on the
134 VGGFace dataset (38) consisting of 2.6M images of human faces where the training objective was
135 to classify the images into 2622 distinct faces as accurately as possible (Figure 1A). We will refer
136 to the network trained on ImageNet as ‘Scene-expert’ and the network pretrained on VGGFace
137 as ‘Face-expert’. The main purpose of this study was to examine the effectiveness of feature-
138 based attention (FBA) to faces and to scenes in the two networks when they are engaged in
139 performing challenging scene recognition and face recognition tasks.

140



141

142 *Figure 1. Model and study design. (A) The VGG16 convolutional neural network model (feature units and dimensions are labelled*
 143 *across the respective layers). One model is pretrained on ImageNet database (Scene-expert) and the other on VGGFace database*
 144 *(Face-expert). (B) The final layer of each network was replaced with a series of binary classifiers (logistic regression, one for each*
 145 *category) which were trained based on the datasets used in this study. (C) Regular and superimposed images from each category,*
 146 *sized at 224x224 pixels (input dimensions of VGG16). Superimposed images (bottom) were composed by transparently*
 147 *superimposing two images from either the same or the different categories. The regular images are used for training the binary*
 148 *classifiers which were then tested on both regular and superimposed images to identify the presence or absence of a certain*
 149 *category. (D) 5-fold cross-validation performance of the two models, image-category wise, for the Face-expert (right) and Scene-*
 150 *expert (left) models. Images were used from publicly available datasets (39-41).*

151 Experimental Paradigms and Model Performance

152

153 The recognition task consisted of identifying whether a particular object feature was present in
 154 an input image. For example, during the face recognition task, the models were asked to
 155 determine the presence/absence of a face in the input. For the model to perform these binary

156 classification tasks, the final classification layer of the VGG16 network was replaced with a series
157 of binary classifiers, each specific for recognizing the presence or absence of a particular object
158 category (Figure 1B). There were two types of input: regular images and superimposed images.
159 Regular images consisted of images belonging to one of the two categories without any
160 distractors: faces (male and female) and scenes (manmade and natural) (Figure 1C). After
161 training, task performance was obtained for the testing data (for more details on the training and
162 testing dataset, see Methods). See Figure 1D. As expected, the Scene-expert model had a higher
163 scene recognition accuracy (97.6%) for scene images compared to the Face-expert model (79.6%;
164 $p < 1e-5$, paired t-test across classification folds), while the Face-expert model performed better
165 (94%) in recognizing faces in regular images than the Scene-expert model (80%; $p < 1e-4$, one-
166 sided paired t-test across classification folds). These performance metrics were summarized in
167 Table 1. Based on these results, the two models can be said to have developed category-specific
168 expertise, with the Scene-Expert model excelling at recognizing scenes and the Face-expert
169 model excelling at recognizing faces.

170

	Scene-Expert Model		Face-Expert Model	
Image type	Regular	Superimposed	Regular	Superimposed
Faces	80%	36.2%	94%	64.4%
Scenes	97.6%	48.4%	79.6%	48%

171

172

Table 1 Performance of the two models for different types of images, regular and superimposed, across categories.

173 For the more challenging tasks, superimposed images were utilized as stimuli; see Figure 1C.
174 There were three types of superimposed images: (1) face over face, (2) scene over scene, and (3)
175 face over scene (which is equivalent to scene over face). For detecting the presence or absence
176 of a face, (1) and (3) were true positives whereas (2) were true negatives. For detecting the
177 presence or absence of a scene, (2) and (3) were true positives whereas (1) were true negatives.
178 The test images were balanced with 50% true positives and 50% true negatives; a total of 40
179 positive images and 40 negative images were used for testing. As expected, on superimposed
180 images, the performance was significantly decreased for both models. The Scene-expert model
181 exhibited a comparable scene recognition accuracy (48.4%) to the Face-expert model (48%; $p >$
182 0.05 , one-sided paired t-test across classification folds). However, the Face-expert model
183 performed significantly better (64.4%) in recognizing faces in superimposed images than the
184 Scene-expert model (36.2%; $p < 1e-3$, one-sided paired t-test across classification folds). A
185 summary of the performance metrics can be found in Table 1. It can be noted that both models
186 experienced a large decrement in performance when recognizing objects in superimposed
187 images. Thus, the presence of distractors made the task very challenging and reduced
188 performance of each expert model, thereby presenting a fruitful opportunity to test where
189 attentional mechanisms can be brought to bear to overcome these challenges.

190

191 [Attention Modulation of Model Performance](#)

192

193 Feature-based attention (FBA) was implemented according to the Feature Similarity Gain Model
194 (FSGM) (42). This model posits that neural activity is modulated in proportion to how strongly a

195 neuron prefers an attended feature (43, 44). When a stimulus falls within the receptive field of a
196 neuron, directing attention to that stimulus results in an increase in neural response, and this
197 increase occurs in a proportional manner, encompassing both preferred and nonpreferred
198 stimuli. From the foregoing, to apply FBA according to FSGM, each neuron's selectivity for
199 different experimental stimuli needs to be determined. This was done by calculating the
200 responses of each feature map to the two types of stimuli: faces and scenes (see Methods). Some
201 examples of the resulting tuning curves are shown in Figure 2. To implement FBA, the activity of
202 a certain neuron was modulated by scaling the slope of its activation function (ReLU) in the
203 network based on its tuning curve (i.e., how strongly a certain neural unit prefers a certain image
204 category). Units that are selective to the attended feature (target category) had their output
205 tuned up, while neurons that are not selective had their output tuned down (distractor category);
206 see Figure 3. We tested the effect of FBA on one layer of each model at a time, while the
207 performance across categories and layers that received attention modulation was recorded
208 individually. Figure 4 and Table 2 shows the change in performances across categories for both
209 model variants when tested with superimposed images. Overall, FBA yielded a favorable
210 influence on model performance. Importantly, an interplay of attention's effect across categories
211 and models was evident upon comparative analysis. Specifically, for the Scene-expert model,
212 attention yielded a more pronounced enhancement in performance for scenes (improvement
213 ranging from 15 - 35% depending on the layer of the model, Figure 4A) as opposed to faces
214 (improvement ranging from 3- 15% depending on the layer of the model). In contrast, the Face-
215 expert model exhibited significant recognition improvement on face categories (15 - 20%
216 depending on the layer, Figure 4B) but inconsistent even negative improvement on scene

217 categories (-12 to 10% depending on the layer). We performed a two-way analysis of variance
218 (ANOVA) including the factors model (Scene-expert and Face-expert), layer (1 to 13) and image
219 category (faces and scenes). The results revealed a significant main effect of model, $F(1, 232) =$
220 187.69 , $p < .001$, a significant main effect of layer, $F(12, 232) = 2.54$, $p = .003$, and a marginally
221 significant main effect of image, $F(1, 232) = 3.29$, $p = .071$. Furthermore, the interaction between
222 layer and model, $F(12, 232) = 6.26$, $p < .001$, and the interaction between image and model, $F(1,$
223 $232) = 636.57$, $p < .001$, were also significant. These findings suggest that the model with
224 perceptual expertise in each image category benefited more from FBA applied to recognize
225 objects from that category.

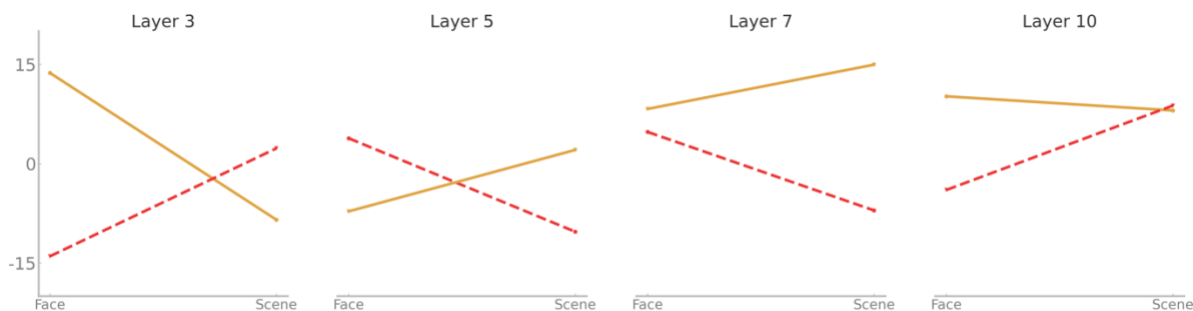
226

Image Category	Model Type	M	SD
Faces	Scene-expert	9.16	6.03
	Face-expert	17.71	5.93
Scenes	Scene-expert	29.23	5.78
	Face-expert	0.34	9.15

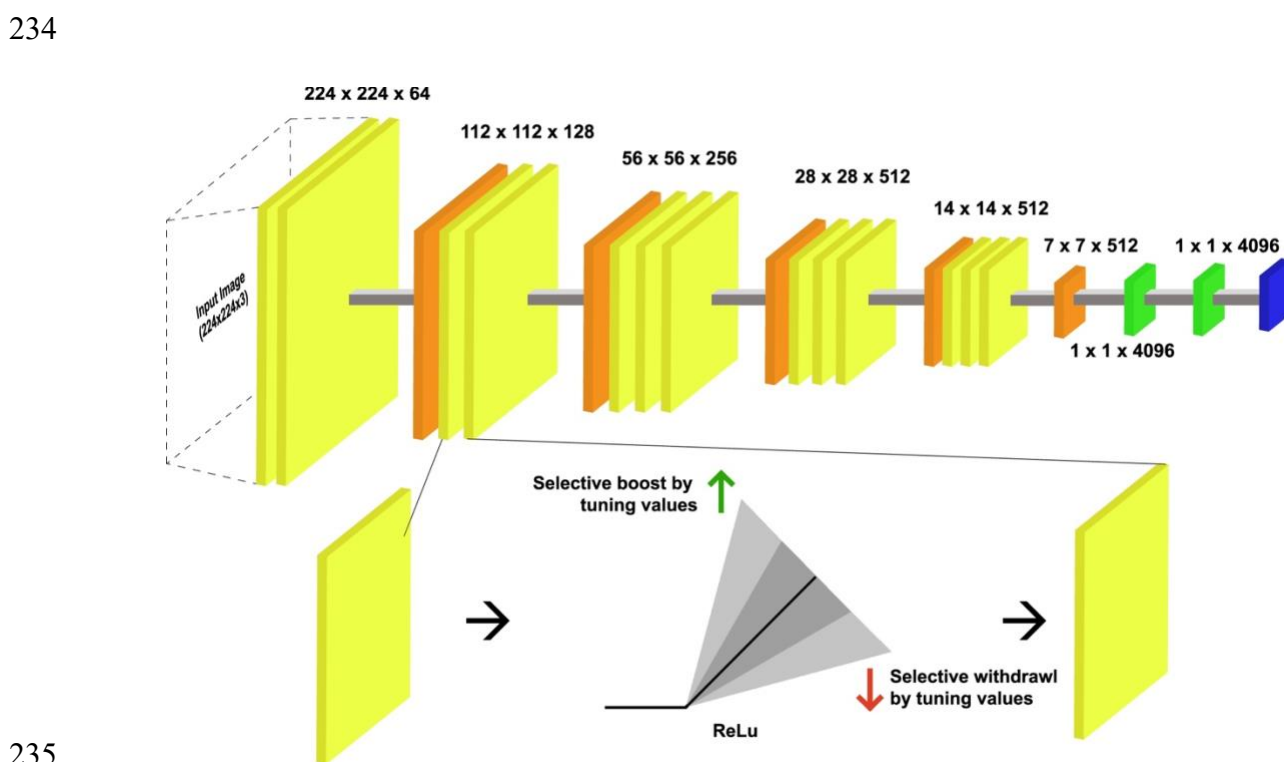
227

228 *Table 2. Summary statistics of performance improvements (i.e., $\Delta\%$) in the Face-expert and Scene-expert model across different*
229 *image categories (faces and scenes) averaged across convolutional layers (1-13).*

230



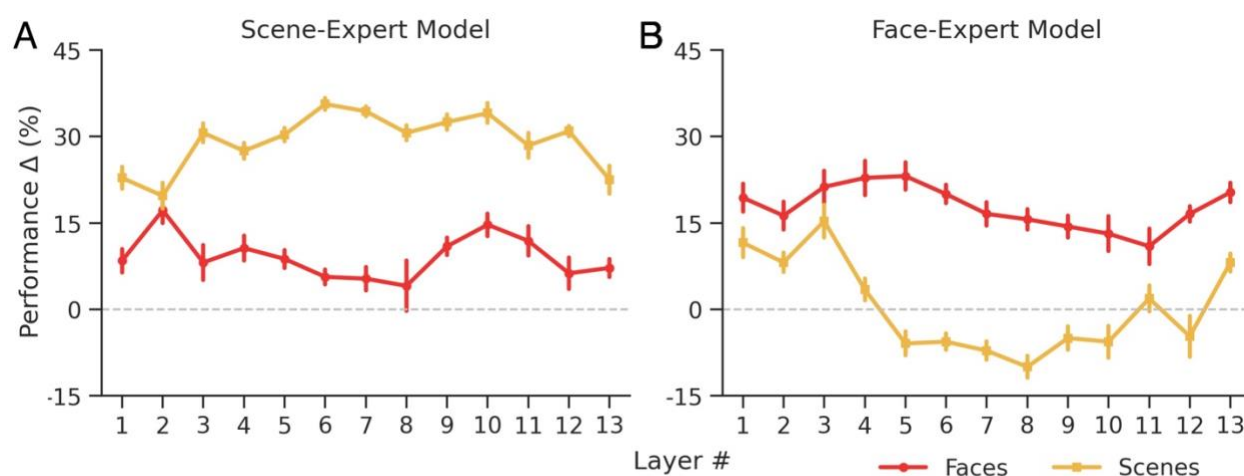
231
 232 *Figure 2. Example tuning curves of units (2 randomly chosen units are shown here) in each layer. From these tuning curves the*
 233 *preference of a neuron towards face or scene is determined.*



235
 236
 237 *Figure 3 Schematic of the FBA implementation in the model. The slope of the Rectified Linear Unit (ReLU) activation function is*
 238 *modulated based on the tuning values of the neuron. If a certain unit in a layer prefers the attended object category, the slope of*
 239 *the ReLU function is tuned-up (green arrow) whereas if a unit does not prefer the attended category, its slope is tuned-down (red*
 240 *arrow). See the Methods section for more information about how FBA was applied.*

241 To ascertain the specificity of attentional modulation observed, we conducted a control
242 experiment involving non-specific modulations of unit activity. Specifically, we applied non-
243 specific modulations to the neurons in both models by randomly shuffling the tuning properties.
244 For each neural unit within a given layer, we derived the tuning properties from randomly
245 permuted sets of tuning values. We did not observe any significant increases in task performance
246 using non-specific scaling. These results rule out the possibility of non-specific interactions
247 influencing the attention-performance relationship. Instead, they affirm the facilitating influence
248 of FBA on task performance, particularly highlighting its capacity to interact in an expertise-
249 specific manner. Furthermore, by employing randomly permuted tuning values instead of equal
250 scaling values (as used in (1)), our findings extend previous research. It is noteworthy that merely
251 permuting the labels of tuning curves does not yield the same effect of attention, underscoring
252 the effectiveness of FBA implemented via the FSGM.

253



254

255 *Figure 4. Outcomes of applying FBA to VGG16 pretrained on ImageNet (A) and VGGFace (B), across categories. Differential*
256 *specificity of categories can be observed in terms of performance increases. (A) For Scene-expert model, FBA increased the*
257 *performance of detecting the presence versus absence of scenes more than detecting the presence versus absence of faces. (B)*

258 *For the Face-expert model, FBA was effective for enhancing the performance of detecting the presence versus absence of faces;*
259 *for detecting the presence versus absence of scenes, the FBA's effect was not very helpful, and could even be negative (i.e.,*
260 *decreasing the performance of the model).*

261 Potential Mechanisms of Enhanced Effectiveness of FBA in Expert Networks

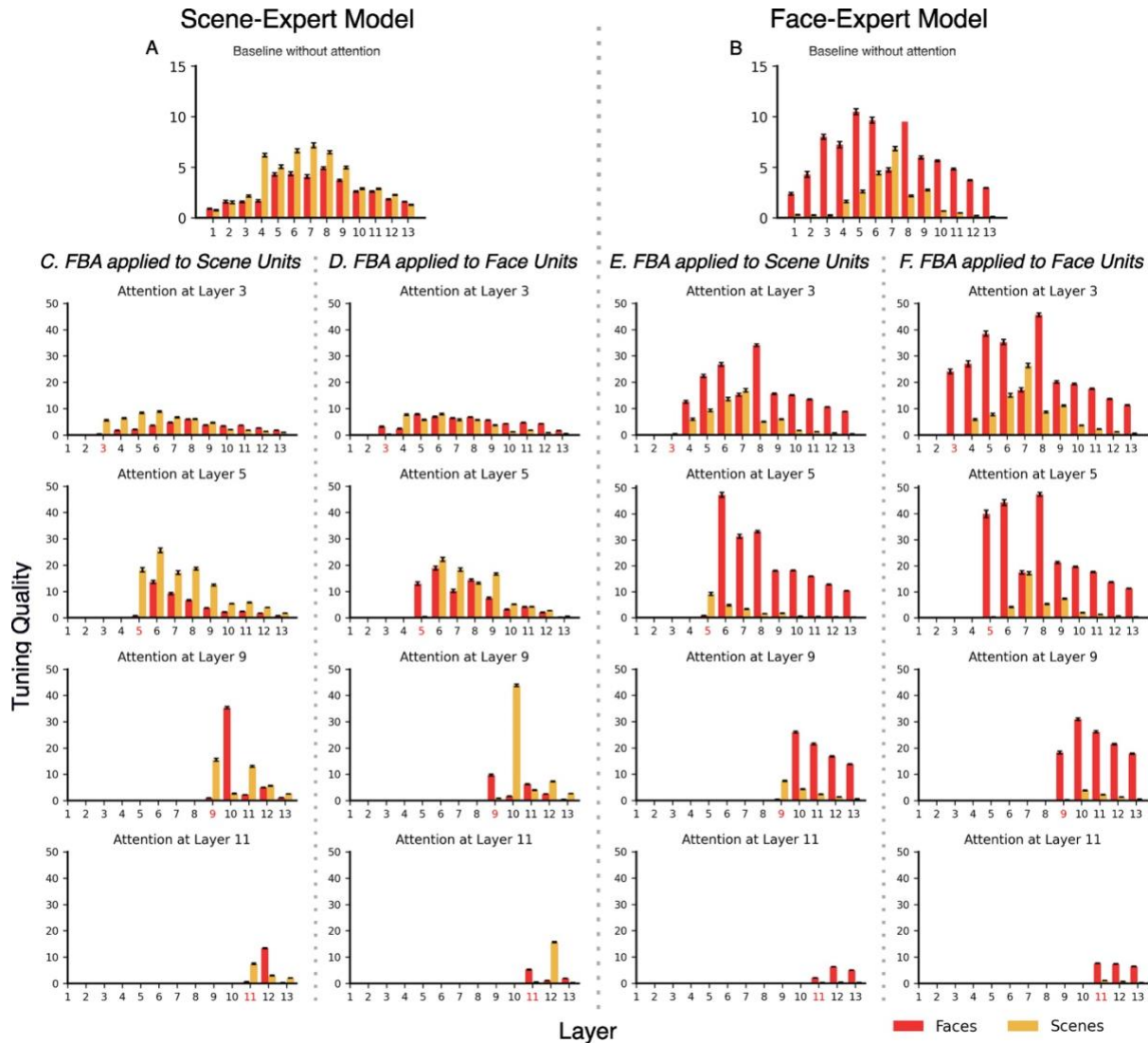
262

263 Studies have shown that attention modulates a neuron's response according to the neuron's
264 tuning properties in early and late layers of the visual system (45, 46). In these studies, the
265 response of a neuron is substantially enhanced when an optimal stimulus is attended, whereas
266 its response to an attended non-optimal stimulus is enhanced to a lesser extent, or even
267 decreased (47, 48). So, we probed if there was a similar effect of attention on neurons' tuning
268 properties in our models, and if it was related to expertise. For this, we analyzed the tuning curve
269 of each neuron and computed its tuning quality, which was its maximum value. Tuning quality
270 provided a quantifiable measure of the strength with which a particular neuron exhibited
271 preference for a specific object category. This was done separately for scene- and face-selective
272 neurons, and then compared to each other. The distribution of tuning quality across layers was
273 analyzed during the baseline condition (without attention), and when attention was applied to
274 neural units that preferred the target category for each layer individually. This enabled us to
275 investigate whether there was any expertise and/or image category preference-related
276 differences that resulted in differential effects of attention in the previous behavioral analysis.
277 Figures 5A and 5B show the results without attention (see Methods for a description of the
278 statistical tests). During the baseline condition, neurons in the Scene-expert model demonstrated
279 a significantly stronger preference for scenes, as evidenced by the higher tuning quality for

280 scenes compared to faces ($p < 0.001$ for layers 3 to 12, one-sided Welch's t-tests across layers,
281 FDR corrected for multiple comparisons). Conversely, the Face-expert model exhibited
282 significantly greater tuning quality for faces compared to scenes, suggesting a stronger
283 preference for faces ($p < 0.001$ for all layers except layer 7, one-sided Welch's t-tests across
284 layers, FDR corrected for multiple comparisons).

285 Given that feature based attention is multiplicative, for the same strength of attention
286 modulation (i.e., the same β value), the face-selective neurons in the Face-expert will have a
287 stronger increase in activity than the scene-selective neurons, resulting in the higher
288 performance improvement in detecting the presence versus absence of faces in the input. The
289 same principle could explain the reason why in the Scene-expert case, for the same strength of
290 attention modulation, there is a higher performance improvement for detecting the presence
291 versus absence of scenes in the input. We examine these ideas next.

292



293
 294 *Figure 5. Tuning quality across layers in the Scene-expert (A) and Face-expert (B) network divided into face (red) and scene (yellow)*
 295 *selective neurons during baseline when attention was not applied. Bars indicate the tuning quality distribution of neurons across*
 296 *layers 1 through 13. Tuning quality of neurons that prefer scenes is higher than that that prefer faces in the Scene-expert Model*
 297 *(A) and vice-versa in the Face-expert Model (B). (C-F) Tuning quality distribution divided based on FBA applied to scene and face*
 298 *selective neural units when attention is applied at different layers (row-wise layers 3,5,9 & 11 shown) in the Scene-expert Model*
 299 *(C-D), and the Face-expert Model (E-F).*

300
 301 In Figures 5C - F, we show the results when attention was applied to investigate whether
 302 attention modulated the tuning quality across different image types. We analyzed the effect

303 separately for scenes and faces in the two models while attention was applied at different layers,
304 showing representative results from layers 3, 5, 9, and 11 in the figure. Applying attention to
305 different layers resulted in stronger tuning quality of neural units compared to baseline, in
306 manner that is in accordance with the principles of the FSGM model. However, the degree of
307 modulation varied across models and image categories. In the Scene-expert model, as illustrated
308 in Figure 5C, applying attention separately to scene-selective units enhanced their tuning quality
309 (yellow bars), surpassing that of face-selective units (red bars). This effect began at the layer
310 where attention was applied and persisted through higher layers, maintaining the higher tuning
311 quality of scene-selective units throughout the model. In contrast, when attention was directed
312 to face-selective units (Figure 5D), their tuning quality (red bars) did not exceed that of scene-
313 selective units (yellow bars). Although face-selective units experienced improvement in tuning
314 quality at the layer where attention was applied (in line with the FSGM principle), this
315 enhancement did not sustain or carry over to higher layers, unlike the consistent propagation
316 observed for scene-selective units. Similarly, in the Face-expert model (Figure 5E), attention
317 selectively improved the tuning quality of face-selective units, but this effect did not significantly
318 influence other types of units (Figure 5F). Thus, in each expert model, the tuning quality of
319 neurons was enhanced by attention only when their preference aligned with the category of
320 expertise specific to the model. Conversely, attention did not improve the tuning quality of
321 neurons whose category preference differed from the model's expertise category.

322

323

324 Representational Similarity Reveals Feature Separation in Models

325

326 In the previous section, we examined the interplay between effects of attention and expertise-
327 specific task performance improvement by analyzing the quality of neuron-level tuning in our
328 models. However, apart from competition and bias that has been observed at single-unit level,
329 the expertise-attention interaction may also arise due to differences in network-wide neural
330 representations. Prior research has demonstrated that top-down bias and task performance is
331 context dependent, with greater competition occurring between stimuli that are more similar or
332 closer to each other in terms of neural representation (47, 49). Therefore, we investigated the
333 impact of dissimilarity in neural representations of different categories on the underlying effects
334 of attentional bias.

335

336 First, we applied representational similarity analysis (RSA) to assess the degree of dissimilarity
337 and competition between different image types (faces and scenes) across layers of each model,
338 taking into account the model's domain of expertise. Next, we conducted the same analysis on
339 images when FBA was applied to neurons in different layers, selectively targeting either face- or
340 scene-preferring neurons. This approach enabled us to examine how attention modulates the
341 multivariate neural representation of each image type within its corresponding expert model.
342 Finally, we compared the effect of attention on RSA dissimilarity and competition between face
343 and scene representations within each model, highlighting how attentional mechanisms interact
344 with expertise to shape multivariate neural representations.

345 The representational dissimilarity matrices (RDMs) for face-expert and scene-expert models
346 across different layers are shown in Figure 6A (top-left panel). The patterns of dissimilarity
347 change across layers, indicating that the models process and differentiate each image category
348 from the other image category at each processing stage (layer) of the network. A theoretical
349 RDM, as illustrated in Figure 6B, represents an idealized categorical structure, where two
350 categories of images (faces and scenes) are perfectly separated. It serves as a reference for
351 understanding the relation between face and scene representations in the two models. We
352 assessed the relation between RDMs obtained from each layer in the two models (baseline
353 without attention) and the theoretical RDM using rank-ordered Spearman correlations. This
354 yielded 26 (13 layers x 2 models) RSA correlation measures, each depicting the degree of
355 dissimilarity or separation between face and scene features encoded in the model layers (Figure
356 6C shows mean RSA correlation values with ± 1 SEM, bootstrapped across 100 iterations, baseline
357 without attention). Throughout the layers, we observed modest dissimilarity measures ($0.1 \pm$
358 0.03) except for early layers in the Face-expert model, which exhibit higher dissimilarity between
359 the features. These values are reasonable and within the bounds of what has been reported in
360 similar studies (50, 51).

361
362 Next, we computed the relation between theoretical RDM and the data based RDMs from each
363 layer (in a model) when attention was separately applied to neural units that either prefer faces
364 (Figure 6E) or scenes (Figure 6F). The mean RSA correlation values, along with ± 1 SEM, were
365 bootstrapped across 100 iterations and shown for layers 2, 4, 7, and 9 as representative
366 examples. In both models, we observed consistently positive mean RSA correlations across all

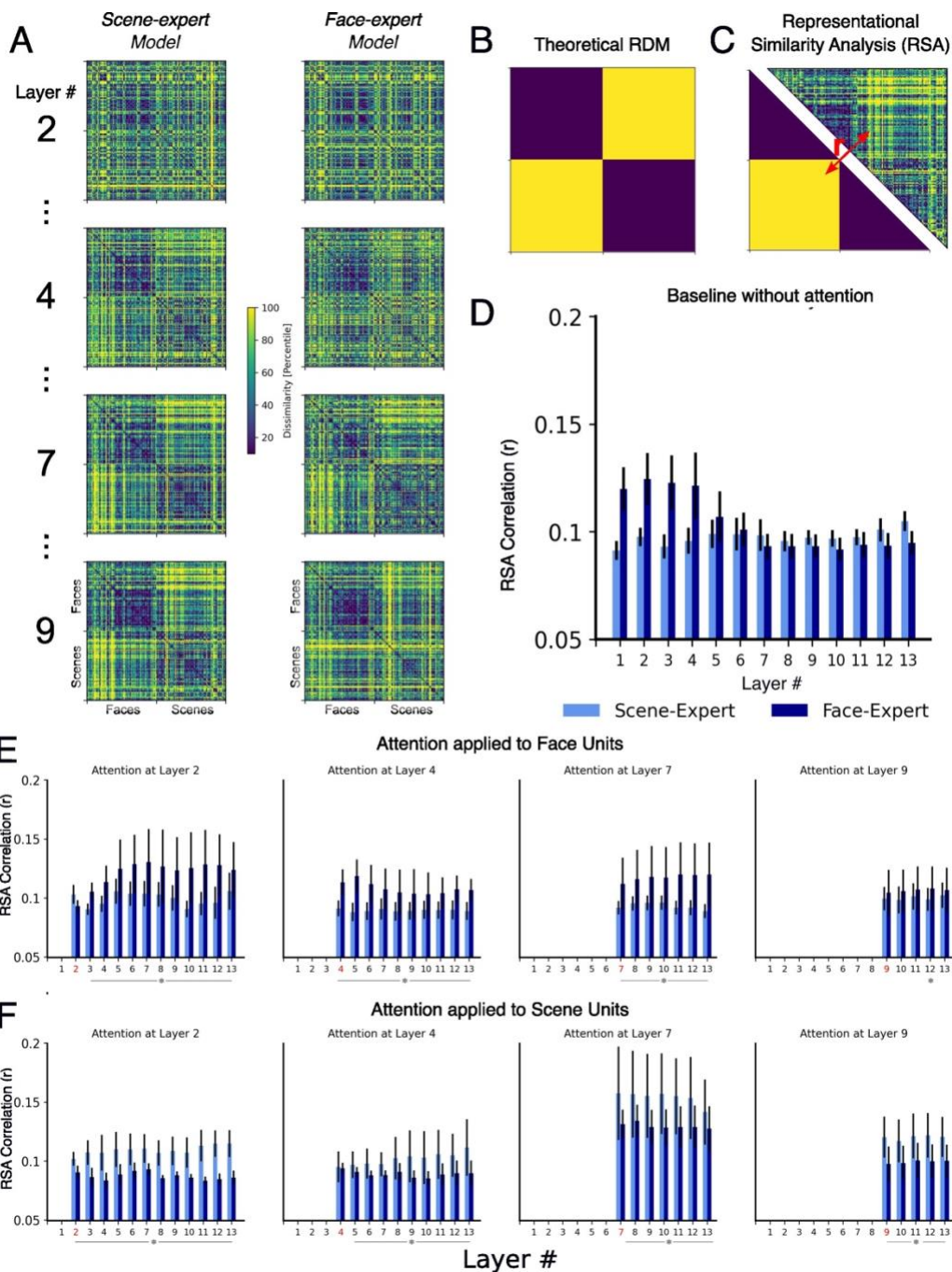
367 layers, with representational similarity patterns varying based on model expertise and the type
368 of neurons to which attention was applied. When attention was applied to face-selective
369 neurons, the similarity values were higher in the Face-expert model compared to the Scene-
370 expert model. This indicates the effectiveness of attentional mechanisms in inducing greater
371 separability between face representations and scene representations in the Face-expert model
372 ($p < 0.05$, one-tailed paired t-tests, FDR corrected for multiple comparisons across layers).
373 Similarly, when attention was applied to scene-selective neurons, RSA correlation values were
374 greater in the Scene-expert model compared to the Face-expert model ($p < 0.05$, one-tailed
375 paired t-tests, FDR corrected for multiple comparisons across layers). Therefore, the Scene-
376 expert model was more efficient in distinguishing scenes from faces when attention was applied
377 to scene-selective neurons only. Similarly, the Face-expert model was able to identify faces from
378 scenes efficiently when attention was applied to face-selective neurons only.

379

380 The expertise in each model allows for finer discrimination in the multivariate neural
381 representations between object features belonging to their category of expertise. Since the
382 effectiveness of attention depends on the neural representation of stimulus features, which is in
383 turn dependent on the expertise of the model, it is more effective when domain-specific objects
384 are distinctly encoded in various layers. This finding supports the feature similarity gain model of
385 attention, demonstrating that attention selectively enhances the processing of features relevant
386 to the model's expertise.

387

388



389

390 *Figure 6. Representational Similarity Analysis (RSA) across different categories of images and models. (A) Representational*

391 *dissimilarity matrices (RDMs). For each layer within a model, separate RDMs were constructed for all scene and face images by*

392 *one minus the Pearson r correlation between each pair of image-evoked multivariate neural activations. Here, RDMs are shown*

393 *for layer 2, 5, 9 and 12 for each model. (B) Theoretical RDM representing the ideal degree of separation between Scene (Manmade*

394 *and Natural) and Face (Male and Female) images. (C) RSA analysis, performed by calculating the rank-ordered Spearman*

395 *correlation between the off-diagonal triangular values of the theoretical RDM and layer RDMs. (D) RSA analysis for each model,*

396 *layer-wise during the baseline condition when attention was not applied to any neural units in Scene-expert (light blue) and Face-*
397 *expert (deep blue) models. (E, F) Representational Similarity (in Spearman rho correlation) when attention was applied to Face-*
398 *selective units (E) and Scene-selective units (F) at layers 2,4,7 & 9 (left to right, highlighted with red layer labels on the x-axis).*
399 *Error bars indicate ± 1 SEM obtained from bootstrapping technique using 100 samples. * = $p < 0.05$, one-tailed paired t tests and*
400 *FDR corrected for multiple comparison across layers.*

401 Discussion

402

403 In this work, we investigated whether and how perceptual expertise interacts with the effect of
404 attention in a convolutional neural network (CNN) model of the primate ventral visual system.
405 Specifically, we used neuron-level tuning and multi-neuron network-wide neural representation
406 to examine the mechanisms underlying the observed interaction between model expertise and
407 top-down attention. Our results complement the previous neuroimaging and
408 electrophysiological studies suggesting that top-down attention control interacts with neural
409 pathways and brain regions associated with perceptual expertise and enhance performance in
410 an expertise-dependent fashion (27, 29, 52-54).

411

412 The interplay of expertise and attentional bias

413

414 We trained a VGG16 neural network model to specialize in recognizing either scenes or faces.
415 The Scene-expert model was trained with scene images, while the Face-expert model was trained
416 with face images. During the experiment, the images were presented in two forms: regular
417 images, which were single images from a category (faces or scenes), and superimposed images,

418 which were superimposed images from the same or different categories. The results established
419 the presence of category-based expertise in the models. Specifically, the Face-expert model
420 demonstrated superior performance in recognizing the presence or absence of faces in the single
421 image input, while the Scene-expert model excelled at detecting the presence or absence of
422 scenes in the single image input. When presented with superimposed images, the recognition
423 performance declined significantly, indicating that the models were not effective at dealing with
424 distractors in the superimposed images. Next, we applied feature-based attention to neurons
425 based on their tuning profiles. The Scene-expert model showed significantly larger improvements
426 in recognition performance with attention when detecting the presence vs absence of scenes as
427 compared to detecting the presence versus absence of faces. Similarly, the Face-expert model
428 showed significant improvement with attention when detecting the presence versus absence of
429 faces than detecting the presence vs absence of scenes. Expertise in a specific category allowed
430 the models to develop specialized neural representations that are optimized to interact with
431 attention and improve task performance.

432

433 [Expert Versus Novice: Unit Level Analysis of Object Based Attentional Enhancement](#)

434

435 To understand why attention is more effective when combined with expertise, we analyzed the
436 tuning quality of artificial neurons, which reflects how strongly they prefer an object category.
437 Our findings indicate that the Scene-expert model exhibits stronger tuning quality for scenes than
438 for faces, while the Face-expert model shows stronger tuning quality for faces than for scenes.
439 According to the Feature-Specific Gain Modulation (FSGM) model, attention modulates the

440 neuron's firing by a multiplicative factor applied to the tuning function (43, 44, 55). Given that
441 the tuning function here is computed over the two categories of images, this operation will lead
442 to the enhancement of the attended category (target) and the suppression of the ignored
443 category (distractor). Given the baseline tuning quality findings mentioned above, the
444 multiplicative nature of the attention amplification further ensures that the attention
445 enhancement is stronger when it is directed to the image category in which the network is an
446 expert. Thus, attention is more effective in experts because it operates on neurons that are
447 already highly tuned towards the object of expertise. This selective tuning, in addition to already
448 sharpened tuning due to expertise, likely underlie the enhanced effect of attention when the
449 attended object matches the expertise of a model. Neurophysiological studies of experts and
450 novices consistently report this phenomenon, highlighting how experts often exhibit increased
451 activation levels within brain regions relevant to their task and sometimes recruit additional areas
452 involved in domain-specific processing (23, 26, 56, 57). For instance, Gauthier et al. (1999) found
453 that car experts showed greater activation in the fusiform face area (FFA) when recognizing cars
454 compared to novices, indicating specialized processing. Similarly, Maguire et al. (2002) observed
455 that London taxi drivers, who are experts in navigation, had larger hippocampal volumes and
456 showed increased activation in this region during navigational tasks. Grabner et al. (2006)
457 reported that individuals with high mathematical expertise exhibited greater activation in brain
458 areas associated with arithmetic processing, and McGugin et al. (2014) found that bird and car
459 experts showed enhanced activation in category-specific regions when recognizing birds and
460 cars, respectively. These studies collectively suggest that expertise enhances the neural efficiency
461 and sharpens the tuning curve in favor of features that match its domain of expertise.

462

463 Population Level Analysis: Representational Similarity of Targets and Distractors

464

465 It has become increasingly clear that visual perception relies on multivariate representations of
466 visual inputs. Stronger separation between the neural representations of targets and distractors
467 is crucial for the effectiveness of attentional selection of task-related stimulus information.
468 Previous studies have demonstrated that efficiency of these mechanisms is heavily influenced by
469 individual differences in the representation of task-related information (58-61). Specifically,
470 unique category representations, formed through top-down attentional templates, are pivotal in
471 guiding the search for targets and suppressing distractors. These templates are shaped by an
472 individual's experience and expertise, meaning that a person's ability to effectively utilize
473 attentional mechanisms depends on how well their mental representations align with the task at
474 hand.

475

476 According to the biased competition model of attention, visual stimuli compete for neural
477 representation in the brain, and the more similar the targets and the distractors, the stronger the
478 competition (47, 62, 63). Conversely, the more dissimilar the targets and the distractors, the
479 weaker the competition, the better the behavioral performance. Our results are consistent with
480 this. In our study, Representational Similarity Analysis (RSA) revealed that at the multivariate
481 neural representational level, attentional mechanisms are effective primarily for object
482 categories that align with the expertise of a model. When attention is directed towards neurons
483 within a specific layer of a model, this focus sharpens population neural tuning by increasing the

484 representational distance between targets and distractors. As the neural features of targets and
485 distractors become more dissimilar, the competition between them diminishes, thereby
486 enhancing the potency of attentional mechanisms.

487

488 Recent work by Doostani et al. (47) has shown that sharpened neuronal tuning amplifies the
489 influence of top-down mechanisms, especially in difficult scenarios where targets and distractors
490 share common features. Our findings offer direct evidence on the strength of the attentional bias
491 towards the target such that it increases as neuronal tuning sharpens. Our findings indicate that
492 this effect is more significant in models with expertise in specific object categories, as expertise
493 increases the dissimilarity between targets and distractors, thereby sharpening the tuning curve.
494 This implies that specialized knowledge enhances the ability of attentional mechanisms to
495 distinguish between similar features.

496

497 [Relation with Prior Literature](#)

498

499 The association between perceptual expertise and attention mechanisms has been studied in the
500 past. It has been reported that expertise has a facilitatory effect on categorization and involves
501 deployment of top-down mechanisms to engage object processing as: (1) Engagement of
502 attention transcends lower-level features of the object and prioritizes them based on their
503 content (64, 65). Bukach et al. (2006) demonstrated that experts in face recognition could
504 categorize faces more efficiently than novices, suggesting that expertise enhances the ability to
505 process complex visual stimuli. Similarly, van der Linden et al. (2014) found that expert

506 radiologists could detect abnormalities in medical images more accurately, indicating that
507 expertise involves sophisticated attentional mechanisms. (2) Expertise entails top-down activity
508 crucial for evaluating and recognizing pertinent stimulus features (29, 66). Harel et al. (2010)
509 showed that experts in visual search tasks could quickly identify target objects among distractors,
510 highlighting the role of top-down processes in expertise. Reddy et al. (2007) found that expert
511 athletes could anticipate the actions of opponents more effectively, demonstrating the
512 importance of top-down attention in dynamic environments. (3) Studies have also shown that
513 experts exhibit differences in attentional selection and interacting with objects that are relevant
514 to them. Stokes (2021) reported that expert musicians could focus on relevant musical elements
515 more efficiently than novices, indicating differences in attentional selection. Kundel et al. (2007)
516 found that expert radiologists took less time to fixate on abnormalities in medical images,
517 suggesting that expertise leads to faster identification of relevant features. Vogt & Magnussen
518 (2016) observed that expert chess players had shorter saccades and fewer fixations when
519 analyzing chess positions, indicating more efficient visual processing. This suggests that experts
520 possess a rapid sensitivity to holistic features of the stimulus array, indicating that their
521 attentional processes are more efficient and diagnostically relevant to the task. Our findings
522 extend upon prior work, by delving into the intricate interplay between biologically inspired
523 attention mechanisms and the intricate architecture of neural networks, and their predisposition
524 to knowledge. Also, our work helps by aligning network behaviors with established neural
525 processes, making it easier to decipher the rationale behind network decisions, fostering
526 robustness, transparency, and development of more interpretable DNNs.

527 Summary

528

529 This study was aimed at understanding the association between attention mechanisms and
530 perceptual expertise. We used deep neural networks as models of the ventral visual pathway of
531 the brain for this purpose. Our findings indicate that when a deep neural network is subjected to
532 complex tasks, its performance can be enhanced by introducing attentional mechanisms, and the
533 effectiveness of this attention enhancement is closely tied to the perceptual expertise developed
534 in the model, i.e., attention is more effective in a model when the task is within the model's
535 domain of expertise. Mechanisms at the individual neuronal response level and at the neural
536 population level were investigated. Methodologically, investigating neural mechanisms of
537 perception in deep learning models represents a convergence of AI and neuroscience, offering
538 (1) a computational platform in which manipulations not practical in empirical experiments can
539 be carried out and (2) the potential to build more efficient, adaptive, and human-like artificial
540 intelligence systems. As our understanding of both artificial and biological neural networks
541 advances, we can expect more refined and effective models to emerge that can better bridge the
542 gap between artificial and natural intelligence.

543

544

545

546 Methods

547

548 The Model

549

550 VGG16 is used as a model of the ventral visual stream (36). VGG16, as shown in Figure 1A, is a
551 feed forward convolutional neural network with 13 convolutional layers followed by 3 fully
552 connected layers. In this work, two pretrained VGG 16 models were considered: one pretrained
553 on the ImageNet dataset and thus an expert in object recognition and the other pretrained on
554 the VGGFace dataset and was thus an expert in face recognition (37, 38). For the model trained
555 on the ImageNet dataset, the last layer of the network outputted the labels of 1000 object
556 categories (92.7% accuracy), whereas for the model trained on the VGGFace dataset, the last
557 layer outputted the labels of 2622 individual faces (98.95% accuracy). In this work, the network
558 weights connecting the first 15 layers of the VGG 16 models were taken from (36) and from (38),
559 respectively, and kept unchanged (frozen). The last layer was replaced by a layer with output
560 suitable for our paradigms (see below). The weights connecting the paradigm-specific last layer
561 and the layer preceding it were trained using the datasets described below according to the goals
562 of the paradigm.

563

564

565

566 Image Category Detection Task

567

568 The models were presented with images containing objects from two categories: faces (male
569 faces & female faces) and scenes (natural scenes & manmade scenes). The model's task was to
570 identify in the image the presence or absence of objects from one out of the two categories (e.g.,
571 is there a face in the image?). To accomplish this binary classification, we replaced the final
572 SoftMax layer of the original pretrained VGG16 models with a layer containing a series of binary
573 classifiers, with each classifier consisting of two units capable of signaling the "presence" or
574 "absence" of the to-be-detected category (Figure 1B). Specifically, the Face-expert has two
575 associated binary classifiers with one detecting the presence and absence of a face in the input
576 and one detecting the presence and absence of a scene in the input, and the Scene-expert also
577 has two associated binary classifiers that perform the same tasks.

578

579 To train the binary classifiers, we used a separate set of training and testing data, sourced from
580 different image repositories (39-41), which did not overlap with ImageNet or VGGFace datasets.
581 The dataset consisted of 200 faces and 200 scenes (224x224 pixel RGB images); see Figure 1C for
582 examples. We used 160 faces and 160 scenes for training and the remaining 40 faces and 40
583 scenes for testing. During the training of the binary classifiers, depending on the category to be
584 detected, there were always 160 true positives from the category along with 160 true negatives
585 from the other category. The classification performances reported here were achieved by
586 implementing logistic regression in the binary classifiers.

587

588 Once the binary classifiers were trained, they were tested separately on regular as well as
589 superimposed images. When the input to the network model was a regular image, the network's
590 task was the same as the task it was trained on, namely, to detect the presence or absence of a
591 particular category in the input image. To challenge the model, besides using regular images, we
592 introduced superimposed images to make the task more difficult. The superimposed images
593 were created by transparently superimposing two images either from the same category or from
594 different categories (by taking the mathematical average of the corresponding pixels of the two
595 images). There were three types of superimposed images: (1) face over face, (2) scene over scene,
596 and (3) face over scene (which is equivalent to scene over face). For detecting the presence or
597 absence of a face, (1) and (3) were true positives whereas (2) were true negatives. For detecting
598 the presence or absence of a scene, (2) and (3) were true positives whereas (1) were true
599 negatives. The test images were balanced with 50% true positives and 50% true negatives; a total
600 of 40 positive images and 40 negative images were used for testing.

601

602 It is expected that the detection accuracy for superimposed images would decline significantly
603 relative to regular images. This scenario then provided us the opportunity to apply feature-based
604 attention to enhance the detection performance and test whether such enhancement depends
605 on the expertise of the network. We expected that for the Face-expert, applying attention to face
606 would enhance face detection performance more than applying attention to scene would scene
607 detection performance, and that for the Scene-expert, the opposite is true.

608

609 Attention Modulation of Neuronal Responses

610

611 The main goal of this study was to examine how attention interacts with perceptual expertise in
612 a deep neural network model of the ventral visual stream. Specifically, we implemented the
613 feature similarity gain modulation (FSGM) model of feature-based attention, following the
614 procedures of (1). To apply this attention mechanism, we calculated the extent to which each
615 neuron in a model (Scene-expert or Face-expert) preferred a certain object category, i.e., their
616 tuning values. Here the term neuron was used to refer to a filter or a feature map in the model.
617 Since feature-based attention is a spatially global phenomenon (1, 67), the responses from all the
618 neural units within a filter were averaged to become the response of the filter or neuron. We
619 implemented attention by modulating the slope of the ReLu function of the units within a filter
620 according to the filter's tuning function.

621

622 Calculation of Tuning Values

623

624 To determine the tuning function of each neuron, we presented the model with the regular
625 images from the two categories (the same set of images of 160 faces and 160 scenes) used for
626 training the binary classifiers) and measured the relative activity levels of the units within the
627 filter. As indicated above, considering that feature-based attention is nonspatial (1, 67), we
628 treated the activity levels of all units within a filter identically and calculate its tuning by z-scoring
629 their activity across categories using the following equations:

630

631

$$\bar{r}^{lk} = \frac{1}{N} \sum_{n=1}^N r^{lk}(n)$$

632

633

634

$$f_c^{lk} = \frac{\frac{1}{N_c} \sum_{c \in [1,2]} r^{lk}(n) - \bar{r}^{lk}}{\sqrt{\frac{1}{N} \sum_{n=1}^N (r^{lk}(n) - \bar{r}^{lk})^2}}$$

635

636

637 Specifically, for the k^{th} neuron in the l^{th} layer, $r^{lk}(n)$ is defined as the mean activity of all units in

638 the filter in response to n^{th} image. Finally, by taking the mean across all these values for the

639 training set images ($N_c = 160$ images per category, $c = c^{\text{th}}$ category with $c = 1$ or 2 because there

640 are 2 categories in total; $N = 2 \times 160 = 320$), we get the mean activity of the neuron \bar{r}^{lk} . The

641 tuning value for each neuron for a given category is the z-scored mean activity with respect to

642 the mean activity of the unit for all images. Put it simply, tuning value for a certain category for a

643 neural unit is the average activity of the unit in response to the images from the category

644 subtracting the mean activity across all images and divided by the standard deviation of the

645 activity across all images. Calculated across the two categories, we get a 2-dimensional vector of

646 values f_c^{lk} , which is the tuning function of each neuron used to implement attention. To find the

647 preferred category of a neural unit, we designate the category with the larger tuning value as its

648 preferred category. Tuning quality for a certain neuron is defined as the maximum tuning value

649 of that neural unit i.e., $\max(|f_c^{lk}|)$. Tuning quality is a measure of the extent of how strong a

650 certain neuron prefers its most preferred category.

651

652 Implementation of Feature Based Attention

653

654 We implemented FSGM attention model in its multiplicative and bidirectional form across the
655 layers in the two networks. To apply attention at a unit from neuron k in layer l and for category
656 c , we modulated the slope of the corresponding rectified linear unit (ReLU) by the tuning value
657 of that category c , weighted by a strength parameter β (varied from 0 to 20 in increments of 0.1.

658

$$659 \quad x_{ij}^{lk} = (1 + \beta f_c^{lk})[(x - 1)_{mn}^{lk}]_+$$

660

661 where x_{ij}^{lk} is the unit response at $(i, j)^{th}$ spatial location in the k^{th} neuron of the l^{th} layer, $[\]_+$ is the
662 ReLU function. $(x - 1)_{mn}^{lk}$ represents the activity of the $(m, n)^{th}$ unit from the preceding layer.

663

664 Tuning Quality Analysis

665

666 As mentioned previously, tuning quality was defined as the magnitude of the maximum tuning
667 value of a neuron: $\max(|f_c^{lk}|)$. Consequently, it was a measure of the relative strength of the
668 neuron's preference towards its favored object category. It is reasonable to expect that for the
669 Face-expert, the tuning quality for face stimuli will be higher than that for scene stimuli, whereas
670 for the Scene-expert, the tuning quality for scene stimuli will be higher than that for face stimuli.
671 For quantitative analysis, since there were unequal numbers of neurons selective to each
672 category in a layer, we performed a comprehensive statistical assessment. This involved

673 subjecting the layer-wise tuning quality distributions of the two models focusing on the
674 categories – faces and scenes, to Shapiro-Wilk normality and Levene’s variance tests. The results
675 of the Shapiro-Wilk test revealed a substantial deviation from the assumption of normality,
676 consistently observed across layers (except layers 1-3 and layer 13) and model variations ($p <$
677 0.001). Furthermore, acknowledging this deviation, Levene’s test indicated the presence of
678 uneven variances in both distributions ($p < 0.001$, spanning all layers in both models). Therefore,
679 the two distributions were compared using Welch’s t-tests, FDR corrected for multiple
680 comparisons.

681

682 Representational Similarity Analysis (RSA)

683

684 For each model and its individual layers, we calculated the output activity of the neurons present
685 in each layer for every image. These neural representations were then analyzed using a
686 Representational Similarity Analysis (RSA) that consisted of two primary steps. Firstly, we
687 generated separate Representational Dissimilarity Matrices (RDMs) for each model type (Figure
688 6A), obtaining one RDM per layer. These matrices were based on the distinct patterns of
689 activation that each image elicited, categorized separately for each object category. Secondly,
690 employing a theoretical RDM illustrated in Figure 6B, we captured the idealized maximum
691 possible divergence between the two categories. Lastly, we computed the representational
692 similarity by calculating the rank-ordered Spearman correlation between the theoretical RDM
693 and the RDM for each layer of the VGG16 model.

694

695 For each model, we subjected it to all possible regular images and extracted the activation
696 patterns across all neurons for each layer. This was done separately for each face (total=160) and
697 scene image (total=160). The activation pattern for each image was then transformed into a one-
698 dimensional vector. To assess dissimilarity between these vectors, we calculated one minus the
699 Pearson's correlation coefficient for every pair of vectors. This process resulted in the generation
700 of 320 x 320 sized Representational Dissimilarity Matrices (RDMs) for each layer, model variant,
701 and image category (face and scene), amounting to a total of 13 x 2 matrices. Within each RDM,
702 the cells contained dissimilarity values $(1 - r)$ representing the dissimilarity in neural
703 representations between pairs of images. Subsequently, a theoretical RDM was constructed to
704 match the dimensions of the VGG16 model layer's RDMs (320 x 320). In this theoretical RDM, the
705 cells in the top-left 160 x 160 and bottom-right 160 x 160 sections along the diagonal were
706 assigned a value of 0, while all other cells were assigned a value of 1 (Figure 6B). This
707 configuration indicated minimum dissimilarity (0) for images within the same category (e.g., face
708 vs face or scene vs scene), and maximum dissimilarity (1) for images belonging to different
709 categories (e.g., face vs. scene or scene vs. face). This method allowed for a systematic
710 comparison of dissimilarity between different categories and facilitated the evaluation of the
711 model's ability to distinguish between faces and scenes at various layers.

712

713 Finally, representational similarities were computed as the rank-ordered Spearman correlation
714 between each RDM from layers of the two models and the theoretical RDM. This process resulted
715 in a set of similarity values (13x2), corresponding to each layer and model of expertise. To cross
716 validate across our test set, we used a bootstrapping technique to assess the statistical

717 significance of these values. We generated 100 test image sets with replacement, recalculating
718 representational similarity for each sample. This process yielded an empirical distribution of
719 these values, along with bootstrapped mean estimations and 95% confidence intervals. To
720 identify significant differences in the mean similarity values across object categories, we
721 employed a p-value threshold of 0.05. We rejected the null hypothesis if the confidence interval
722 did not include 0. Furthermore, we subjected the results across layers to FDR correction to
723 account for multiple comparisons across layers.

724

725 [RSA with attention](#)

726

727 We repeated the same RSA analysis procedure described above when attention was applied at
728 each layer individually to neurons that prefer faces or scenes (Figure 6E & F). Specifically,
729 attention was modeled by amplifying the activation of specific units that exhibited a stronger
730 response to either faces or scenes within each layer of a model (separately for the scene-expert
731 and face-expert model). The attention-modulated activations were then processed similarly to
732 the unmodulated activations: we calculated the output activity of each regular image, generated
733 RDMs, and performed RSA by correlating these attention-modulated RDMs with the theoretical
734 RDM.

735

736

737 Acknowledgements

738

739 This work was supported by NIH grant MH117991, NSF grant BCS2318886, and NSF grant
740 BCS2318984. We are grateful to Joy Geng, John Henderson, Ruogu Fang, Randall O'Reilly, Lee
741 Miller, and the members of our labs for their helpful comments and advice.

742

743 References

- 744 1. Lindsay GW, Miller KD. How biological attention mechanisms improve task performance in a large-scale visual system
745 model. *eLife*2018. p. 1-29.
- 746 2. Xu K, Ba J, Kiros R, Cho K, Courville A, Salakhudinov R, et al. Show, Attend and Tell: Neural Image Caption Generation
747 with Visual Attention. In: Francis B, David B, editors. Proceedings of the 32nd International Conference on Machine Learning;
748 Proceedings of Machine Learning Research: PMLR; 2015. p. 2048--57.
- 749 3. Cao C, Liu X, Yang Y, Yu Y, Wang J, Wang Z, et al. Look and Think Twice: Capturing Top-Down Visual Attention
750 with Feedback Convolutional Neural Networks. 2015 IEEE International Conference on Computer Vision (ICCV)2015. p. 2956-
751 64.
- 752 4. Yang X, Yan J, Wang W, Li S, Hu B, Lin J. Brain-inspired models for visual object recognition: an overview. *Artificial*
753 *Intelligence Review*. 2022;55(7):5263-311.
- 754 5. Kanwisher N, Gupta P, Dobs K. CNNs reveal the computational implausibility of the expertise hypothesis. *iScience*.
755 2023;26(2).
- 756 6. Cadena SA, Denfield GH, Walker EY, Gatys LA, Tolia AS, Bethge M, et al. Deep convolutional models improve
757 predictions of macaque V1 responses to natural images. *PLOS Computational Biology*. 2019;15(4):e1006897.
- 758 7. Bonner MF, Epstein RA. Computational mechanisms underlying cortical responses to the affordance properties of visual
759 scenes. *PLOS Computational Biology*. 2018;14(4):e1006111.
- 760 8. Yamins DLK, Hong H, Cadieu CF, Solomon EA, Seibert D, DiCarlo JJ. Performance-optimized hierarchical models
761 predict neural responses in higher visual cortex. Proceedings of the National Academy of Sciences of the United States of
762 America2014. p. 8619-24.
- 763 9. Mohsenzadeh Y, Mullin C, Lahner B, Oliva A. Emergence of Visual Center-Periphery Spatial Organization in Deep
764 Convolutional Neural Networks. *Scientific Reports*. 2020;10(1).
- 765 10. Wallis TSA, Funke CM, Ecker AS, Gatys LA, Wichmann FA, Bethge M. A parametric texture model based on deep
766 convolutional features closely matches texture appearance for humans. *Journal of Vision*. 2017;17(12):5.
- 767 11. Kuperwajs I, Schütt HH, Ma WJ. Using deep neural networks as a guide for modeling human planning. *Scientific*
768 *Reports*. 2023;13(1).
- 769 12. Peterson JC, Abbott JT, Griffiths TL. Evaluating (and Improving) the Correspondence Between Deep Neural Networks
770 and Human Representations. *Cognitive Science*. 2018;42(8):2648-69.
- 771 13. Jang H, McCormack D, Tong F. Noise-trained deep neural networks effectively predict human vision and its neural
772 responses to challenging images. *PLoS Biol*. 2021;19(12):e3001418.
- 773 14. Kell AJE, Yamins DLK, Shook EN, Norman-Haignere SV, McDermott JH. A Task-Optimized Neural Network
774 Replicates Human Auditory Behavior, Predicts Brain Responses, and Reveals a Cortical Processing Hierarchy. *Neuron*: Elsevier
775 Inc.; 2018. p. 630-44.e16.
- 776 15. Ivet Rafegasa MV, Lú is A. Alexandreb, Guillem Ariasa. Understanding Trained CNNs by Indexing Neuron Selectivity.
777 2019.
- 778 16. Ratan Murty NA, Bashivan P, Abate A, DiCarlo JJ, Kanwisher N. Computational models of category-selective brain
779 regions enable high-throughput tests of selectivity. *Nat Commun*. 2021;12(1):5540.
- 780 17. VanRullen R. Reconstructing faces from fMRI patterns using deep generative neural networks. 2019.
- 781 18. MCGugin RW, Van Gulick AE, Tamber-Rosenau BJ, Ross DA, Gauthier I. Expertise Effects in Face-Selective Areas are
782 Robust to Clutter and Diverted Attention, but not to Competition. *Cerebral Cortex*. 2015;25(9):2610-22.
- 783 19. Bukach CM, Phillips WS, Gauthier I. Limits of generalization between categories and implications for theories of
784 category specificity. *Attention, Perception & Psychophysics*. 2010;72(7):1865-74.
- 785 20. Brefczynski-Lewis JA, Lutz A, Schaefer HS, Levinson DB, Davidson RJ. Neural correlates of attentional expertise in
786 long-term meditation practitioners. Proceedings of the National Academy of Sciences. 2007;104(27):11483-8.
- 787 21. Wong YK, Folstein JR, Gauthier I. The nature of experience determines object representations in the visual system.
788 *Journal of Experimental Psychology: General*. 2012;141(4):682-98.
- 789 22. Zhang T, Dong M, Wang H, Jia R, Li F, Ni X, et al. Visual expertise modulates baseline brain activity: a preliminary
790 resting-state fMRI study using expertise model of radiologists. *BMC Neuroscience*. 2022;23(1).
- 791 23. Gauthier I, Tarr MJ, Anderson AW, Skudlarski P, Gore JC. Activation of the middle fusiform 'face area' increases with
792 expertise in recognizing novel objects. *Nature Neuroscience*. 1999;2(6):568-73.
- 793 24. Wong AC-N, Palmeri TJ, Gauthier I. Conditions for Facelike Expertise With Objects. *Psychological Science*.
794 2009;20(9):1108-17.
- 795 25. Xu Y. Revisiting the Role of the Fusiform Face Area in Visual Expertise. *Cerebral Cortex*. 2005;15(8):1234-42.
- 796 26. MCGugin RW, Newton AT, Gore JC, Gauthier I. Robust expertise effects in right FFA. *Neuropsychologia*. 2014;63:135-
797 44.
- 798 27. Gauthier I, Skudlarski P, Gore JC, Anderson AW. Expertise for cars and birds recruits brain areas involved in face
799 recognition. *Nature Neuroscience*. 2000;3(2):191-7.
- 800 28. Stokes D. On perceptual expertise. *Mind & Language*. 2021;36(2):241-63.

- 801 29. Harel A, Gilaie-Dotan S, Malach R, Bentin S. Top-Down Engagement Modulates the Neural Expressions of Visual
802 Expertise. *Cerebral Cortex*. 2010;20(10):2304-18.
- 803 30. Kok EM, Sorger B, Van Geel K, Gegenfurtner A, Van Merriënboer JJG, Robben SGF, et al. Holistic processing only?
804 The role of the right fusiform face area in radiological expertise. *PLOS ONE*. 2021;16(9):e0256849.
- 805 31. Martens F, Bulthé J, van Vliet C, Op de Beeck H. Domain-general and domain-specific neural changes underlying visual
806 expertise. *NeuroImage*. 2018;169:80-93.
- 807 32. Bilalic M, Langner R, Ulrich R, Grodd W. Many Faces of Expertise: Fusiform Face Area in Chess Experts and Novices.
808 *Journal of Neuroscience*. 2011;31(28):10206-14.
- 809 33. Stokes D. On perceptual expertise. *Mind & Language*. 2020;36(2):241-63.
- 810 34. Richler JJ, Wong YK, Gauthier I. Perceptual Expertise as a Shift From Strategic Interference to Automatic Holistic
811 Processing. *Current Directions in Psychological Science*. 2011;20(2):129-34.
- 812 35. Kanwisher N, Gupta P, Dobs K. CNNs Reveal the Computational Implausibility of the Expertise Hypothesis. *iScience*.
813 2023.
- 814 36. Simonyan K, Zisserman A. Very deep convolutional networks for large-scale image recognition. *arXiv preprint*
815 *arXiv:14091556*. 2014.
- 816 37. Deng J, Dong W, Socher R, Li L-J, Kai L, Li F-F. ImageNet: A large-scale hierarchical image database. 2009 IEEE
817 Conference on Computer Vision and Pattern Recognition 2009. p. 248-55.
- 818 38. Parkhi OM, Vedaldi A, Zisserman A. Deep Face Recognition. *Proceedings of the British Machine Vision Conference*
819 2015. p. 41.1-12.
- 820 39. Burge J, Geisler WS. Optimal defocus estimation in individual natural images. *Proceedings of the National Academy of*
821 *Sciences*. 2011;108(40):16849-54.
- 822 40. Wennekers T, Dhamecha TI, Singh R, Vatsa M, Kumar A. Recognizing Disguised Faces: Human and Machine
823 Evaluation. *PLoS ONE*. 2014;9(7).
- 824 41. Paterson K, Brodeur MB, Guérard K, Bouras M. Bank of Standardized Stimuli (BOSS) Phase II: 930 New Normative
825 Photos. *PLoS ONE*. 2014;9(9).
- 826 42. Treue S, Trujillo JCM. Feature-based attention influences motion processing gain in macaque visual cortex. *Nature*.
827 1999;399(6736):575-9.
- 828 43. McAdams CJ, Maunsell JHR. Effects of Attention on Orientation-Tuning Functions of Single Neurons in Macaque
829 Cortical Area V4. *The Journal of Neuroscience*. 1999;19(1):431-41.
- 830 44. Lee J, Maunsell JHR. The Effect of Attention on Neuronal Responses to High and Low Contrast Stimuli. *Journal of*
831 *Neurophysiology*. 2010;104(2):960-71.
- 832 45. Compte A, Wang X-J. Tuning Curve Shift by Attention Modulation in Cortical Neurons: a Computational Study of its
833 Mechanisms. *Cerebral Cortex*. 2006;16(6):761-78.
- 834 46. Haenny PE, Schiller PH. State dependent activity in monkey visual cortex. *Experimental Brain Research*.
835 1988;69(2):225-44.
- 836 47. Doostani N, Hossein-Zadeh G-A, Cichy RM, Vaziri-Pashkam M. Attention Modulates Human Visual Responses to
837 Objects by Tuning Sharpening. 2023.
- 838 48. Ling S, Liu T, Carrasco M. How spatial and feature-based attention affect the gain and tuning of population responses.
839 *Vision Research*. 2009;49(10):1194-204.
- 840 49. Cohen MA, Konkle T, Rhee JY, Nakayama K, Alvarez GA. Processing multiple visual objects is limited by overlap in
841 neural channels. *Proceedings of the National Academy of Sciences*. 2014;111(24):8955-60.
- 842 50. Kiat JE, Luck SJ, Beckner AG, Hayes TR, Pomaranski KI, Henderson JM, et al. Linking patterns of infant eye movements
843 to a neural network model of the ventral stream using representational similarity analysis. *Developmental Science*. 2022;25(1).
- 844 51. Diedrichsen J, Khaligh-Razavi S-M, Kriegeskorte N. Deep Supervised, but Not Unsupervised, Models May Explain IT
845 Cortical Representation. *PLoS Computational Biology*. 2014;10(11).
- 846 52. Peelen MV, Kastner S. A neural basis for real-world visual search in human occipitotemporal cortex. *Proceedings of the*
847 *National Academy of Sciences*. 2011;108(29):12125-30.
- 848 53. Noah S, Powell T, Khodayari N, Olivan D, Ding M, Mangun GR. Neural Mechanisms of Attentional Control for Objects:
849 Decoding EEG Alpha When Anticipating Faces, Scenes, and Tools. *Journal of Neuroscience* 2020. p. 4913-24.
- 850 54. Folstein JR, Monfared SS, Maravel T. The effect of category learning on visual attention and visual representation.
851 *Psychophysiology*. 2017;54(12):1855-71.
- 852 55. Reynolds JH, Chelazzi L. Attentional Modulation of Visual Processing. *Annual Review of Neuroscience*.
853 2004;27(1):611-47.
- 854 56. Grabner RH, Neubauer AC, Stern E. Superior performance and neural efficiency: The impact of intelligence and
855 expertise. *Brain Research Bulletin*. 2006;69(4):422-39.
- 856 57. Maguire EA, Valentine ER, Wilding JM, Kapur N. Routes to remembering: the brains behind superior memory. *Nature*
857 *Neuroscience*. 2002;6(1):90-5.
- 858 58. Wolfe JM. Guided Search 2.0 A revised model of visual search. *Psychonomic Bulletin & Review*. 1994;1(2):202-38.
- 859 59. Williams M, Becker SI. Determinants of Dwell Time in Visual Search: Similarity or Perceptual Difficulty? *PLoS ONE*.
860 2011;6(3).
- 861 60. Hout MC, Goldinger SD. Target templates: the precision of mental representations affects attentional guidance and
862 decision-making in visual search. *Attention, Perception, & Psychophysics*. 2014;77(1):128-49.

- 863 61. Lee J, Geng JJ. Idiosyncratic Patterns of Representational Similarity in Prefrontal Cortex Predict Attentional
864 Performance. *The Journal of Neuroscience*. 2017;37(5):1257-68.
- 865 62. Sabine Kastner MAP, Peter De Weerd, Robert Desimone, and Leslie G. Ungerleider. Increased Activity in Human Visual
866 Cortex during Directed Attention in the Absence of Visual Stimulation. 1999.
- 867 63. Reynolds JH, Chelazzi L, Desimone R. Competitive Mechanisms Subserve Attention in Macaque Areas V2 and V4. *The*
868 *Journal of Neuroscience*. 1999;19(5):1736-53.
- 869 64. Bukach CM, Gauthier I, Tarr MJ. Beyond faces and modularity: the power of an expertise framework. *Trends in*
870 *Cognitive Sciences*. 2006;10(4):159-66.
- 871 65. van der Linden M, Wegman J, Fernández G. Task- and Experience-dependent Cortical Selectivity to Features Informative
872 for Categorization. *Journal of Cognitive Neuroscience*. 2014;26(2):319-33.
- 873 66. Reddy L, Kanwisher N. Category Selectivity in the Ventral Visual Pathway Confers Robustness to Clutter and Diverted
874 Attention. *Current Biology*. 2007;17(23):2067-72.
- 875 67. Zhang W, Luck SJ. Feature-based attention modulates feedforward visual processing. *Nature Neuroscience*.
876 2008;12(1):24-5.
- 877 *ADDIN ADDIN*