

Understanding Biological Visual Attention Using Convolutional Neural Networks

Grace W. Lindsay^{a,b}, Kenneth D. Miller^{a,b}

^a *Center for Theoretical Neuroscience, College of Physicians and Surgeons, Columbia University, New York, New York, USA*

^b *Mortimer B. Zuckerman Mind Brain Behavior Institute, College of Physicians and Surgeons, Columbia University, New York, New York, USA*

Abstract

Covert visual attention has been shown repeatedly to enhance performance on tasks involving the features and spatial locations to which it is deployed. Many neural correlates of covert attention have been found, but given the complexity of the visual system, connecting these neural effects to performance changes is challenging. Here, we use a deep convolutional neural network as a large-scale model of the visual system to test the effects of applying attention-like neural changes. Particularly, we explore variants of the feature similarity gain model (FSGM) of attention—which relates a cell’s tuning to its attentional modulation. We show that neural modulation of the type and magnitude observed experimentally can lead to performance changes of the type and magnitude observed experimentally. Furthermore, performance enhancements from attention occur for a diversity of tasks: high level object category detection and classification, low level orientation detection, and cross-modal color classification of an attended orientation. Utilizing the full observability of the model we also determine how activity should change to best enhance performance and how activity changes propagate through the network. Through this we find that, for attention applied at certain layers, modulating activity according to tuning performs as well as attentional modulations determined by backpropagation. At other layers, attention applied according to tuning does not successfully propagate through the network, and has a weaker impact on performance than attention determined by backpropagation. This thus highlights a discrepancy between neural tuning and function.

1. Introduction

1 Covert visual attention, applied according to spatial location or visual features, has
2 been shown repeatedly to enhance performance on challenging visual tasks [11]. To ex-
3 plore the neural mechanisms behind this enhancement, neural responses to the same
4 visual input are compared under different task conditions. Such experiments have
5 identified numerous neural modulations associated with attention, including changes
6 in firing rates, noise levels, and correlated activity [91, 15, 24, 57], however, the extent
7 to which these changes are responsible for behavioral effects is debated. Therefore,
8 theoretical work has been used to link sensory processing changes to performance
9 changes. While offering helpful insights, much of this work is either based on small,
10 hand-designed models [68, 79, 94, 12, 31, 100, 30] or lacks direct mechanistic inter-
11 pretability [99, 9, 90]. Here, we utilize a large-scale model of the ventral visual stream
12 to explore the extent to which neural changes like those observed in the biology can

13 lead to performance enhancements on realistic visual tasks. Specifically, we use a deep
14 convolutional neural network trained to perform object classification to test variants
15 of the feature similarity gain model of attention [92].

16 Deep convolutional neural networks (CNNs) are popular tools in the machine learn-
17 ing and computer vision communities for performing challenging visual tasks [75].
18 Their architecture—comprised of layers of convolutions, nonlinearities, and response
19 pooling—was designed to mimic the retinotopic and hierarchical nature of the mam-
20 malian visual system [75]. Models of a similar form have been used in neuroscience to
21 study the biological underpinnings of object recognition for decades [26, 76, 85]. Re-
22 cently it has been shown that when these networks are trained to successfully perform
23 object classification on real-world images, the intermediate representations learned are
24 remarkably similar to those of the primate visual system [102, 39, 38]. Specifically,
25 deep CNNs are state-of-the-art models for capturing the feedforward pass of the ven-
26 tral visual stream [40, 36, 10]. Many different studies have now built on this fact to
27 further compare the representations [93, 51, 44] and behavior [45, 27, 73, 77, 50] of
28 CNNs to that of biological vision. A key finding has been the correspondence between
29 different areas in the ventral stream and layers in the deep CNNs, with early convolu-
30 tional layers able to capture the representation of V1 and deeper layers relating to V4
31 and IT [29, 23, 83]. Given that CNNs reach near-human performance on visual tasks
32 and have architectural and representational similarities to the visual system, they are
33 particularly well-positioned for exploring how neural correlates of attention can impact
34 behavior.

35 We focus here on attention’s ability to impact activity levels (rather than noise or
36 correlations) as these findings are straightforward to implement in a CNN. Further-
37 more, by measuring the effects of firing rate manipulations alone, we make clear what
38 behavioral enhancements can plausibly be attributable to them.

39 One popular framework to describe attention’s effects on firing rates is the feature
40 similarity gain model (FSGM). This model, introduced by Treue & Martinez-Trujillo,
41 claims that a neuron’s activity is multiplicatively scaled up (or down) according to
42 how much it prefers (or doesn’t prefer) the properties of the attended stimulus [92,
43 56]. Attention to a certain visual attribute, such as a specific orientation or color,
44 is generally referred to as feature-based attention (FBA) and its effects are spatially
45 global: that is, if a task performed at one location in the visual field activates attention
46 to a particular feature, neurons that represent that feature across the visual field will
47 be affected [104, 81]. Overall, this leads to a general shift in the representation of the
48 neural population towards that of the attended stimulus [17, 35, 71]. Spatial attention
49 implies that a particular portion of the visual field is being attended. According to the
50 FSGM, spatial location is treated as an attribute like any other. Therefore, a neuron’s
51 modulation due to attention can be predicted by how well its preferred features and
52 spatial receptive field align with the features and location of the attended stimulus.
53 The effects of combined feature and spatial attention have been found to be additive
54 [33].

55 While the FSGM does describe many findings, its components are not uncontroversial.
56 For example, it is questioned whether attention impacts responses multiplicatively
57 or additively [6, 3, 52, 60], and whether or not the activity of cells that do not prefer
58 the attended stimulus is actually suppressed [7, 68]. Furthermore, only a handful of
59 studies have looked directly at the relationship between attentional modulation and
60 tuning [56, 80, 13, 97]. Another unsettled issue is where in the visual stream attention

61 effects can be seen. Many studies of attention focus on V4 and MT/MST [91], as
62 these areas have reliable attentional effects. Some studies do find effects at earlier
63 areas [66], though they tend to be weaker and occur later in the visual response [37].
64 Therefore, a leading hypothesis is that attention signals, coming from prefrontal areas
65 [65, 63, 4, 42], target later visual areas, and the feedback connections that those areas
66 send to earlier ones causes the weaker effects seen there later [8, 52].

67 In this study, we define the FSGM of attention mathematically and implement it
68 in a deep CNN. By testing different variants of the model, applied at different layers
69 in the network and for different tasks, we can determine the ability of these neural
70 changes to change behavior. Given the complexity of these large nonlinear networks,
71 the effects of something like FSGM are non-obvious. Because we have full access to all
72 units in the model, we can see how neural changes at one area propagate through the
73 network, causing changes at others. This provides a fuller picture of the relationship
74 between neural and performance correlates of attention.

75 2. Methods

76 2.1. Network Model

77 This work uses a deep convolutional neural network (CNN) as a model of the
78 ventral visual stream. Convolutional neural networks are feedforward artificial neural
79 networks that consist of a few basic operations repeated in sequence, key among
80 them being the convolution. The specific CNN architecture used in the study comes
81 from [86] (VGG-16D) and is shown in Figure 1A. A previous variant of this work used
82 a smaller network [48].

83 Here, the activity values of the units in each convolutional layer are the result of
84 applying a 2-D spatial convolution to the layer below, followed by positive rectification
85 (rectified linear 'ReLU' nonlinearity):

$$x_{ij}^{lk} = [(W^{lk} \star X^{l-1})_{ij}]_+ \quad (1)$$

86 where W^{lk} is the k^{th} convolutional filter at the l^{th} layer. The application of each filter
87 results in a 2-D feature map (the number of filters used varies across layers and is given
88 in parenthesis in Figure 1A). x_{ij}^{lk} is the activity of the unit at the i, j^{th} spatial location
89 in the k^{th} feature map at the l^{th} layer. X^{l-1} is thus the activity of all units at the
90 layer below the l^{th} layer. The input to the network is a 224 by 224 pixel RGB image,
91 and thus the first convolution is applied to these pixel values. For the purposes of this
92 study the convolutional layers are most relevant, and will be referred to according to
93 their numbering in Figure 1A.

94 Max pooling layers reduce the size of the feature maps by taking the maximum
95 activity value of units in a given feature map in non-overlapping 2x2 windows.

96 The final three layers of this network are each fully-connected to the layer below
97 them, with the number of units per layer given in parenthesis in Figure 1A. Therefore,
98 connections exist from all units from all feature maps in the last convolutional layer
99 (layer 13) to all 4096 units of the next layer, and so on. This network was pre-trained
100 [25] using backpropagation on the ImageNet classification task, which involves doing
101 1000-way object categorization (for details see [86]). The final layer of the network
102 thus contains 1000 units upon which a softmax classifier is used to output a ranked
103 list of category labels for a given image. Looking at the top-5 error rate (wherein an

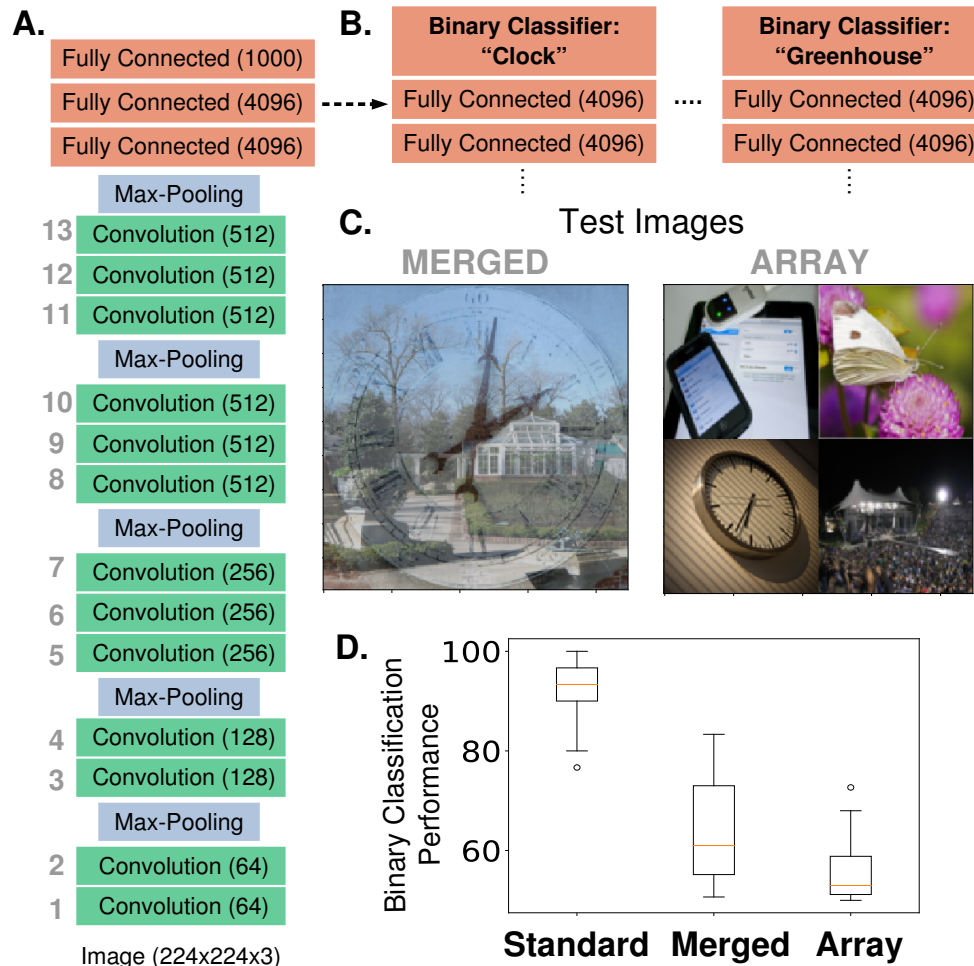


Figure 1: Network Architecture and Feature-Based Attention Task Setup. A.) The model used is a pre-trained deep neural network (VGG-16) that contains 13 convolutional layers (labeled in gray, number of feature maps given in parenthesis) and is pre-trained on the ImageNet dataset to do 1000-way object classification. All convolutional filters are 3x3. B.) Modified architecture for feature-based attention tasks. To perform our feature-based attention tasks, the final layer that was implementing 1000-way softmax classification is replaced by binary classifiers (logistic regression), one for each category tested (2 shown here). These binary classifiers are trained on standard ImageNet images. C.) Test images for feature-based attention tasks. Merged images (left) contain two transparently overlaid ImageNet images of different categories. Array images (right) contain four ImageNet images on a 2x2 grid. Both are 224 x 224 pixels. These images are fed into the network and the binary classifiers are used to label the presence or absence of the given category. D.) Performance of binary classifiers. Box plots describe values over 20 different object categories (median marked in red, box indicates lower to upper quartile values and whiskers extend to full range with outliers marked as dots). Standard images are regular ImageNet images not used in the binary classifier training set.

104 image is correctly labeled if the true category appears in the top five categories given
105 by the network), this network achieves 92.7% accuracy.

106 *2.2. Object Category Attention Tasks*

107 The tasks we use to probe the effects of feature-based attention in this network
108 involve determining if a given object category is present in an image or not, similar to
109 tasks used in [88, 72, 41]. To have the network perform this specific task, we replaced
110 the final layer in the network with a series of binary classifiers, one for each category
111 tested (Figure 1B). We tested a total of 20 categories: paintbrush, wall clock, seashore,
112 paddlewheel, padlock, garden spider, long-horned beetle, cabbage butterfly, toaster,
113 greenhouse, bakery, stone wall, artichoke, modem, football helmet, stage, mortar,
114 consomme, dough, bathtub. Binary classifiers were trained using ImageNet images
115 taken from the 2014 validation set (and were therefore not used in the training of
116 the original model). A total of 35 unique true positive images were used for training
117 for each category, and each training batch was balanced with 35 true negative images
118 taken from the remaining 19 categories. The results shown here come from using
119 logistic regression as the binary classifier, though trends in performance are similar if
120 support vector machines are used. Experimental results suggest that classifiers trained
121 on unattended and isolated object images are appropriate for reading out attended
122 objects in cluttered images [105].

123 Once these binary classifiers are trained, they are then used to classify more chal-
124 lenging test images. These test images are composed of multiple individual images
125 (drawn from the 20 categories) and are of two types: "merged" and "array". Merged
126 images are generated by transparently overlaying two images, each from a different
127 category (specifically, pixel values from each are divided by two and then summed).
128 Array images are composed of four separate images (all from different categories) that
129 are scaled down to 112 by 112 pixels and placed on a two by two grid. The images that
130 comprise these test images also come from the 2014 validation set, but are separate
131 from those used to train the binary classifiers. See examples of each in Figure 1C. Test
132 image sets are balanced (50% do contain the given category and 50% do not, 150 total
133 test images per category). Both true positive and true negative rates are recorded and
134 overall performance is the average of these rates.

135 To test the effects of spatial attention, only the "array" images are used. The task is
136 to identify the category of the object at the attended location. Therefore, performance
137 is measured using the original 1000-way classifier, with the category of the image in
138 the attended quadrant as the true label (200 images were tested per quadrant).

139 *2.3. Object Category Gradient Calculations*

140 When neural networks are trained via backpropagation, gradients are calculated
141 that indicate how a given weight in the network impacts the final classification. We
142 use this same method to determine how a given unit's activity impacts the final clas-
143 sification. Specifically, we input a "merged" image (wherein one of the images belongs
144 to the category of interest) to the network. We then use gradient calculations to deter-
145 mine the changes in activity that would move the 1000-way classifier toward classifying
146 that image as belonging to the category of interest (i.e. rank that category highest).
147 We average these activity changes over images and over all units in a feature map.
148 This gives a single value per feature map:

$$g_c^{lk} = -\frac{1}{N_c} \sum_{n=1}^{N_c} \frac{1}{HW} \sum_{i=1, j=i}^{H, W} \frac{\partial E(n)}{\partial x_{ij}^{lk}(n)} \quad (2)$$

149 where H and W are the spatial dimensions of layer l and N_c is the total number of
150 images from the category (here $N_c = 35$, and the merged images used were generated
151 from the same images used to generate tuning curves, described below). $E(n)$ is the
152 error of the classifier in response to image n , which is defined as the difference between
153 the activity vector of the final layer (after the soft-max operation) and a one-hot
154 vector, wherein the correct label is the only non-zero entry. Because we are interested
155 in activity changes that would decrease the error value, we negate this term. The
156 gradient value we end up with thus indicates how the feature map's activity would
157 need to change to make the network more likely to classify an image as the desired
158 category. Repeating this procedure for each category, we obtain a set of gradient
159 values (one for each category, akin to a tuning curve), for each feature map: \mathbf{g}^{lk} . Note
160 that, as these values result from applying the chain rule through layers of the network,
161 they can be very small, especially for the earliest layers. For this study, the sign and
162 relative magnitudes are of more interest than the absolute values.

163 *2.4. Oriented Grating Attention Tasks*

164 In addition to attending to object categories, we also test attention on simpler
165 stimuli. In the orientation detection task, the network detects the presence of a given
166 orientation in an image. Again, the final layer of the network is replaced by a series
167 of binary classifiers, one for each of 9 orientations (0, 20, 40, 60, 80, 100, 120, 140,
168 and 160 degrees. Gratings had a frequency of .025 cycles/pixel). The training sets
169 for each were balanced (50% had only the given orientation and 50% had one of 8
170 other orientations) and composed of full field (224 by 224 pixel) oriented gratings of
171 various colors (to increase the diversity of the training images, they were randomly
172 degraded by setting blocks of pixels ranging uniformly from 0% to 70% of the image
173 to 0 at random). Test images were each composed of two oriented gratings of different
174 orientation and color (color options: red, blue, green, orange, purple). Each of these
175 gratings were of size 112 by 112 pixels and placed randomly in a quadrant while the
176 remaining two quadrants were black (Figure 6A). Again, the test sets were balanced
177 and performance was measured as the average of the true positive and true negative
178 rates (100 test images per orientation).

179 These same test images were used for a cross-modal attention task wherein the
180 network had to classify the color of the grating that had the attended orientation. For
181 this, the final layer of the network was replaced with a 5-way softmax color classifier.
182 This color classifier was trained using the same full field oriented gratings used to train
183 the binary classifiers (therefore, the network saw each color at all orientation values).
184 The test sets contained images that all had the attended orientation as one of the two
185 gratings (125 images per orientation). Performance was measured as the percent of
186 trials wherein the color classifier correctly ranked the color of the attended grating
187 highest (top-1 error).

188 Finally, for one analysis, a joint feature and spatial attention task was used. This
189 task is almost identical to the setup of the orientation detection task, except that the
190 searched-for orientation would only appear in one of the four quadrants. Therefore,
191 performance could be measured when applying feature attention to the searched-for
192 orientation, spatial attention to the quadrant in which it could appear, or both.

193 2.5. How Attention is Applied

194 This study aims to test variations of the feature similarity gain model of attention,
195 wherein neural activity is modulated by attention according to how much the neuron
196 prefers the attended stimulus. To replicate this in our model, we therefore must first
197 determine the extent to which units in the network prefer different stimuli ("tuning
198 values"). When attention is applied to a given category, for example, units' activities
199 are modulated according to these values. We discuss below the options for how exactly
200 to implement that modulation.

201 2.5.1. Tuning Values

202 To determine tuning to the 20 object categories used, we presented the network
203 with images of each object category (the same images on which the binary classifiers
204 were trained) and measured the relative activity levels.

205 Specifically, for the k^{th} feature map in the l^{th} layer, we define $r^{lk}(n)$ as the activity in
206 response to image n , averaged over all units in the feature map (i.e., over the spatial
207 dimensions). Averaging these values over all images in the training sets ($N_c = 35$
208 images per category, 20 categories. $N=700$) gives the mean activity of the feature map
209 \bar{r}^{lk} :

$$\bar{r}^{lk} = \frac{1}{N} \sum_{n=1}^N r^{lk}(n) \quad (3)$$

210 Tuning values are defined for each object category, c as:

$$f_c^{lk} = \frac{\frac{1}{N_c} \sum_{n \in c} r^{lk}(n) - \bar{r}^{lk}}{\sqrt{\frac{1}{N} \sum_{i=1}^N (r^{lk}(n) - \bar{r}^{lk})^2}} \quad (4)$$

211 That is, a feature map's tuning value for a given category is merely the average
212 activity of that feature map in response to images of that category, with the mean
213 activity under all image categories subtracted and standard deviation divided. These
214 tuning values determine how the feature map is modulated when attention is applied
215 to the category. Taking these values as a vector over all categories, \mathbf{f}_{lk} , gives a tuning
216 curve for the feature map. We define the overall tuning quality of a feature map as
217 its maximum absolute tuning value: $\max(|\mathbf{f}_{lk}|)$. To determine expected tuning quality
218 by chance, we shuffled the responses to individual images across category and feature
219 map at a given layer and calculated tuning quality for this shuffled data.

220 We define the category with the highest tuning value as that feature map's most
221 preferred, and the category with the lowest (most negative) value as the least or anti-
222 preferred.

223 We apply the same procedure to generate tuning curves for orientation and for
224 color by using the full field gratings used to train the orientation detection and color
225 classification classifiers. The orientation tuning values were used when applying at-
226 tention in these tasks. The color tuning curves were generated only to measure color
227 tuning and its quality in the network.

228 When measuring how correlated tuning values are with gradient values, shuffled
229 comparisons are used. To do this shuffling, correlation coefficients are calculated from
230 pairing each feature map's tuning values with a random other feature map's gradient
231 values.

232 2.5.2. Gradient Values

233 In addition to applying attention according to tuning, we also attempt to generate
234 the "best possible" attentional modulation by utilizing gradient values. These gradient
235 values are calculated slightly differently from those described above (2.3), because they
236 are meant to represent how feature map activity should change in order to increase
237 overall task performance, rather than just increase the chance of classifying an image
238 as a certain object or orientation.

239 The error functions used to calculate gradient values for the category and orienta-
240 tion detection tasks were for the binary classifiers associated with each object/orientation.
241 A balanced set of test images was used. Therefore a feature map's gradient value for
242 a given object/orientation is the averaged activity change that would increase binary
243 classification performance for that object/orientation. Note that on images that the
244 network already classifies correctly, gradients are zero. Therefore, the gradient values
245 are driven by the errors: false negatives (classifying an image as not containing the
246 category when it does) and false positives (classifying an image as containing the cat-
247 egory when it does not). In our detection tasks, the former error is more prevalent
248 than the latter, and thus is the dominant impact on the gradient values.

249 The same procedure was used to generate gradient values for the color classification
250 task. Here, gradients were calculated using the 5-way color classifier: for a given
251 orientation, the color of that orientation in the test image was used as the correct label,
252 and gradients were calculated that would lead to the network correctly classifying the
253 color. Averaging over many images of different colors gives one value per orientation
254 that represents how a feature map's activity should change in order to make the
255 network better at classifying the color of that orientation.

256 In both of the orientation tasks, the test images used for gradient calculations
257 (50 images per orientation) differed from those used to assess performance. For the
258 object detection task, images used for gradient calculations were merged images (45
259 per category) drawn from the same pool as, but different from, those used to test
260 detection performance.

261 2.5.3. Spatial Attention

262 In the feature similarity gain model of attention, attention is applied according
263 to how much a cell prefers the attended feature, and location is considered a feature
264 like any other. In CNNs, each feature map results from applying the same filter at
265 different spatial locations. Therefore, the 2-D position of a unit in a feature map
266 represents more or less the spatial location to which that unit responds. Via the max-
267 pooling layers, the size of each feature map shrinks deeper in the network, and each
268 unit responds to a larger area of image space, but the "retinotopy" is still preserved.
269 Thus, when we apply spatial attention to a given area of the image, we enhance the
270 activity of units in that area of the feature maps (and, as we discuss below, possibly
271 decrease the activity of units in other areas). In this study, spatial attention is tested
272 using array images, and thus attention is applied to a given quadrant of the image.

273 2.5.4. Implementation Options

274 The values discussed above determine how strongly different feature maps or units
275 should be modulated under different attentional conditions. We will now lay out the
276 different implementation options for that modulation.

277 First, the modulation can be multiplicative or additive. That is, when attending
278 to category c , the slope of the rectified linear units can be multiplied by a weighted

279 function of the tuning value for category c :

$$x_{ij}^{lk} = (1 + \beta f_c^{lk})[(I_{lk}^{ij})]_+ \quad (5)$$

280 with I_{lk}^{ij} representing input to the unit coming from layer $l - 1$. Alternatively, a
281 weighted version of the tuning value can be added before the rectified linear unit:

$$x_{ij}^{lk} = [I_{ij}^{lk} + \mu_l \beta f_c^{lk}]_+ \quad (6)$$

282 Strength of attention is varied via the weighting parameter, β . For the additive effect,
283 manipulations are multiplied by μ_l , the average activity level across all units of layer
284 l in response to all images (for each of the 13 layers respectively: 20, 100, 150, 150,
285 240, 240, 150, 150, 80, 20, 20, 10, 1). When gradient values are used in place of tuning
286 values, we normalize them by the maximum value at a layer, to be the same order of
287 magnitude as the tuning values: $\mathbf{g}^l / \max(|\mathbf{g}^l|)$.

288 Note that for feature-based attention all units in a feature map are modulated the
289 same way, as feature attention has been found to be spatially global. In the case of
290 spatial attention, object category tuning values are not used. Rather, the tuning value
291 term is set to +1 if the i, j position of the unit is in the attended quadrant and to -1
292 otherwise. For feature attention tasks, β ranged from 0 to a maximum of 11.85 (object
293 attention) and 0 to 4.8 (orientation attention). For spatial attention tasks, it ranged
294 from 0 to 2.

295 Next, we chose whether attention only enhances units that prefer the attended
296 feature/location, or also decreases activity of those that don't prefer it. For the latter,
297 the tuning values are used as-is. For the former, the tuning values are positively-
298 rectified: $[\mathbf{f}^{lk}]_+$.

299 Combining these two factors, there are four implementation options: additive
300 positive-only, multiplicative positive-only, additive bidirectional, and multiplicative
301 bidirectional.

302 The final option is the layer in the network at which attention is applied. We try
303 attention at all convolutional layers individually and simultaneously (when applying
304 simultaneously the strength range tested is a tenth of that when applying to a single
305 layer).

306 Note that when gradient values were used, only results from using multiplicative
307 bidirectional effects are reported (when tested on object category detection, multi-
308 plicative effects performed better than additive when using gradient values).

309 2.6. Signal Detection Calculations

310 For the joint spatial-feature attention task, we calculated criteria (c , "threshold")
311 and sensitivity (d') using true (TP) and false (FP) positive rates as follows [53]:

$$c = -.5(\Phi^{-1}(TP) + \Phi^{-1}(FP)) \quad (7)$$

312 where Φ^{-1} is the inverse cumulative normal distribution function. c is a measure of
313 the distance from a neutral threshold situated between the mean of the true negative
314 and true positive distributions. Thus, a positive c indicates a stricter threshold (fewer
315 inputs classified as positive) and a negative c indicates a more lenient threshold (more

316 inputs classified as positive).

$$d' = \Phi^{-1}(TP) - \Phi^{-1}(FP) \quad (8)$$

317 This measures the distance between the means of the distributions for true negative
318 and two positives. Thus, a larger d' indicates better sensitivity.

319 When necessary, a correction was applied wherein false positive rates of 0 were set
320 to .01 and true positive rates of 1 were set to .99.

321 *2.7. "Recording" Procedures*

322 We examined the effects that applying attention at certain layers in the network
323 (specifically 2, 6, 8, 10, and 12) has on activity of units at other layers. We do this for
324 many different circumstances, using multiplicative bidirectional attention with $\beta = .5$
325 unless otherwise stated.

326 *2.7.1. Unimodal Task Recording Setup*

327 This recording setup is designed to mimic the analysis of [56]. Here, the images
328 presented to the network are full-field oriented gratings of all orientation-color combi-
329 nations. Feature map activity is measured as the spatially averaged activity of all units
330 in a feature map in response to an image. Activity in response to a given orientation
331 is further averaged over all colors. Each feature map's preferred (most positive tuning
332 value) and anti-preferred (most negative tuning value) orientations are determined.
333 Activity is recorded when attention is applied to the preferred or anti-preferred orien-
334 tation and activity ratios are calculated. According to the FSGM, the ratio of activity
335 when the preferred orientation is attended over when the anti-preferred is attended
336 should be greater than one and the same regardless of whether the image is of the pre-
337 ferred or anti-preferred orientation. According to the feature matching (FM) model,
338 the ratio of the activity when attending the presented orientation over attending an
339 absent orientation should be greater than one and similar regardless of whether the
340 orientation is preferred or not. We measure all of these ratios, and the fraction of total
341 feature maps which show FM behavior, when attention is applied according to tuning
342 values or gradient values.

343 As in [56], we also look at a measure of activity changes across all orientations.
344 We calculate the ratio of activity when attention is applied to a given orientation
345 (and the orientation is present in the image) over activity in response to the same
346 image when no attention is applied. These ratios are then organized according to
347 orientation preference: the most preferred is at location 0, then the average of next
348 two most preferred at location 1, and so on with the average of the two least preferred
349 orientations at location 4 (the reason for averaging of pairs is to match [56] as closely
350 as possible). Fitting a line to these points gives a slope and intercept for each feature
351 map. FSGM predicts a negative slope and an intercept greater than one.

352 We also calculate the same activity ratios described above when the images pre-
353 sented are standard (single image) ImageNet images from each of the 20 categories
354 (activity is averaged over 5 images per category). Attention is applied according to
355 object category tuning values or to gradient values for binary classification as described
356 in 2.5.2.

357 *2.7.2. Cross-modal Task Recording Setup*

358 Cross-modal tasks involve attending to one modality (here, space or orientation)
359 and reading out another (category or color, respectively). Specifically, in the first task,
360 activity is recorded when spatial attention is applied to a given quadrant. Here, the
361 activity for each feature map is averaged only over units in the quadrant that matches
362 the attended quadrant. The images used are array images with 6 examples of each
363 object category in the attended quadrant (for a total of 120 images). Activity ratios are
364 calculated as the activity when the recorded quadrant is attended over activity when
365 no attention is applied. The average ratio for each category is organized according to
366 category preference for each feature map and a line is fit to these points. The intercept
367 (measured here as the true intercept minus one) and difference (slope multiplied by
368 the number of categories minus one, 19) are calculated for each feature map. FSGM
369 predicts a positive intercept and zero slope, because responses to all categories should
370 be scaled equally by spatial attention.

371 The second cross-modal task setup involves measuring color encoding in different
372 attention conditions. Here, images similar to those used in the orientation detection
373 and color classification tasks are used. Specifically, images are generated that have two
374 oriented gratings in two of the four quadrants. One is oriented at 160 degrees and the
375 other nearly orthogonal at 80. All pairs of colors are generated for the two gratings
376 (thus the two gratings may have the same color, which is a difference from the stimuli
377 used in the orientation tasks). Activity is organized according to the color of the 160
378 degree grating (and averaged over the colors of the 80 degree grating), in order from
379 most to least preferred color for each feature map. Lines were fit to these points in
380 two cases: when attention was directed to 80 degrees and when it was directed to 160
381 degrees. We then asked if attention to 160 degrees led to better encoding of the color of
382 the 160 degree stimulus compared to attention to 80 degrees. We considered a feature
383 map to have better color encoding of the 160 degree grating if its mean increased (a
384 stronger overall signal, measured as the activity value at the middle of the line) and
385 if its slope became more negative (stronger differentiation between colors). Results
386 are similar if only the latter condition is used. We measure the encoding changes for
387 two separate populations of feature maps: those that prefer 160 degrees and those
388 that anti-prefer it (most negative tuning value). Stimuli at 160 degrees were chosen as
389 the focus of this analysis because across all layers there are roughly equal numbers of
390 feature maps that prefer and anti-prefer it. Percent of feature maps that have better
391 encoding were measured when attention was applied according to orientation tuning
392 values or color classification gradient values.

393 In all cases, lines are fit using the least squares method, and any activity ratios
394 with zero in the denominator were discarded.

395 *2.8. Experimental Data*

396 Model results were compared to previously published data coming from several
397 studies. In [55], a category detection task was performed using stereogram stimuli
398 (on object present trials, the object image was presented to one eye and a noise mask
399 to another). The presentation of the visual stimuli was preceded by a verbal cue
400 that indicated the object category that would later be queried (cued trials) or by
401 meaningless noise (uncued trials). After visual stimulus presentation, subjects were
402 asked if an object was present and, if so, if the object was from the cued category
403 (categories were randomized for uncued trials). In Experiment 1, the object images

404 were line drawings (one per category) and the stimuli were presented for 1.5 sec. In
405 Experiment 2, the object images were grayscale photographs (multiple per category)
406 and presented for 6 sec. True positives were counted as trials wherein a given object
407 category was present and the subject correctly indicated its presence when queried.
408 False positives were trials wherein no category was present and subjects indicated that
409 the queried category was present.

410 In [54], a similar detection task is used. Here, subjects detect the presence of an
411 uppercase letter that is (on target present trials) presented rapidly and followed by
412 a mask. Prior to the visual stimulus, a visual or audio cue indicated a target letter.
413 After the visual stimulus, the subjects were required to indicate whether any letter
414 was present. True positives were trials in which a letter was present and the subject
415 indicated it (only uncued trials or validly cued trials—where the cued letter was the
416 letter shown—were considered here). False positives were trials where no letter was
417 present and the subject indicated that one was.

418 The task in [41] is also an object category detection task. Here, an array of several
419 images was flashed on the screen with one image marked as the target. All images
420 were color photographs of objects in natural scenes. In certain blocks, the subjects
421 knew in advance which category they would later be queried about (cued trials). On
422 other trials, the queried category was only revealed after the visual stimulus (uncued).
423 True positives were trials in which the subject indicated the presence of the queried
424 category when it did exist in the target image. False positives were trials in which
425 the subject indicated the presence of the cued category when it was not in the target
426 image. Data from trials using basic category levels with masks were used for this
427 study.

428 Finally, we include one study using macaques wherein both neural and performance
429 changes were measured [58]. In this task, subjects had to report a change in orientation
430 that could occur in one of two stimuli. On cued trials, the change occurred in the cued
431 stimulus in 80% of trials and the uncued stimulus in 20% of trials. On neutrally-cued
432 trials, subjects were not given prior information about where the change was likely
433 to occur (50% at each stimulus). Therefore performance could be compared under
434 conditions of low (uncued stimuli), medium (neutrally cued stimuli), and high (cued
435 stimuli) attention strength. Correct detection of an orientation change in a given
436 stimulus (indicated by a saccade) is considered a true positive and a saccade to the
437 stimulus prior to any orientation change is considered a false positive. True negatives
438 are defined as correct detection of a change in the uncued stimulus (as this means the
439 subject correctly did not perceive a change in the stimulus under consideration) and
440 false negatives correspond to a lack of response to an orientation change.

441 In cases where the true and false positive rates were not published, they were
442 obtained via personal communications with the authors.

443 **3. Results**

444 The ability to manipulate activities as well as measure performance on complicated
445 visual tasks make CNNs a great testing ground for theories of attention. CNNs trained
446 on visual object recognition learn representations that are similar to those of the
447 ventral stream. The network used in this study was explored in [29], where it was
448 shown that early convolutional layers of this CNN are best at predicting activity of
449 voxels in V1, while late convolutional layers are best at predicting activity of voxels in

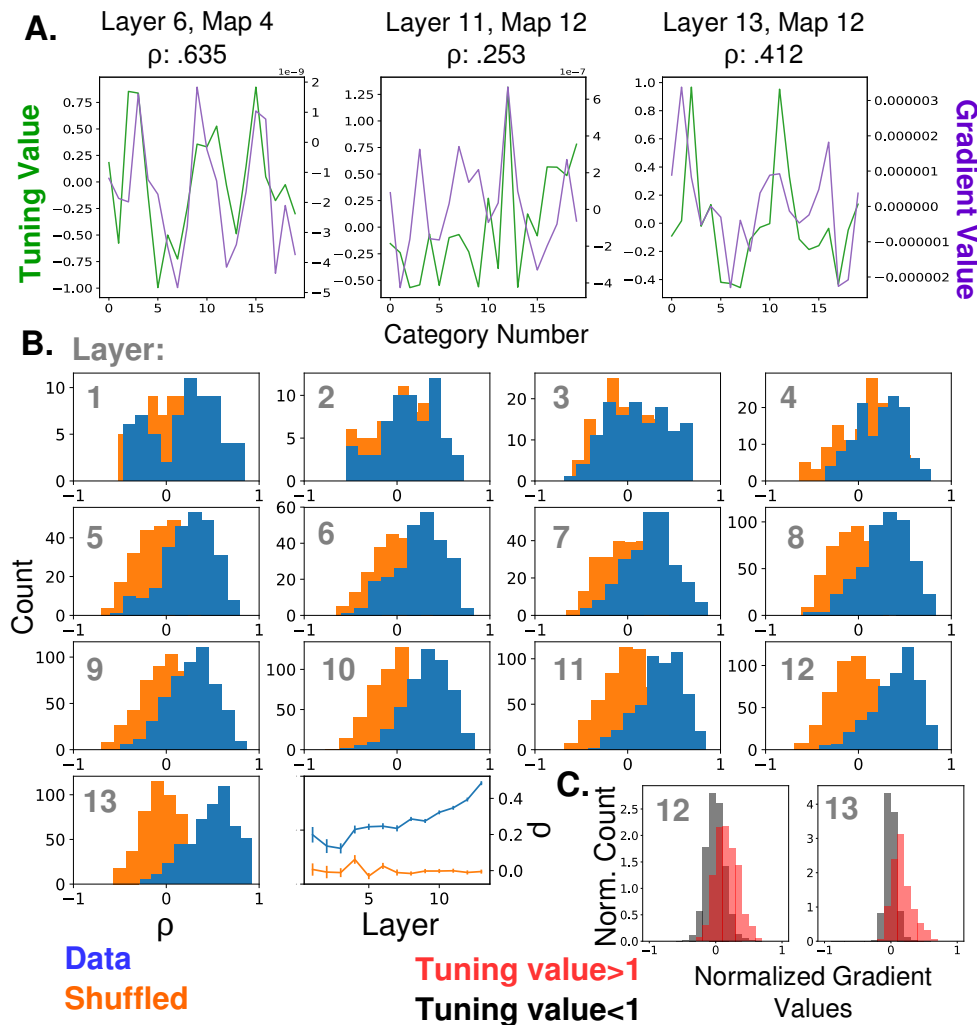


Figure 2: Relationship Between Feature Map Tuning and Gradients. A.) Example tuning values (green, left axis) and gradient values (purple, right axis) of three different feature maps from three different layers (identified in titles, layers as labeled in Fig 1A) over the 20 tested object categories. Correlation coefficients between tuning curves and gradient values given in titles. B.) Histograms of correlation coefficients across all feature maps at each layer (blue) along with shuffled comparisons (orange). Final subplot shows average correlation coefficients across layers (errorbars +/- S.E.M.). C.) Distributions of gradient values when tuning is strong. In red, histogram of gradient values associated with tuning values larger than one, across all feature maps in layer 12 (left) and 13 (right). For comparison, histograms of gradient values associated with tuning values less than one are shown in black (counts are separately normalized for visibility, as the population in black is much larger than that in red).

450 the object-selective lateral occipital area (LO). In addition, CNN architecture makes
451 comparison to biological vision straightforward. For example, the application of a
452 given convolutional filter results in a feature map, which is a 2-D grid of artificial
453 neurons that represent how well the bottom-up input aligns with the filter at each
454 location. Therefore a "retinotopic" layout is built into the structure of the network,
455 and the same visual features are represented across that retinotopy (akin to how cells
456 that prefer different orientations exist at all locations across the V1 retinotopy). We
457 utilize these properties to test variants of the feature similarity gain model (FSGM)
458 on a diverse set of visual tasks that are challenging for the network. We also take
459 advantage of the full observability of this network model to compare the FSGM to
460 "optimal" attentional manipulation, as determined by backpropagation calculations.

461 *3.1. The Relationship between Tuning and Classification*

462 The feature similarity gain model of attention posits that neural activity is modu-
463 lated by attention in proportion to how strongly a neuron prefers the attended features,
464 as assessed by its tuning. However, the relationship between a neuron's tuning and its
465 ability to influence downstream readouts remains a difficult one to investigate biolog-
466 ically. We use our hierarchical model to explore this question directly. We do so by
467 calculating gradient values, which we compare to tuning curves (see Methods Sections
468 2.3 and 2.5.1 for details). These gradient values indicate the way in which activity of a
469 feature map should change in order to make the network more likely to classify an im-
470 age as being of a certain object category. If there is a correspondence between tuning
471 and classification, a feature map that prefers a given object category (that is, responds
472 strongly to it compared to other categories) should also have a high positive gradient
473 value for that category. In Figure 2A we show gradient values and tuning curves for
474 three example feature maps. In Figure 2B, we show the distribution of correlation co-
475 efficients between tuning values and gradient values for all feature maps at each of the
476 13 convolutional layers. As can be seen in the final subplot, on average, tuning curves
477 show higher than expected correlation with gradient values at all layers (compared to
478 shuffled controls). Furthermore, this correlation increases with later layers. While the
479 correlation between tuning and gradient values suggests that a feature map's response
480 is indicative of its functional role, the correspondence is not perfect. In Figure 2C,
481 we show the gradient values of feature maps at layers 12 and 13, segregated according
482 to tuning value. In red are gradient values that correspond to tuning values greater
483 than one (for example, category 12 for the feature map in the middle pane of Figure
484 2A). As these distributions show, strong tuning values can be associated with weak or
485 even negative gradient values. Negative gradient values indicate that increasing the
486 activity of that feature map makes the network less likely to categorize the image as
487 the given category. Therefore, even feature maps that strongly prefer a category (and
488 are only a few layers from the classifier) still may not be involved in its classification,
489 or even be inversely related to it.

490 *3.2. Feature-based Attention Improves Performance on Challenging Object Classifica-* 491 *tion Tasks*

492 To determine if manipulation according to tuning values can enhance performance,
493 we created challenging visual images composed of multiple objects for the network to
494 classify. These test images are of two types: merged (two object images transparently
495 overlaid, such as in [84]) or array (four object images arranged on a grid) (see Figure

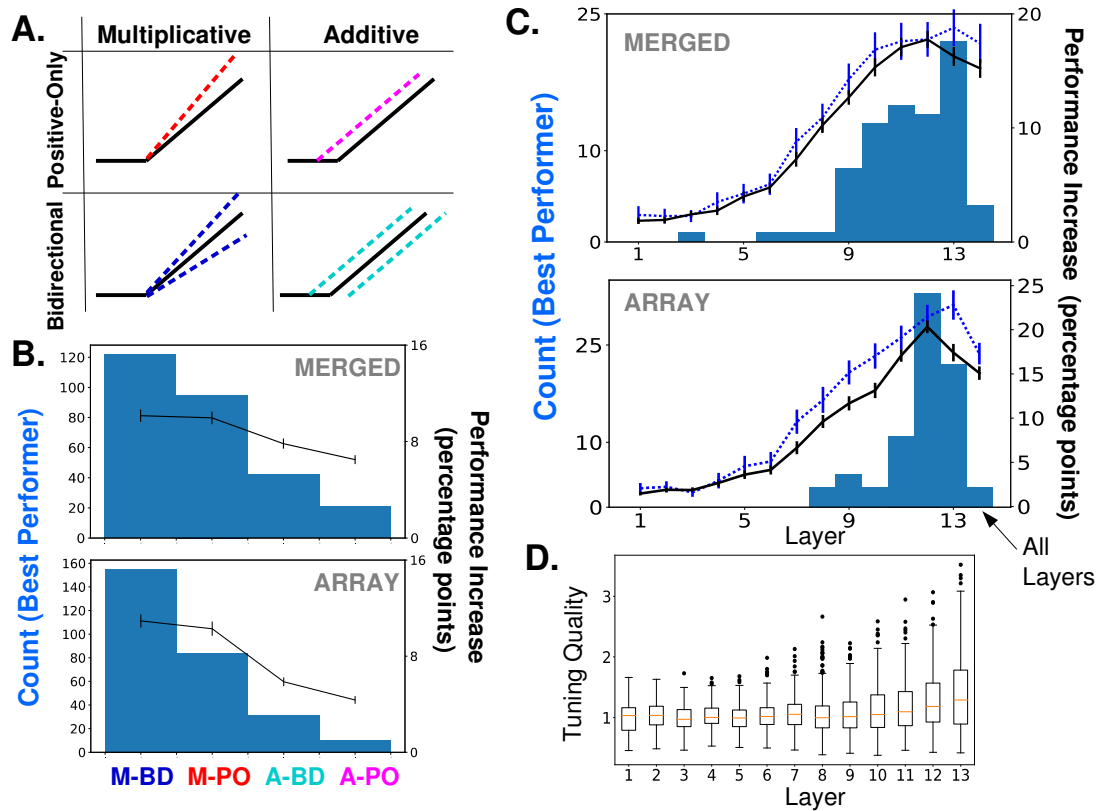


Figure 3: Effects of Applying Feature-Based Attention on Object Category Tasks. A.) Schematics of how attention can modulate the activity function. Feature-based attention modulates feature maps according to their tuning values but this modulation can scale the activity multiplicatively or additively, and can either only enhance feature maps that prefer the attended category (positive-only) or also decrease the activity of feature maps that do not prefer it (bidirectional). B.) Considering the combination of attention applied to a given category at a given layer as an instance (20 categories * 14 layer options = 280 instances), histograms (left axis) show how often the given option is the best performing, for merged (top) and array (bottom) images. Average increase in binary classification performance for each option also shown (right axis, averaged across all instances, errorbars +/- S.E.M.) C.) Comparison of performance effects of layer options. Considering each instance as the combination of attention applied to a given category using a given implementation option (20 categories * 4 implementation options = 80 instances), histograms show how often applying attention to the given layer is the best performing, for merged (top) and array (bottom) images. The final column corresponds to attention applied to all layers simultaneously with the same strength (strengths tested are one-tenth of those when strength applied to individual layers). Average increase in binary classification performance for each layer also shown in black (right axis, errorbars +/- S.E.M.). Average performance increase for MBD option only shown in blue. In all cases, best performing strength from the range tested is used for each instance. D.) Tuning quality across layers. Tuning quality is defined per feature map as the maximum absolute tuning value of that feature map. Box plots show distribution across feature maps for each layer. Average tuning quality for shuffled data: $.372 \pm .097$ (this value does not vary significantly across layers)

496 1C for an example of each). The task for the network is to detect the presence or
497 absence of a given object category in these images. It does so using a series of binary
498 classifiers trained on standard images of these objects, which replace the last layer
499 of the network (Figure 1B). The performance of these classifiers on the test images
500 indicates that this is a challenging task for the network (Figure 1D), and thus a good
501 opportunity to see the effects of attention. Without attention, the average performance
502 of the binary classifiers across all categories is 64.4% on merged images and 55.6%
503 on array (compared to a chance performance of 50%, as the test sets contained the
504 attended category 50% of the time).

505 We implement feature-based attention in this network by modulating the activity
506 of feature maps according to how strongly they prefer the attended object category
507 (see Methods 2.5.1). While tuning values determine the relative strength and direction
508 of the modulation, there are still options regarding how to implement those changes.
509 We test additive effects (wherein attention alters the activity of a feature map by
510 the same amount regardless of its activity level) and multiplicative effects (attention
511 changes the slope of the activity function). We also consider the situation where
512 attention only increases the activity of feature maps that prefer the attended category
513 (i.e., have a positive tuning value), or when attention also decreases the activity of
514 feature maps that do not prefer the attended category. Taken together this leads
515 to four implementation options: additive positive-only, multiplicative positive-only,
516 additive bidirectional, and multiplicative bidirectional (see Figure 3A for depictions of
517 each, and Methods 2.5.4 for details). A final option is the choice of convolutional layer
518 at which these manipulations are applied.

519 To determine which of these attention mechanisms is best, attention is applied
520 to each object category and the performance of the binary classifier associated with
521 that category is compared with and without the different activity manipulations. The
522 results of this are shown in Figure 3B and C (the best performing strength, including
523 0 if necessary, is assumed for each category. See Methods for details).

524 As Figure 3B shows, multiplicative bi-directional effects are best able to enhance
525 performance, measured in terms of the number of times that the multiplicative bidirec-
526 tional option beats out the other three options when compared for the same category
527 and layer (blue histogram). The second best option is multiplicative positive-only,
528 then additive bidirectional, and additive positive-only. This ordering is the same when
529 looking at the average increase in performance (black line), however, the differences
530 between multiplicative bi-directional and multiplicative positive-only performance are
531 not significant. Furthermore, these trends are identical regardless of whether tested
532 on merged (top) or array (bottom) images, though the differences are starker for array
533 images.

534 Figure 3C shows a similar analysis but across layers at which attention is applied.
535 Again, the trends are the same for merged and array images and show a clear increase
536 in performance as attention is applied at later layers in the network (numbering is as
537 in 1A). Across all implementation options, attention at layer 12 best increases average
538 performance (black lines). However this is driven by the additive implementations.
539 We show the average performance increase with layer for multiplicative bi-directional
540 effects alone (blue dotted line). For this the final layer is best, leading to an 18.8%
541 percentage point increase in binary classification on the merged image task and 22.8%
542 increase on the array task.

543 The trends in performance track trends in tuning quality shown in 3D. That is,

544 layers with better object category tuning lead to better performance when attention is
545 applied at them. They also track the correlation between tuning values and gradient
546 values, as that correlation increases with later layers.

547 Overall, the best performing options for implementing attention—multiplicative
548 bidirectional effects applied at later layers—are in line with what has been observed
549 biologically and described by the feature similarity gain model [92, 57].

550 *3.3. Strength of Attention Influences True and False Positive Tradeoff*

551 As mentioned above, strength is a relevant variable when implementing attention.
552 Specifically, the strength parameter, which we call β , scales the tuning values to deter-
553 mine how strongly attention modulates activities (in the case of additive effects, this
554 value is further multiplied by the average activity level of the layer before being added
555 to the response). We tested a range of β values and the analysis in Figure 3 assumes
556 the best-performing β for each combination of category, layer, and implementation
557 option. Here, we look at how performance changes as the strength varies.

558 Figure 4A (top) plots the increase in binary classification performance (averaged
559 across all categories) as a function of strength for the four different implementation
560 options, when attention is applied at layer 12 for merged images (results similar for
561 array images). From this we can see that not only is the multiplicative bidirectional
562 manipulation the best performing, it also reaches its peak at a lower strength than the
563 other options.

564 On the bottom of Figure 4A, we show the best performing strength (calculated
565 for each category individually and averaged) across layers, and when applied to all
566 layers simultaneously. It is clear from this analysis that multiplicative bidirectional
567 effects consistently require lower strength to reach maximum performance than other
568 options. Furthermore, the fact that the best performing strengths occur below the
569 peak strength tested ($\beta = 11.85$ for individual layers and $\beta = 1.19$ for all layers
570 simultaneously) indicates that any performance limitations are not due to a lack of
571 strength. The best performing strength for additive attention at layer 13 is surprisingly
572 high. To understand why this may be, it is important to remember that, when using
573 additive attention, the attention value added to each unit's response is the product
574 of the relevant tuning value, β , and the average activity level of the layer. This is
575 necessary because average activity levels vary by 2 orders of magnitude across layers.
576 The variability of activity across feature maps, however, is much higher at layer 13
577 compared to layers 1 through 12. This makes the mean activity level used to calculate
578 attention effects less reliable, which may contribute to why higher β values are needed.

579 Performance can change in different ways with attention. In Figure 4B we break the
580 binary classification performance down into true and false positive rates. Here, each
581 colored line indicates a different category and increasing dot size indicates increasing
582 strength of attention (multiplicative bidirectional effects used). True and false positive
583 rates in the absence of attention have been subtracted such that all categories start
584 at (0,0). Ideally, true positives would increase without an equivalent increase (and
585 possibly with a decrease) in false positive rates. If they increase in tandem (i.e.,
586 follow the black dotted lines) then attention would not have a net beneficial effect on
587 performance.

588 Looking at the effects of applying attention at different layers (layer labeled in
589 gray), we can see that attention at lower layers is less effective at moving the per-
590 formance in this space, and that movement is in somewhat random directions. As

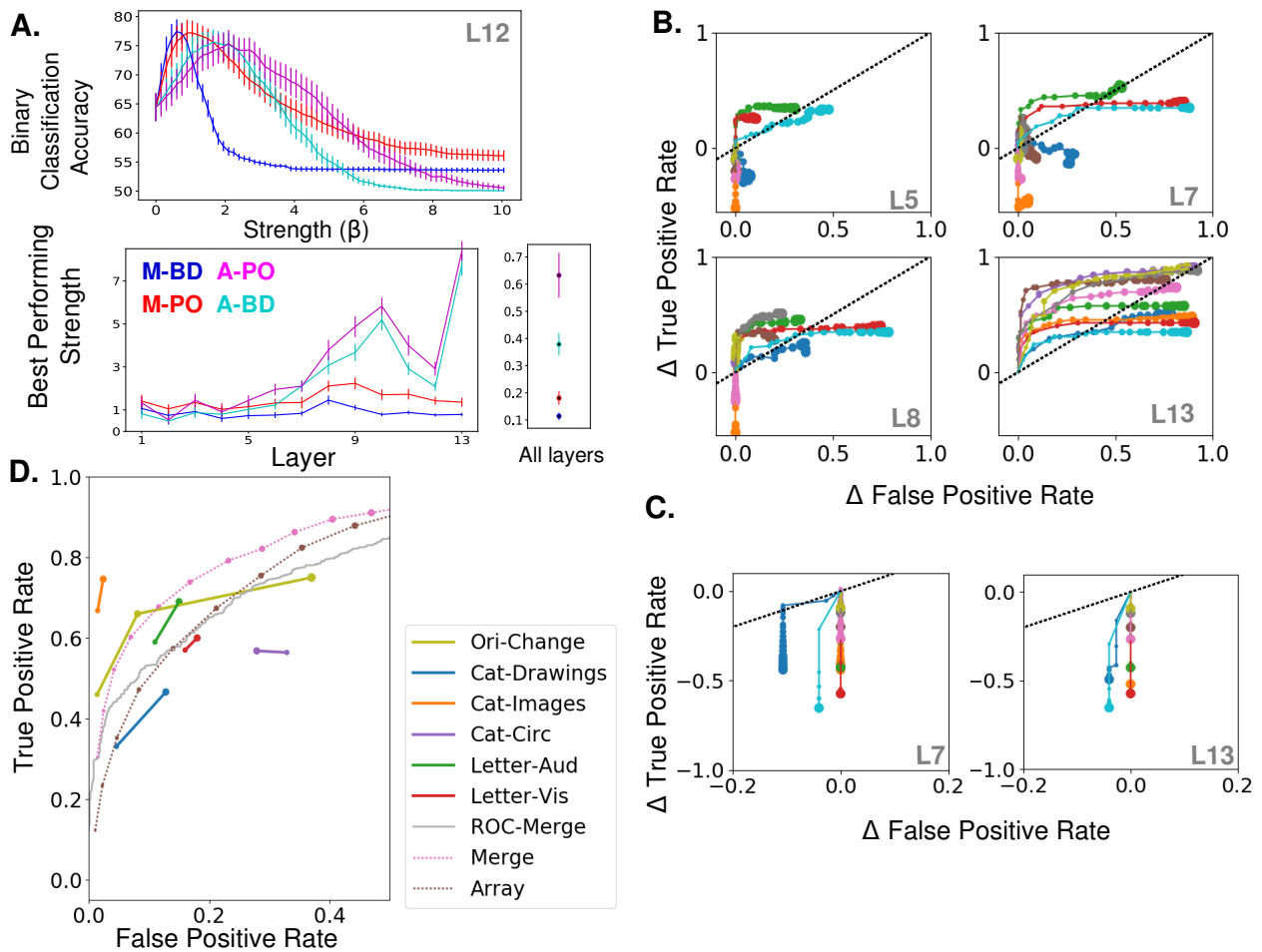


Figure 4: Effects of Varying Attention Strength in Feature-Based Attention Tasks. A.) Effect of strength on different implementation options. On the top, performance averaged over categories (errorbars \pm S.E.M.) shown as a function of the strength parameter, β , for each implementation option. Attention is applied to layer 12 and on merged images. The location of the peak for each category individually is the best performing strength for that category. On the bottom, the best performing strength averaged across categories (errorbars \pm S.E.M.) at each layer for each implementation option. When applied at all layers simultaneously, the range of attention strength tested was smaller. Color scheme as in Figure 1A. B.) and C.) multiplicative bidirectional attention is used, on merged images. B.) Effect of strength increase in true- and false-positive rate space for each of four layers (layer indicated in bottom right of each panel). Each line represents performance changes that arise from applying attention to a different category (only 10 categories shown for visibility), with each increase in dot size representing a .15 increase in strength. Baseline (no attention) values are subtracted for each category such that all start at (0,0) and the layer attention is applied to is indicated in gray. The black dotted line represents equal changes in true and false positive rates. C.) Effect of strength increase in true- and false-positive rate space when tuning values are negated. Same as B, but with sign of attention effects switched (only attention at layer 7 and 13 shown). D.) Comparisons from experimental data. The true and false positive rates from four previously published studies are shown for conditions of increasing attentional strength (solid lines). True and false positive rates are shown for merged and array images (dotted lines, averaged over categories) when attention is applied with increasing strengths (starting at 0, each increasing dot size equals .15 increase in β) at layer 13 (multiplicative bidirectional effects). Receiver operator curve for merged images shown in gray. Cat-Drawings=[55], Exp. 1; Cat-Images=[55],Exp. 2; Objects=[41], Letter-Aud.=[54], Exp. 1; Letter-Vis.=[54], Exp. 2. Ori-Change=[58]. See Methods for details of experiments.

591 attention is applied at later layers, true positive rates are more likely to increase and
592 the increase in false positive rates is delayed. Thus, when attention is applied with
593 modest strength at layer 13, most categories see a substantial increase in true posi-
594 tives with only modest increases in false positives. As strength continues to increase
595 however, false positives increase substantially and eventually lead to a net decrease in
596 overall classifier performance (i.e., cross the black dotted line). Without attention the
597 false negative rate is $69.7 \pm 21.8\%$ and decreases to $19.9 \pm 10\%$ using the best perform-
598 ing strength for each category. Without attention the false positive rate is $1.4 \pm 3.1\%$
599 and increases to $13.7 \pm 7.7\%$ using the best performing strength for each category.

600 To confirm that these behavioral enhancements result from the targeted effects of
601 attention, rather than a non-specific effect of activity manipulation, we apply multi-
602 plicative bi-directional attention using negated tuning values. Because tuning values
603 sum to zero over all feature maps and categories, using negated tuning values doesn't
604 change the overall level of positive and negative modulation applied to the network.
605 Applying attention this way, however, leads to unambiguously different results. Figure
606 4C shows these results, plotted in the same format as Figure 4B, for attention at layers
607 7 and 13. Using negated tuning values leads to a decrease in true and false positive
608 values with increasing attention strength. Thus, attention appears to function as a
609 knob that can turn true and false positives up or down in an intuitive way.

610 It would be useful to know how the magnitude of neural activity changes in our
611 model compare to those used by the brain. Experimentally, the strength of attention
612 can be manipulated by controlling the presence and/or validity of cues [58], switching
613 attention from the non-preferred to preferred stimulus can have large effects on firing
614 rate (111% increase in MT [46]). Before the presentation of a target array, cells in
615 IT showed a 40% increase in firing when the to-be-detected object was preferred
616 versus non-preferred [13]. Of most direct relevance to this study, however, is the
617 modulation strength when switching from no or neutral attention to specific feature-
618 based attention, rather than switching attention from a non-preferred to a preferred
619 stimulus. In [56], neurons in MT showed an average increase in activity of 7% when
620 attending their preferred motion direction (and similar decrease when attending the
621 non-preferred) versus a neutral attention condition.

622 In our model, when $\beta = .75$ (roughly the value at which performance with multi-
623 plicative bidirectional effects peaks at later layers), given the magnitude of the tuning
624 values (average magnitude: .38), attention scales activity by an average of 28.5%. This
625 value refers to how much activity is modulated in comparison to a the $\beta = 0$ condi-
626 tion. This $\beta = 0$ condition is probably more comparable to passive or anesthetized
627 viewing, as task engagement has been shown to scale neural responses generally [70].
628 This complicates the relationship between modulation strength in our model and the
629 values reported in the data.

630 To allow for a more direct comparison, in Figure 4D, we have collected the true
631 and false positive rates obtained experimentally during different object detection tasks
632 (explained in detail in Methods), and plotted them in comparison to the model results.
633 The first five studies plotted in Figure 4D come from human studies. In all of these
634 studies, uncued trials are those in which no information about the upcoming visual
635 stimulus is given, and therefore attention strength is assumed to be low. In cued
636 trials, the to-be-detected category is cued before the presentation of a challenging
637 visual stimulus, allowing attention to be applied to that object or category. The
638 tasks range from detecting simple, stereotyped stimuli (e.g. letters) to highly-varied

639 photographic instances of a given category. Not all changes in performance were
640 statistically significant, but we plot them here to show general trends.

641 The majority of these experiments show a concurrent increase in both true and false
642 positive rates as attention strength is increased. The rates in the uncued conditions
643 (smaller dots) are generally higher than the rates produced by the $\beta = 0$ condition
644 in our model, which suggests that neutrally cued conditions do indeed correspond to
645 a value of $\beta > 0$. We can determine the average β value for the neutral and cued
646 conditions by projecting the data values onto the nearest point on the model line
647 (each dot on the model line corresponds to an increase in β of .15). Specifically, we
648 project the values from the four datasets whose experiments are most similar to our
649 merged image task (Cat-Drawings, Cat-Images, Letter-Aud, and Letter-Vis) onto the
650 model line generated from using the merged images. Through this, we find that the
651 average β value for the neutral conditions is .39 and for the attended conditions .53.
652 Because attention scales activity by $1 + \beta f_c^{lk}$ (where f_c^{lk} is the tuning value and the
653 average tuning value magnitude is .38), these changes correspond to a $\approx 5\%$ change
654 in activity. Thus, the size of observed performance changes is broadly consistent with
655 the size of observed neural changes.

656 Among the experiments used, the one labeled "Cat-Images" is an outlier, as it has
657 much higher true positive and lower true negative rates than the model can achieve
658 simultaneously. This experimental setup is the one most similar to the merged im-
659 ages used in the model (subjects are cued to attend a given category and grayscale
660 category images are presented with a concurrent noise mask), however, the images
661 were presented for 6 seconds. This presumably allows for several rounds of feedback
662 processing, which our purely feedforward model cannot capture. Notably though, true
663 and false positive rate still increase with attention in this task.

664 Another exception is the experiment labeled as "Cat-Circ", which has a larger
665 overall false positive rate and shows a decrease in false positives with stronger attention.
666 In this study, a single target image is presented in a circular array of distractor images,
667 and the subject may be cued ahead of time as to which object category will need to
668 be detected in that target image. The higher false positive rates in this experiment
669 may be attributable to the fact that the distractors were numerous and were pixelated
670 versions of real images. Attention's ability to decrease false positives, however, suggests
671 a different mechanism than the one modeled here. The reason for this difference is not
672 clear. However, in this experiment, the cued trials were presented in blocks wherein
673 the same category was to be detected in each trial, whereas for the uncued trials, the
674 to-be-detected category changed trialwise. The block structure for the attended trials
675 may have allowed for a beneficial downstream adaptation to the effects of attention,
676 which reined in the false positive rate.

677 The last dataset included in the plot (Ori-Change) differs from the others in sev-
678 eral ways. First, it comes from a macaque study that also measured neural activity
679 changes, which allows for a direct exploration of the relationship between neural and
680 performance effects. The task structure is different as well: subjects had to detect an
681 orientation change in one of two stimuli. For cued trials, the change occurs at the cued
682 stimulus on 80% of trials. Attention strength could thus be low (for the uncued stimuli
683 on cued trials), medium (for both stimuli on neutrally-cued trials), or high (for the
684 cued stimuli on cued trials). While this task includes a spatial attention component,
685 it is still useful as a test of feature-based attention effects. Previous work has demon-
686 strated that, during a change detection task, feature-based attention is deployed to the

687 pre-change features of a stimulus [16, 59]. Therefore, because the pre-change stimuli
688 are of differing orientations, the cueing paradigm used here controls the strength of
689 attention to orientation as well. So, while this task differs somewhat from the one
690 performed by the model, it can still offer broad insight into how the magnitude of
691 neural changes relates to the magnitude of performance changes.

692 We plot the true positive (correct change detection) and false positive (premature
693 response) rates as a function of strength as the yellow line in 4D. Like the other
694 studies, this study shows a concurrent increase in both true and false positive rates
695 with increasing attention strength. According to recordings from V4 taken during
696 this task, average firing rates increase by 3.6% between low and medium levels of
697 attention. To achieve the performance change observed between these two levels the
698 model requires a roughly 12% activity change. This gap may indicate the role of
699 other biologically observed effects of attention (e.g., on Fano Factor and correlations)
700 in performance enhancement, or the smaller effect in the data may be due to the
701 averaging of both positive and negative changes (because the stimuli were optimized
702 for a subset of the recorded neurons, positive changes would be expected on average).
703 Firing rates increased by 4.1% between medium and high attention strength conditions.
704 For the model to achieve the observed changes in true positive rates alone between
705 these levels requires a roughly 6% activity change. However, the data shows a very
706 large increase in false positives between these two attention strengths, which would
707 require a roughly 20% activity change in the model. This high rate of false positives
708 points to a possible effect of attention downstream of sensory processing.

709 Finally, we show in this plot the change in true and false positive rates when the
710 threshold of the final layer binary classifier is varied (a receiver operating characteristic
711 analysis. No attention was applied during this analysis). The gray line in Figure
712 4D shows this analysis for merged images. Comparing this to the effect of varying
713 attention strength (pink line), it is clear that varying the strength of attention applied
714 at the final convolutional layer has more favorable performance effects than altering
715 the classifier threshold. This points to the role of attentional modulation in sensory
716 areas, rather than targeting only downstream "readout" areas.

717 Overall, the findings from these studies suggest that much of the change in true
718 and false positive rates observed experimentally could be attributed to moderately-
719 sized changes in neural activity in sensory processing areas. However, it is clear that
720 the details of the experimental setup are relevant, both for the absolute performance
721 metrics and how they change with attention [68].

722 An analysis of performance changes in the context of signal detection theory (sen-
723 sitivity and criteria) will come later.

724 *3.4. Spatial Attention Increases Object Categorization Performance*

725 In addition to feature-based attention, we also test the effects of spatial attention
726 in this network. For this, we use our array images, and the task of the network
727 is to correctly classify the object category in the attended quadrant of the image.
728 Therefore, the original final layer of the network which performs 1000-way object
729 categorization is used (Figure 5A). The same implementation and layer options were
730 tested and compared to 1000-way classification performance without attention (see
731 Methods 2.5.4). However, tuning values were not used; rather, because the spatial
732 layout of activity is largely conserved in CNNs, an artificial neuron was assumed to
733 "prefer" a given quadrant of the image if that unit was in the corresponding quadrant

734 of the feature map.

735 In Figure 5B, the performance (classification was considered correct if the true
736 category label appeared in the top five categories outputted by the network, but trends
737 are the same for top-1 error) is shown as a function of attention strength for each of
738 the four options. The layer at which attention is applied is indicated by the line color.
739 Because tuning values are not used for the application of spatial attention, the β value
740 can be interpreted directly as the amount of activity modulation due to attention
741 (recall that for multiplicative effects rates are multiplied by $1 + \beta$).

742 Using experimentally-observed performance changes to relate our model to data
743 (as we did in Figure 4D) is more challenging for the spatial attention case because the
744 specific tasks used are more varied. Using the performance on trials with a neutral
745 spatial cue as a baseline, we report the impact of spatial attention as the factor by
746 which performance increases on trials with valid spatial cues. Experimentally, spatial
747 attention scales performance by $\approx 19\%$ on a color recognition task [28], $\approx 16\%$ on an
748 orientation categorization task [20], $\approx 10\%$ on an orientation classification task [78] and
749 a gap detection task [64], and $\approx 3.3\%$ on a red line detection task [89]. Spatial attention
750 effects range in magnitude but have been shown to increase neural activity by $\approx 20\%$ in
751 several studies [61, 18] when calculated for attend-in versus attend-out conditions. If
752 we assume that attend-in and attend-out conditions scale activity in opposite directions
753 (bi-directional effects) but with equal magnitude from a baseline [58], then spatially
754 cued trials should have a roughly 10% change in activity compared to neutral trials.
755 As mentioned above, the $\beta = 0$ condition in our model is not necessarily comparable
756 to a neutrally-cued condition experimentally, so it is unclear what performance level in
757 our model should be used as a baseline. However, going from $\beta = 0$ to $\beta = .1$ enhances
758 performance from 14% correct to an average (across attention at each layer) of 17.4%
759 correct. This is a 24.2% increase in accuracy stemming from a 22% change in activity
760 on attend-in versus attend-out conditions. Again, these simple calculations suggest
761 that the experimentally-observed magnitude of neural modulations could indeed lead
762 to the observed magnitude of behavioral changes.

763 It is also of note that performance in the case of multiplicative bidirectional effects
764 plateaus around $\beta = 1$, yet for multiplicative positive-only effects it continues to climb.
765 This suggests that the suppressing of the three non-attended quadrants is a strong
766 driver of the performance changes when using multiplicative bidirectional effects, as
767 this suppression is complete at $\beta = 1$ (i.e., activity is 100% silenced at that value).
768 While it is not believed that spatial attention leads to complete silencing of cells
769 representing unattended locations, these results highlight the potential importance of
770 scaling such activity downward.

771 Figure 5C and D summarize the performance enhancements that result from differ-
772 ent options (assuming the best performing strengths, as in Figure 3B and C). Unlike
773 feature-based attention, spatial attention is relatively insensitive to the layer at which
774 it is applied, but is strongly enhanced by using multiplicative bidirectional effects com-
775 pared to others. This discrepancy makes sense when we consider that spatial attention
776 tasks are cross-modal—that is, they involve attending to one dimension (space) and
777 reading out another (object category)—whereas the object detection tasks used above
778 are unimodal—the same dimension (object category) is attended to and read out. In
779 a cross-modal task it is not valuable just to amplify the attended attribute, but rather
780 to amplify the information carried by the attended attribute. Assuming the absolute
781 difference in rates across cells is relevant for encoding object identity, multiplicative

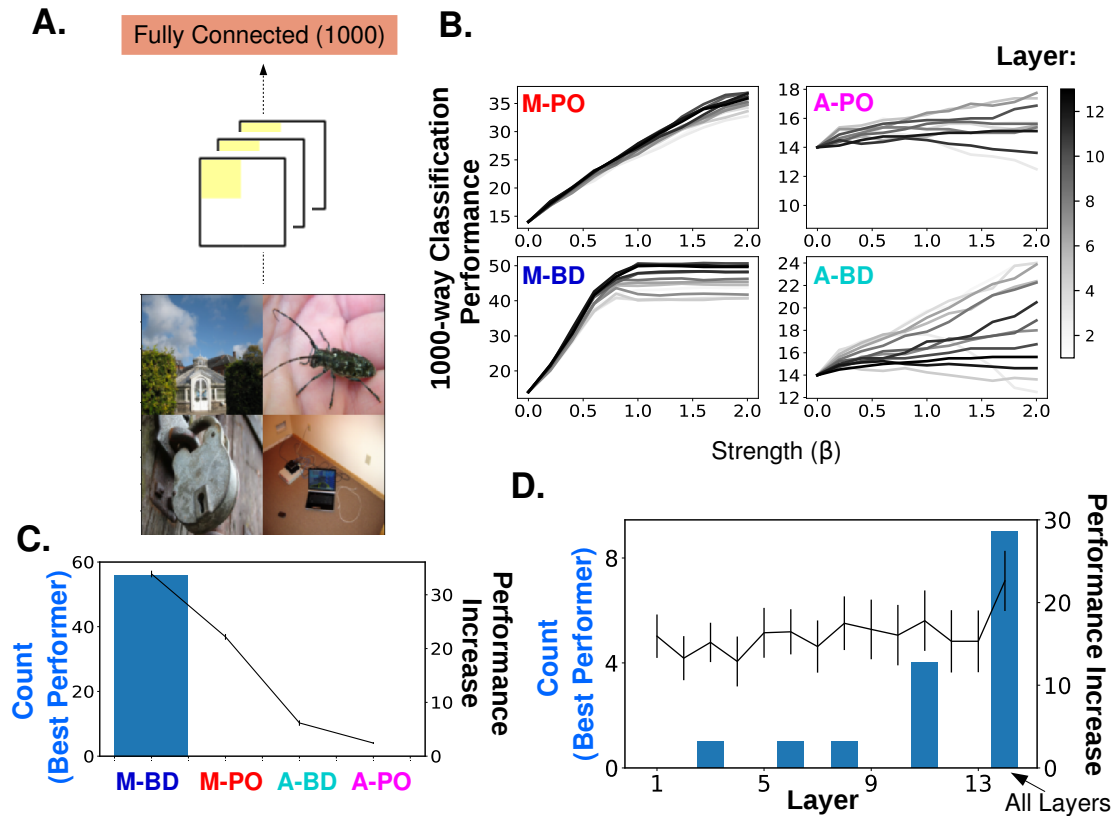


Figure 5: Spatial Attention Task and Results. A.) Array images were used to test spatial attention. Performance was measured as the ability of the original 1000-way classifier to identify the category in the attended quadrant (measured as top-5 error). Attention was applied according to the spatial layout of the feature maps (for example, when attending to the upper left quadrant of the image, units in the upper left quadrant of the feature maps are enhanced). B.) 1000-way classification performance as a function of attention strength, when applied at different layers (indicated by line darkness) and for each of the four attention options. C.) Comparison of performance effects of attention options (using best performing strength). Histograms (left axis) show how often the given option is the best performing (over 4 quadrants * 14 layer options = 56 instances). Average increase in 1000-way classification performance for each option also shown (right axis, errorbars +/- S.E.M.). D.) Histograms (over 4 quadrants * 4 implementation options = 16 instances) show how often the applying attention to the given layer is the best performing. The final column corresponds to attention applied to all layers simultaneously (strength at one-tenth that of strength applied to individual layers). Average increase in 1000-way classification performance for each layer also shown (right axis, errorbars +/- S.E.M.).

782 effects amplify these informative differences and can thus aid in object classification
783 in the attended quadrant. In a system with noise, attention's benefits would depend
784 on the extent to which it simultaneously enhanced the non-informative noise. Exper-
785 imentally, attention leads to a decrease in mean-normalized variance in firing across
786 trials [15].

787 Another difference between feature-based and spatial attention is the effect of ap-
788 plying attention at all layers simultaneously. When applying attention at all layers,
789 the β values tested are one-tenth that of when attention is applied at individual lay-
790 ers. Despite this weakened strength, applying attention at all layers leads to better
791 performance in the spatial attention task than applying it to any layer individually.
792 In the feature-based attention task, this is not the case (Figure 3C). This difference is
793 explored more directly later.

794 *3.5. Feature-based Attention Enhances Performance on Orientation Detection and* 795 *Color Classification Tasks*

796 Some of the results presented above, particularly those related to the layer at
797 which attention is applied, may be influenced by the fact that we are using an object
798 categorization task. To see if results are comparable using simpler stimuli, we created
799 an orientation detection task (Figure 6A), wherein binary classifiers trained on full
800 field oriented gratings are tested using images that contain two gratings of different
801 orientation and color. The performance of these binary classifiers without attention
802 is above chance (distribution across orientations shown in inset of Figure 6A). The
803 performance of the binary classifier associated with vertical orientation (0 degrees) was
804 abnormally high (92% correct without attention, other orientations average 60.25%)
805 and this orientation was excluded from further analysis for the detection task.

806 Attention is applied according to orientation tuning values of the feature maps
807 (tuning quality by layer is shown in Figure 6C) and tested across layers (using multi-
808 plicative bidirectional effects). We find that the trend in this task is similar to that of
809 the object task: applying attention at later layers leads to larger performance increases
810 (14.4% percentage point increase at layer 10). This is despite the fact that orientation
811 tuning quality peaks in the middle layers.

812 We also explore a cross-modal attention task that is in line with the style of cer-
813 tain attention experiments in neuroscience and psychology [80, 67, 98]. Specifically,
814 the task for the network is to readout the color of the stimulus in the image with
815 the attended orientation (Figure 6B, mean 5-way classification performance without
816 attention: 42.89%). Thus, attention is applied according to orientation tuning values,
817 but the final layer of the network is a 5-way color classifier. This is akin to studies
818 where the task of the subject is, for example, to report a speed change in random dots
819 that are moving in the attended direction. Interestingly, in this case attention applied
820 at earlier layers (specifically layers 2-6, best performance increase is 7.8 percentage
821 points at layer 2) performs best. Color tuning quality is stronger at earlier layers as
822 well (layers 1-3 particularly).

823 The β values that lead to peak performance in the detection task at later layers
824 ranges from .5 to 1. Given that β scales the tuning values and the average tuning
825 value magnitude at later layers is .32, the average modulation strength (compared
826 to the $\beta = 0$ condition) is 16%-32%. For the color classification task the successful
827 modulation at earlier layers ranges from 13-28%. Therefore the two different tasks
828 require similar modulations.

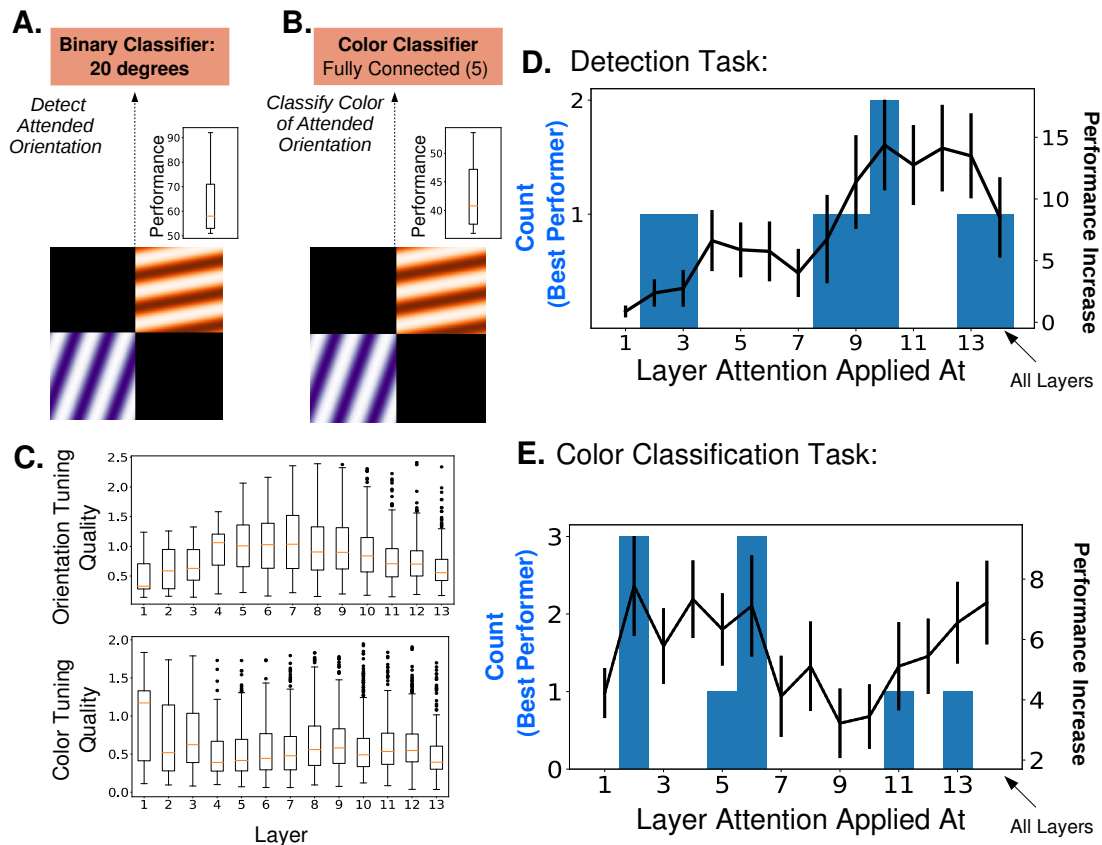


Figure 6: Attention Tasks and Results Using Oriented Gratings. A.) Orientation detection task. Like with the object category detection tasks, separate binary classifiers trained to detect each of 8 different orientations replaced the final layer of the network. Test images included 2 oriented gratings of different color and orientation located at two of 4 quadrants. Insets show performance over 9 orientations without attention B.) Color classification task. The final layer of the network is replaced by a single 5-way color classifier. The same test images are used as in the detection task and performance is measured as the ability of the classifier to identify the color of the attended orientation. Inset shows performance over 9 orientations without attention (chance is 25%) C.) Orientation tuning quality (top) and color tuning quality (bottom) as a function of layer. D.) Comparison of performance on detection task when attention (determined by orientation tuning values) is applied at different layers. Histogram of best performing layers in blue, average increase in binary classification performance in black. E.) Comparison of performance on color classification task when attention (determined by orientation tuning values) is applied at different layers. Histogram of best performing layers in blue, average increase in 5-way classification performance in black. Errorbars are +/- S.E.M.

829 3.6. Gradient Values Offer Performance Comparison

830 Previously, we used gradient values to determine if object category tuning values
831 were related to classification behavior. Here, we use a similar procedure to obtain
832 gradient values that tell us how feature map activity should change in order to make
833 the network better at the tasks of orientation detection and color classification (see
834 Methods 2.5.2). We then use these values in place of the orientation tuning values
835 when applying attention, and compare the performances.

836 In Figure 7A, we first show the extent to which these gradient values correlate with
837 the tuning values. On the left, an example feature map's tuning curve (green) along
838 with curves generated from gradient values for the orientation detection task (solid
839 purple) and color classification task (dashed purple). The middle and right panels
840 show the average correlation coefficients between tuning curves and the respective
841 gradient values across layers. Correlation with orientation detection gradients peaks
842 at later layers, while correlation with color classification gradients peaks at early layers.
843 In Figure 7B, the solid lines and histograms document the performance using gradient
844 values. For comparison, the dashed lines give the performance improvement from
845 using the tuning values. In the orientation detection task, gradient values perform
846 better than tuning values at earlier layers, but the performance difference vanishes
847 at later layers (where the tuning values and gradient values are most correlated).
848 Thus, tuning values can actually reach the same performance level as the gradient
849 values suggesting that, while they are not identical to the values determined by the
850 gradient calculations, they are still well-suited for increasing detection performance.
851 The performance for color classification using gradient values has the reverse pattern.
852 It is most similar to the performance using tuning values at earlier layers (where the
853 two are more correlated), and the performance gap is larger at middle layers. At all
854 layers, the mean performance using gradient values is larger than that using tuning
855 values.

856 The results of applying this procedure to the object category detection task are
857 discussed later (Figure 8E).

858 3.7. Feature-based Attention Primarily Influences Criteria and Spatial Attention Pri- 859 marily Influences Sensitivity

860 Signal detection theory is frequently used to characterize the effects of attention
861 on performance [96]. Here, we use a joint feature-spatial attention task to explore
862 effects of attention in the model. The task uses the same 2-grating stimuli described
863 above. The same binary orientation classifiers are used and the task of the model is to
864 determine if a given orientation is in a given quadrant. Performance is then measured
865 when attention is applied according to orientation, space, or both (effects are combined
866 additively), and two key signal detection measurements are computed. Criteria is a
867 measure of how lenient is the threshold that's used to mark an input as a positive.
868 Sensitivity is a measure of how separate the two populations of positive and negatives
869 are.

870 Figure 7C shows how these values, along with the overall binary classification
871 performance, vary with the strength and type of attention applied at two example
872 layers. Performance is best when both spatial and feature-based attention are applied
873 simultaneously. The ways in which these two types of attention affect performance can
874 be teased apart by looking at their effects when applied separately. Criteria decreases
875 more when feature-based attention is applied alone than when spatial is. Sensitivity

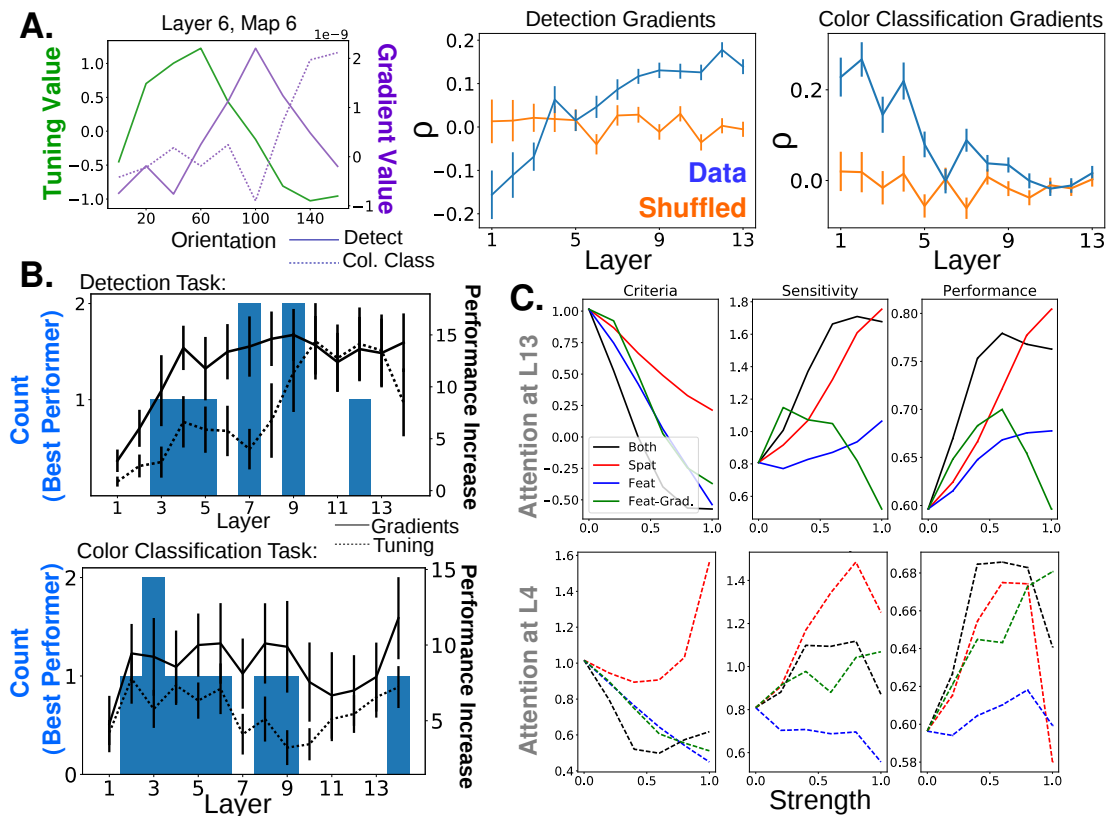


Figure 7: Comparison of Orientation Task Gradient Values to Tuning Values. A.) Correlation of gradient values with tuning values for the detection and color classification tasks. On the left, an example feature maps orientation tuning curve (green) and curves generated from detection gradient values (solid purple) and color classification gradient values (dashed purple). Correlation coefficients with tuning curve are -0.196 and -0.613 , respectively. Average correlation coefficient values between tuning curves and detection gradient curves (middle) and color classification gradient curves (right) across layers (blue). Shuffled correlation values in orange. Errorbars are \pm S.E.M. B.) Comparison of performance on detection task when attention is determined by detection gradient values and applied at different layers (top). Comparison of performance on color classification task when attention is by determined by color classification gradient values and applied at different layers (bottom). Histograms of best performing layers in blue, average increase in binary or 5-way classification performance in black. Errorbars are \pm S.E.M. In both, performance increase when attention is determined by tuning values is shown for comparison (dashed lines). Only multiplicative bidirectional effects are used. C.) Change in signal detection values when attention is applied in different ways (spatial, feature according to tuning, both spatial and feature according to tuning, and feature according to gradient values) for the task of detecting a given orientation at a given quadrant. Top row is when attention is applied at layer 13 and bottom when applied at layer 4 (multiplicative bidirectional effects).

876 increases more for spatial attention alone than feature-based attention alone. These
877 general trends hold regardless of the layer at which attention is applied, though when
878 applied at layer 4, feature-based attention alone actually decreases sensitivity.

879 Applying feature-based attention according to orientation detection gradient values
880 has a very similar effect on criteria as applying it with tuning values. The effect
881 on sensitivity however, is slightly different, as the gradient values are better able to
882 increase sensitivity. Therefore, attending to feature using gradient values leads to
883 slightly better overall performance than when using tuning values in this example.

884 Various impacts of attention on sensitivity and criteria have been found experi-
885 mentally. Task difficulty (an assumed proxy for attentional strength) was shown to
886 increase both sensitivity and criteria [87]. In line with our results, spatial attention has
887 been found to increase sensitivity and (less reliably) decrease criteria [32, 21], and fea-
888 ture attention is known to decrease criteria, with minimal effects on sensitivity [74, 2].
889 A study that looked explicitly at the different effects of spatial and category-based at-
890 tention [88] found that, in line with our results, spatial attention increases sensitivity
891 more than category-based attention (most visible in their Experiment 3c, which uses
892 natural images) and that the effects of the two are additive.

893 The diversity of results in the literature (including discrepancies with our model)
894 may be attributed to different task types and to the fact that attention is known
895 to impact neural activity in various ways beyond pure sensory areas [43]. This idea
896 is borne out by a study that aimed to isolate the neural changes associated with
897 sensitivity and criteria changes [53]. In this study, the authors designed behavioral
898 tasks that encouraged changes in sensitivity or criteria exclusively: high sensitivity was
899 encouraged by associating a given stimulus location with higher overall reward, while
900 high criteria was encouraged by rewarding correct rejects more than hits (and vice versa
901 for low sensitivity/criteria). Differences in V4 neural activity were observed between
902 trials using high versus low sensitivity stimuli. No differences were observed between
903 trials using high versus low criteria stimuli. This indicates that areas outside of the
904 ventral stream (or at least outside V4) are capable of impacting criteria. Importantly,
905 it does not mean that changes in V4 don't impact criteria, but merely that those
906 changes can be countered by downstream processes. Indeed, to create sessions wherein
907 sensitivity was varied without any change in criteria, the authors had to increase the
908 relative correct reject reward (i.e., increase the criteria) at locations of high absolute
909 reward, presumably to counter the decrease in criteria that appeared naturally as a
910 result of attention-induced neural changes in V4 (similarly, they had to decrease the
911 correct reject reward at low reward locations). Our model demonstrates clearly how
912 such effects from sensory areas alone can impact detection performance, which, in turn
913 highlights the role downstream areas play in determining the final behavioral outcome.

914

915 *3.8. Recordings Show How Feature Similarity Gain Effects Propagate*

916 To explore how attention applied at one location in the network impacts activity
917 later on, we apply attention at various layers and "record" activity at others (Figure
918 8A). In particular, we record activity of feature maps at all layers while applying mul-
919 tiplicative bidirectional attention at layers 2, 6, 8, 10, and 12 individually. The results
920 of these recordings show both which features of the activity changes are correlated
921 with performance enhancements as well as how FSGM effects at one area can lead to
922 very different effects at another.

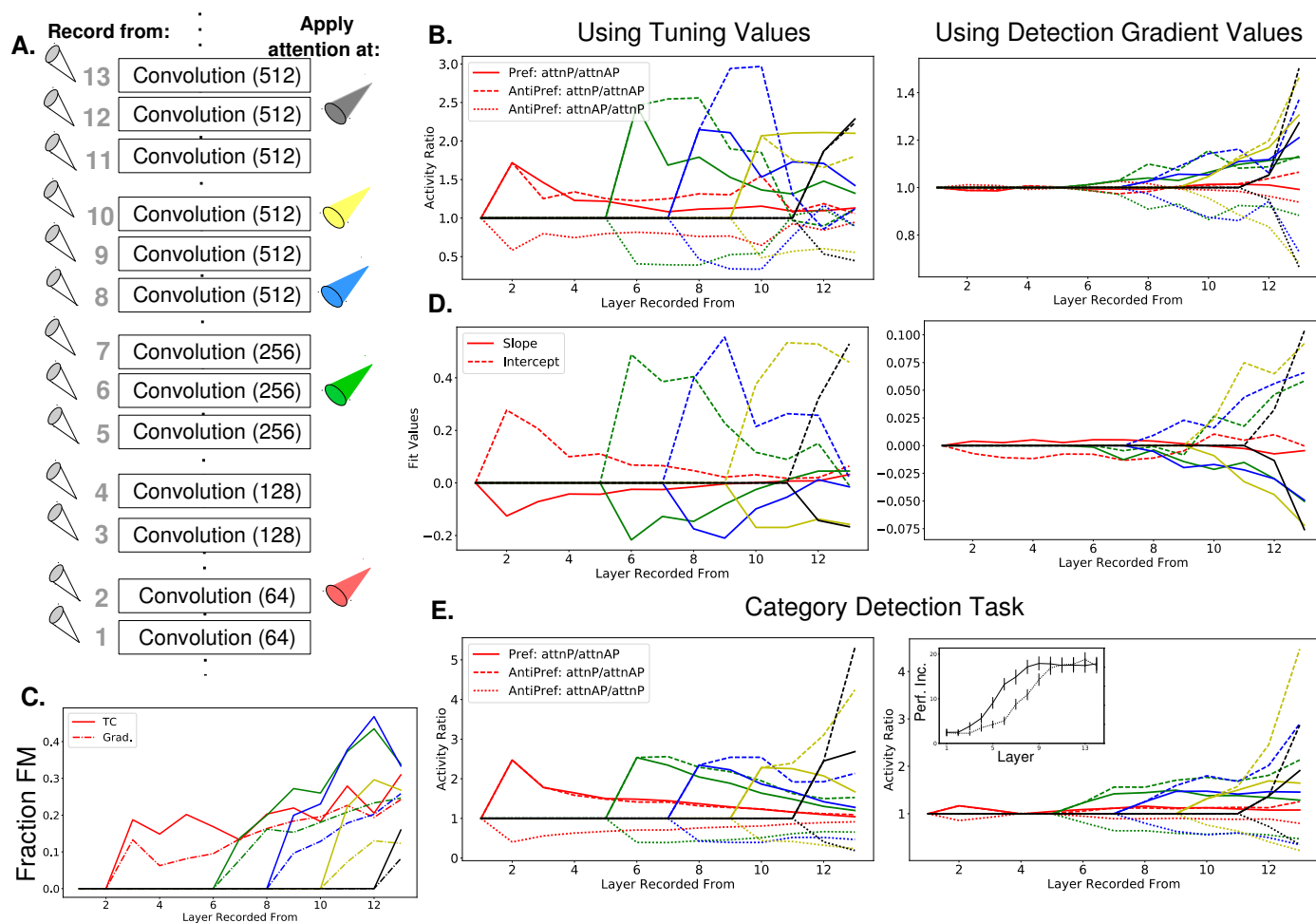


Figure 8: How Activity Changes from Attention Propagate for Unimodal Tasks. A.) Recording setup. The spatially averaged activity of feature maps at each layer was recorded (left) while attention was applied at layers 2, 6, 8, 10, and 12 individually. Activity was in response to a full field orientated grating for (B), (C), and (D) or full field standard ImageNet images for (E). Attention was always multiplicative and bidirectional. B.) Activity ratios for different attention conditions as a function of recorded layer when attention is applied at different layers (given by color as in (A)). Line style indicates whether the stimulus presented is preferred (solid line) or anti-preferred (dashed and dotted lines), and whether the ratio is calculated as activity when the preferred is attended divided by when the anti-preferred is attended (solid and dashed) or the reverse (dotted). Values are medians over all feature maps. Orientation tuning values (left) or orientation detection gradient values (right) are used for applying attention. C.) The fraction of feature maps that display feature matching (FM) behavior, defined as activity ratios greater than one for Pref:AttnP/AttnAP and AntiPref:AttnAP/AttnP when attention is applied according to orientation tuning curve values (solid) or detection gradient values (dashed). D.) Dividing activity when a given orientation is present and attended by activity when no attention is applied gives a set of activity ratios. Ordering these ratios from most to least preferred orientation and fitting a line to them gives the slope and intercept values plotted here (intercept values are plotted in terms of how they differ from 1, so positive values are an intercept greater than 1). Values are medians across all feature maps at each layer with attention applied at layers indicated in (A). E.) Same as in (B) but using object category images, tuning values, and detection gradient values. The inset on the right shows mean performance detection over all 20 categories when attention is applied at different layers using category detection gradient values (solid line, performance using tuning values shown as dotted line for comparison). Errorbars S.E.M.)

923 Activity was recorded in response to multiple different stimuli and attentional
924 conditions. In Figure 8B we explore whether applying feature attention according to
925 the FSGM at one layer continues to have FSGM-like effects at later layers. To do this
926 we use an analysis taken from [56]. Specifically, full field oriented gratings were shown
927 to the network that were either of the preferred (most positive tuning value) or anti-
928 preferred (most negative tuning value) orientation for a given feature map. Attention
929 was also applied either to the preferred or anti-preferred orientation. According to
930 the FSGM, the ratio of activity when the preferred orientation is attended divided
931 by activity when the anti-preferred orientation is attended should be larger than one
932 regardless of whether the orientation of the stimulus is preferred or not (indeed, the
933 ratio should be constant for any stimulus). An alternative model, the feature matching
934 (FM) model, suggests that the effect of attention is to amplify the activity of a neuron
935 whenever the stimulus in its receptive field matches the attended stimulus. In this
936 case, the ratio of activity when the preferred stimulus is attended over when the anti-
937 preferred is attended would only be greater than one when the stimulus is the preferred
938 orientation. If the stimulus is the anti-preferred orientation, the inverse of the that
939 ratio would be greater than one.

940 In Figure 8B, we plot the median value of these ratios across all feature maps at a
941 layer when attention is applied at different layers, indicated by color. When attention
942 is applied directly at a layer according to its tuning values (left), FSGM effects are
943 seen by default. As these activity changes propagate through the network, however,
944 the FSGM effects wear off. Thus, when attention is applied at an early layer, it does
945 not create strong changes in the final convolutional layer and thus cannot strongly
946 impact the classifier. This explains the finding (Figure 6D) that attention works best
947 for the detection task when applied at later layers, as the only way for strong FSGM
948 effects to exist at the final layers is to apply attention near the final layers.

949 The notion that strong FSGM-like effects at the final layer are desirable for in-
950 creasing classification performance is further supported by findings using the gradient
951 values. In Figure 8B(right), we show the same analysis, but while applying atten-
952 tion according to orientation detection gradient values rather than tuning values. The
953 effects at the layer at which attention is applied do not look strongly like FSGM, how-
954 ever FSGM properties evolve as the activity changes propagate through the network,
955 leading to clear FSGM-like effects at the final layer.

956 These results are recapitulated in Figure 8D using a broader analysis also from
957 [56]. Here, the activity of a feature map is calculated when attention is applied to
958 the orientation in the stimulus and divided by the activity in response to the same
959 orientation when no attention is applied. These ratios are organized according to
960 orientation preference (most to least) and a line is fit to them. According to the FSGM
961 of attention, this ratio should be greater than one for more preferred orientations and
962 less than one for less preferred, creating a line with an intercept greater than one
963 and negative slope. As expected, applying attention according to tuning values causes
964 similar changes at the layer at which it is applied in this model (intercept values are
965 plotted in terms of how they differ from one. Comparable average values from [56] are
966 intercept: .06 and slope: 0.0166). Again, these effects wear off as the activity changes
967 propagate through the network. Also gradient values ultimately lead to this kind of
968 change at the final layer (right panel).

969 While Figure 8B and D show FSGM-like effects according to median values across
970 all feature maps, some individual feature maps may show different behavior. In Fig-

971 ure 8C, we calculate the fraction of feature maps at a given layer that show feature
972 matching behavior (defined as having activity ratios greater than one when the stimu-
973 lus orientation matches the attended orientation for both preferred and anti-preferred
974 orientations). As early as one layer post-attention feature maps start showing feature
975 matching behavior, and the fraction grows as activity changes propagate. Interest-
976 ingly, applying attention according to gradient values also causes an increase in the
977 fraction of feature maps with FM behavior, even as the median values become more
978 FSGM-like. The attention literature contains conflicting findings regarding the fea-
979 ture similarity gain model versus the feature matching model [67, 80]. This may result
980 from the fact that FSGM effects can turn into FM effects as they propagate through
981 the network. In particular, this mechanism can explain the observations that feature
982 matching behavior is observed more in FEF than V4 [106] and that match information
983 is more easily readout from perirhinal cortex than IT [69].

984 We explore the propagation of these effects for category-based attention as well. In
985 Figure 8E, we perform the same analysis as 8B, but with attention applied according
986 to object category tuning values and stimuli that are full-field standard ImageNet
987 images. We also calculate gradient values that would increase performance on category
988 detection tasks (the same procedure used to calculate orientation detection gradients).
989 The binary classification performance increase that results from applying attention
990 according to these values is shown in Figure 8E(right, inset, solid line) in comparison
991 to that when applying according to tuning values (dashed line). Like with orientation
992 detection gradient values, applying attention according to these values propagates
993 through the network to result in FSGM-like effects at the final layer. Also as with the
994 orientation findings, the size of the FSGM effects that reach the final layer track with
995 how well applying attention increases performance; for example, applying attention at
996 layer 2 (red lines) does not lead to strong FSGM effects at the final layer and does not
997 strongly increase performance.

998 *3.9. Attention Alters Encoding Properties in Cross-Modal Tasks*

999 The above recordings looked at how encoding of the attended dimension changed
1000 with attention. In cross-modal tasks, such as the spatial attention task and color
1001 classification task, the encoding that is relevant for performance is the that of the
1002 read-out dimension. We therefore measured how that encoding changes with attention
1003 at different layers as well.

1004 For the spatial attention task, we measured category encoding by fitting a line to a
1005 set of activity ratios (see Figure 9A, left). Those activity ratios represent the activity
1006 of a quadrant when a given object category was in it and the quadrant was attended
1007 divided by activity when the same category was in the quadrant and no attention was
1008 applied. Arranging these from most to least preferred category for each feature map
1009 and fitting a line to them gives two values per feature map: the intercept (the ratio
1010 for the most preferred category, measured in terms of its difference from one) and the
1011 difference (the ratio for the most preferred minus the ratio for the least preferred, akin
1012 to the slope). A purely multiplicative effect leads to a positive intercept value and zero
1013 difference. This effect is clearly observed at the layers at which attention is applied in
1014 Figure 9A(right). It also continues with only a small amount of decay as the activity
1015 changes propagate through the network. By the final layer, the median intercept is still
1016 positive. The median difference becomes negative, indicating that preferred categories
1017 are enhanced more than non-preferred. The values at the final layer are fairly similar

1018 regardless of the layer at which attention was applied. This is in line with the fact
1019 that performance with multiplicative spatial attention is only moderately affected by
1020 the layer at which attention is applied (Figure 5B).

1021 We also looked at how color encoding changes when attention is applied to orien-
1022 tation. Here, we use 2-grating stimuli like those in Figure 6B to ask if encoding of
1023 the color of the grating with a given orientation increases when attention is applied
1024 to that orientation versus when it is applied to the orientation of the other grating
1025 (160 and 80 degree gratings were used). Arranging activity levels from most to least
1026 preferred color, we consider the encoding better if both the overall activity level is
1027 higher and the slope is more negative (see Figure 9B, left). We then measure the
1028 percent of feature maps that have better encoding of 160 degrees when attending 160
1029 degrees versus attending 80 degrees. Looking at those feature maps that most prefer
1030 160 degrees (solid lines, Figure 9B, right), nearly all feature maps enhance their color
1031 encoding at the layer at which attention was applied. However this percent decreases
1032 as the activity changes propagate through the network. On the other hand, for feature
1033 maps that anti-(or least) prefer 160 degrees, none have better encoding at the layer at
1034 which attention was applied, but the percent increases as activity changes propagate
1035 through the layers. Essentially, the burden of better encoding becomes evenly spread
1036 across feature maps regardless of preferred orientation.

1037 This helps understand why, when applying attention according to tuning values,
1038 color classification performance is high at early layers, falls off at mid layers, and
1039 then recovers at final layers (Figure 6E, bottom). This is due to the different effects
1040 attention at these layers have on the final layer. When attention is applied at early
1041 layers, fewer final layer feature maps that prefer the attended orientation have better
1042 encoding, but many that don't prefer it do. When applied at late layers, a high percent
1043 of final layer feature maps that prefer the attended orientation have better encoding,
1044 even if those that don't prefer it do not. When attention is applied at middle layers,
1045 the effect on final layer feature maps that prefer the orientation has decayed, but the
1046 effect on those that don't prefer it hasn't increased much yet. Therefore performance
1047 is worse.

1048 The idea that both feature maps that prefer and anti-prefer the attended orienta-
1049 tion should enhance their color encoding is borne out by the gradient results. When
1050 attention is applied according to gradient values (Figure 9B, bottom), the percent of
1051 feature maps with better encoding is roughly equal for both those that prefer and
1052 anti-prefer the attended orientation. Experimentally, MT neurons have been found to
1053 better encode the direction of motion of a stimulus of the attended color as compared
1054 to a simultaneously presented stimulus of a different color [98]. Importantly, this effect
1055 of attention was *not* stronger when the preferred color was attended (indeed, there was
1056 a slight negative correlation between color preference and attention effect strength).
1057 This is not predicted by the FSGM directly, but as our model indicates, could result
1058 from FSGM-like effects at earlier areas, such as V1.

1059 *3.10. Applying Feature-based Attention at Multiple Layers Counteracts Effects*

1060 It is conceivable that feature-based attention applied at a lower layer could be as (or
1061 more) effective in modulating the activity of feature maps at a later layer as applying
1062 attention at that layer directly. In particular, for a given filter at layer l that prefers
1063 the attended category, bidirectional attention applied at layer $l - 1$ could decrease the
1064 activity of units that have negative weights to the filter and increase the activity of

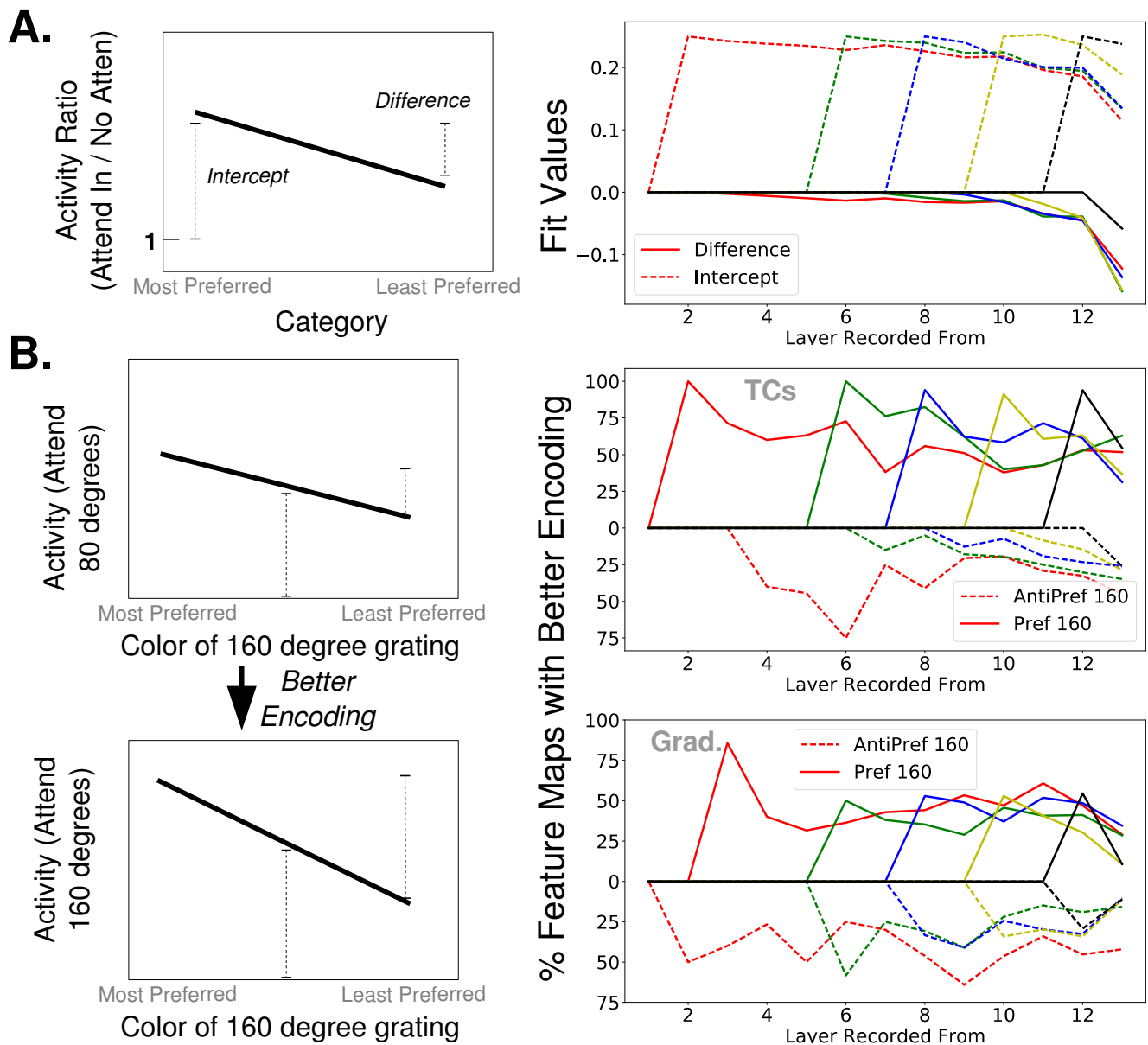


Figure 9: How Activity Changes from Attention Propagate for Cross-modal Tasks. A.) For each feature map, activity averaged over the attended quadrant when attention is applied to it is divided by activity when attention is not applied. Arranging these activity ratios from when the most to least preferred category is present in the quadrant and fitting a line to them results in the intercept and difference values as diagrammed on the left. Specifically, the intercept is the ratio for the most preferred category minus 1 and the difference is the ratio for the most preferred category minus the ratio for the least preferred. On the right, the median fit values across all feature maps are shown for each layer when attention is applied at layers indicated in 8A. B.) Orientated grating stimuli like those in 6B were designed with one grating at 140 degrees and the other at 60. Encoding of the color of the 140 degree grating is measured by fitting a line to the activity (spatially averaged over entire feature map) evoked by when each color is presented in the 140 degree grating (averaged over all colors presented in the 60 degree grating), ordered from most to least preferred. If the intercept (at the middle of this line) and difference increase when attention is applied to 140 degrees compared to attention at 60 degrees, the feature map has better encoding. On the right, the percent of feature maps with better encoding, segregated according to those that prefer 140 degrees (solid line) and those that anti-pref (least prefer) 140 degrees (dashed lines, presented on a mirrored y-axis for visibility). Attention applied according to orientation tuning values (top) or color classification gradients (bottom).

1065 units that have positive weights to the filter (note that in a more biologically-realistic
1066 model, the negatively weighted components would come indirectly from di-synaptic
1067 feedforward inhibition or surround interactions, as feedforward connections are largely
1068 excitatory). For example, if for a given unit in response to a given image the sum
1069 of its positively-weighted inputs is a , and the sum of its negatively-weighted inputs
1070 is b , without any attention, net input is $a - b$. If attention at $l - 1$ scales positively-
1071 weighted inputs up by 20% and negatively-weighted inputs down by 20%, the total
1072 input is now $1.2a - .8b$. These would lead to a greater net activity level than attention
1073 at l itself, which would just scale the net input by 1.2: $1.2(a - b)$. Therefore, given
1074 the same strength, applying attention at layer $l - 1$ could be a more effective way to
1075 modulate activity than applying it at layer l directly. However this assumes a very
1076 close alignment between the preferences of the feature maps at $l - 1$ and the weighting
1077 of the inputs into l .

1078 We investigate this alignment by applying attention to object categories at various
1079 layers and recording at others (stimuli are standard ImageNet images of the attended
1080 category). The ratio of activity when attention is applied at a lower layer is divided
1081 by that when no attention is applied. Feature maps are then divided according to
1082 whether they prefer the attended category (have a tuning value greater than zero) or
1083 don't prefer it (tuning value less than zero). The strength value used is $\beta = .5$, therefore
1084 if attention at lower layers is more effective, we should see activity ratios greater than
1085 1.5 for feature maps that prefer the attended category. The histograms in Figure 10A
1086 (left) show that the majority of feature maps that prefer the attended category (red)
1087 have ratios less than 1.5, regardless of the layer of attention or recording. In many
1088 cases, these feature maps even have ratios less than one, indicating that attention at
1089 a lower layer decreases the activity of feature maps that prefer the attended category.
1090 The misalignment between lower and later layers is starker the larger the distance
1091 between the attended and recorded layers. For example, when looking at layer 12,
1092 attention applied at layer 2 appears to increase and decrease feature map activity
1093 equally, without respect to category preference.

1094 This helps to understand why feature-based attention applied at multiple layers
1095 simultaneously is not particularly effective at enhancing detection performance (Figure
1096 3C). Specifically, if attention at a lower layer decreases the activity of feature maps that
1097 prefer the attended category at a later layer, it is actively counteracting the effects
1098 of attention applied at that layer. In Figure 10A, the effects of applying attention
1099 simultaneously at all layers is shown in black (using the same analysis of Figure 8B. The
1100 results from that figure are also replicated in paler colors for comparison). Attention
1101 is applied at each layer at one-tenth the strength ($\beta = .05$) as when it is applied to
1102 an individual layer. It is clear these effects are not accumulating effectively, as the
1103 activity ratios at the final layer (after passing through 13 layers of $\beta = .05$) are weaker
1104 than effects applied at layer 12 with $\beta = .5$.

1105 Spatial attention, on the other hand, does lead to an effective accumulation of
1106 effects when applied at multiple layers. Figure 10B(left) uses the same analysis as
1107 Figure 9A, and shows the effect of applying spatial attention at all layers (with $\beta =$
1108 $.025$) in black. The effect on the intercept at the tenth layer is equal whether applying
1109 attention at all layers or only at layer 10 with $\beta = .25$. The difference parameter,
1110 however, is more negative when attention is applied at all layers than when attention
1111 is applied at layer 10. This demonstrates something that spatial attention can achieve
1112 at a given layer only when it is applied at a lower one: amplify preferred categories

1113 more than non-preferred. When all activity for all images is scaled multiplicatively
1114 at $l - 1$, some feature maps at layer l may see only a small increase when the image
1115 is of their non-preferred categories, due to the scaling up of their negatively-weighted
1116 inputs. In the cases where this effect is so strong that attention causes a decrease
1117 in activity in response to non-preferred category images (i.e., activity ratio less than
1118 one) while still causing an increase for preferred, attention would have the effect of
1119 sharpening the tuning curve. Tuning curve sharpening as a result of spatial attention
1120 is generally not found experimentally [60, 92].

1121 Activity ratios plotted in Figure 10B(right) are calculated as the activity recorded
1122 from a given quadrant when attention was applied to that quadrant over when no
1123 attention was applied. They are organized according to whether the feature map
1124 prefers or does not prefer the category present in the quadrant. By looking at different
1125 attended and recorded layers, we can see that spatial attention at lower layers can
1126 indeed lead to a higher scaling of feature maps that prefer the presented category, and
1127 that feature maps that do not prefer the presented category can have their activity
1128 decreased due to attention (especially when the gap between attended and recorded
1129 layers is larger).

1130 4. Discussion

1131 In this work, we utilized a deep convolutional neural network (CNN) as a model of
1132 the visual system to probe the relationship between neural activity and performance.
1133 Specifically, we provide a formal mathematical definition of the feature similarity gain
1134 model (FSGM) of attention, the basic tenets of which have been described in several
1135 experimental studies. This formalization allows us to investigate the FSGM's abil-
1136 ity to enhance a CNN's performance on challenging visual tasks. Through this, we
1137 show that neural activity changes matching the type and magnitude of those observed
1138 experimentally can indeed lead to performance changes of the kind and magnitude
1139 observed experimentally. Furthermore, these results hold for a variety of tasks, from
1140 high level category detection to spatial tasks to color classification. The benefit of
1141 these particular activity changes for performance can be analyzed more formally in
1142 a signal detection or Bayesian framework [96, 22, 5, 68, 14], however such analysis is
1143 outside the scope of this work.

1144 A finding from our model is that the layer at which attention is applied can have
1145 a large impact on performance. For detection tasks in particular, attention at early
1146 layers does little to enhance performance while attention at later layers such as 9-
1147 13 is most effective. According to [29], these layers correspond most to areas V4
1148 and LO. Such areas are known and studied for reliably showing attentional effects,
1149 whereas earlier areas such as V1 are generally not [52]. In a study involving detection
1150 of objects in natural scenes, the strength of category-specific preparatory activity in
1151 object selective cortex was correlated with performance, whereas such preparatory
1152 activity in V1 was anti-correlated with performance [71]. This is in line with our
1153 finding that feature-based attention effects at earlier areas can counter the beneficial
1154 effects of that attention at later areas.

1155 While CNNs have representations that are similar to the ventral stream, they lack
1156 many biological details including recurrent connections, dynamics, cell types, and noisy
1157 responses. Preliminary work has shown that these elements can be incorporated into
1158 a CNN structure, and attention can enhance performance in this more biologically-
1159 realistic architecture [49]. Furthermore, while the current work does not include neural

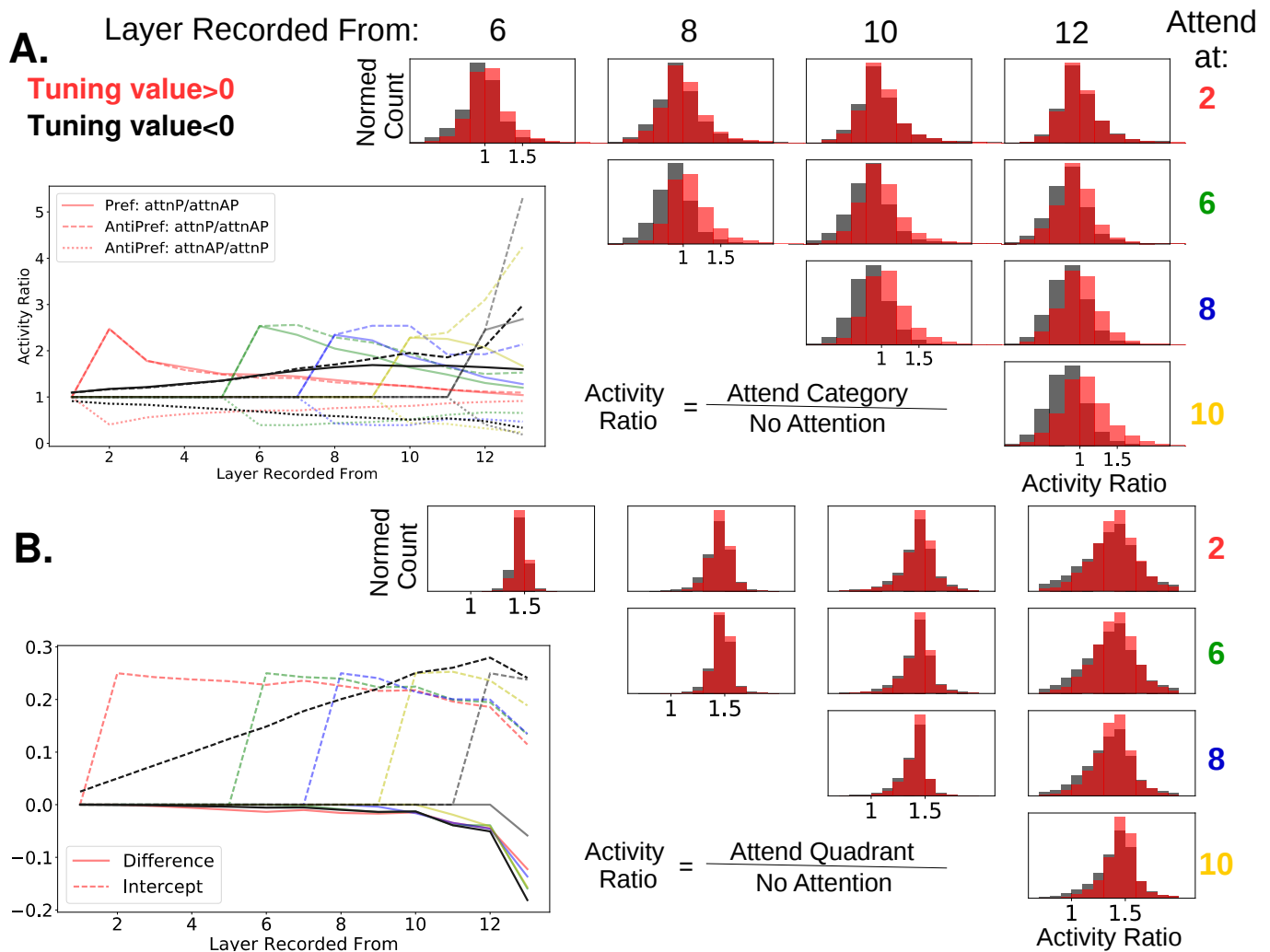


Figure 10: Differences When Applying Attention at All Layers for Feature and Spatial Attention. A.) Feature attention is not enhanced by being applied at multiple layers simultaneously. On the left, activity ratios as described in 8E are reproduced in lighter colors. Black lines show ratios when attention is applied at all layers ($\beta = .05$). On the right activity ratios are shown for when attention is applied at various layers individually and activity is recorded from later layers. In all cases, the category attended was the same as the one present in the input image. Histograms are of ratios of feature map activity when attention is applied to the category divided by activity when no attention is applied, dividing according to whether the feature map prefers (red) or does not prefer (black) the attended category. B.) Attention at multiple layers aides spatial attention. On the left, fit values for lines as described in 9A are shown in paler colors. Black lines are when attention is applied at all layers simultaneously ($\beta = .025$). On the right, histograms of activity ratios are given. Here the activity ratio is activity when attention is applied to the recorded quadrant over when no attention is applied. Feature maps are divided are according whether they prefer (red) or do not prefer (black) the category present in the quadrant.

1160 noise independent of the stimulus, the images used do introduce variable responses.
1161 Take for example, the merged images, wherein a given image from one category is
1162 overlaid with an image from another. This can be thought of as highly structured
1163 noise added to the first image (rather than, for example, pixel-wise Gaussian noise).
1164 Such noise in the signal direction is known to be particularly challenging to overcome
1165 [1].

1166 Another biological detail that this model lacks is "skip connections," when one
1167 layer feeds into both the layer directly above and layers above that. This is seen
1168 frequently in the brain, for example, in connections from V2 to V4 or V4 to parietal
1169 areas [95]. Our results show that the effects on attention at the final convolutional
1170 layer are important for performance changes, suggesting that synaptic distance from
1171 the classifier is a relevant feature—one that is less straight forward to determine in
1172 a network with skip connections. It may be the case though that thinking about
1173 visual areas in terms of their synaptic distance from decision-making areas such as
1174 prefrontal cortex [34] may be more useful for the study of attention than in terms
1175 of their distance from the retina. Finally, a major challenge for understanding the
1176 biological implementation of selective attention is determining how the attention signal
1177 is carried by feedback connections. Feature-based attention in particular appears to
1178 require targeted cell-by-cell modulation, which if implemented directly by top-down
1179 inputs, would require an unrealistic amount of fine tuning. A mechanism wherein
1180 feedback targeting is coarse, but the effects of it are refined by local processing is more
1181 plausible. It may be useful to take inspiration from the machine learning literature on
1182 attention and learning for hypotheses on how the brain does this [101, 47].

1183 While they lack certain biological details, a benefit of using CNNs as a model is
1184 the ability to backpropagate error signals and understand causal relationships. Here
1185 we use this to calculate gradient values that estimate how attention should modulate
1186 activity, and compare these to the tuning values that the FSGM uses. The fact that
1187 these values are correlated and can lead to similar performance changes at task-specific
1188 layers (including similar changes in true and false positive rates, not shown) raises a
1189 question about the nature of biological attention: are neurons really targeted accord-
1190 ing to their tuning, or does the brain use something like gradient values? In [13] the
1191 correlation coefficient between an index of tuning and an index of attentional modula-
1192 tion was .52 for a population of V4 neurons, suggesting factors other than selectivity
1193 influence attention. Furthermore, many attention studies, including that one, use only
1194 preferred and non-preferred stimuli and therefore don't include a thorough investiga-
1195 tion of the relationship between tuning and attentional modulation. [56] use multiple
1196 stimuli to provide support for the FSGM, however the interpretation is limited by
1197 the fact that they only report population averages. Furthermore, those population
1198 averages are closer to the average values in our model when attention is applied ac-
1199 cording to gradient values, rather than tuning values (Figure 8D). [80] investigated the
1200 relationship between tuning strength and the strength of attentional modulation on a
1201 cell-by-cell basis. While they did find a correlation (particularly for binocular disparity
1202 tuning), it wasn't very strong, which leaves room for the possibility that tuning is not
1203 the primary factor that determines attentional modulation.

1204 Another finding from comparing gradient values with tuning values (and doing
1205 "recordings") is that tuning does not always predict how effectively one unit in the
1206 network will impact downstream units or the classifier. In particular, applying at-
1207 tention according to gradient values leads to changes that are hard to interpret when

1208 looked at through the lens of tuning, especially at earlier layers (Figure 8). However
1209 these changes eventually lead to large and impactful changes at later layers. Because
1210 experimenters can easily control the image, defining a cell’s function in terms of how it
1211 responds to stimuli makes practical sense. A recent study looking at the relationship
1212 between tuning and choice probabilities suggests that tuning is not always an indica-
1213 tion of a causal role in classification [103]. Studies that activate specific neurons in
1214 one area and measure changes in another area or in behavioral output will likely be
1215 of significant value for determining function. Thus far, coarse stimulation protocols
1216 have found a relationship between the tuning of neural populations and their impact
1217 on perception [62, 19, 82]. Ultimately though, targeted stimulation protocols and a
1218 more fine-grained understanding of inter-area connections will be needed.

1219 In this study, we used a diversity of attention tasks to see if the same mechanism
1220 could enhance performance universally. While we do find support for the feature simi-
1221 larity gain model’s broad applicability, it is likely the case that the effects of attention
1222 in the brain are influenced substantially by the specifics of the task. Naturally, uni-
1223 modal detection tasks have different challenges than cross-modal readout tasks (such
1224 as detecting a motion change in dots of a certain color). Generally, studies probing
1225 the neural mechanisms of attention care largely about the stimulus that is being at-
1226 tended, and less so about the information the animal needs from that stimulus to do
1227 the task. The task, then, is merely a way to get the subject to attend. However, as we
1228 see in our results, the best attention strategy is dependent on the task. Performance
1229 on our category detection task is only somewhat influenced by the choice of activity
1230 modulation (additive vs. multiplicative, etc), however, performance on the category
1231 classification task depends strongly on the use of multiplicative spatial attention. This
1232 task-dependency is even more stark in the orientation tasks, where the pattern of
1233 performance for attention at different layers is different for the detection and color
1234 classification tasks, even though the attention applied is identical. The effects of at-
1235 tention on firing rates, noise, and correlations may be more similar across studies if
1236 more similar tasks were used.

1237 5. Acknowledgements

1238 We are very grateful to the authors who so readily shared details of their behavioral
1239 data upon request: J. Patrick Mayo, Gary Lupyan, and Mika Koivisto. We further
1240 thank J. Patrick Mayo for helpful comments on the manuscript. GWL was supported
1241 by a Google PhD Fellowship and NIH (T32 NS064929). The authors declare no com-
1242 peting financial interests.

1243 6. References

- 1244 [1] Bruno B Averbeck, Peter E Latham, and Alexandre Pouget. Neural correlations,
1245 population coding and computation. *Nature reviews. Neuroscience*, 7(5):358,
1246 2006.
- 1247 [2] Ji Won Bang and Dobromir Rahnev. Stimulus expectation alters decision crite-
1248 rion but not sensory signal in perceptual decision making. *Scientific reports*, 7
1249 (1):17072, 2017.

- 1250 [3] Jalal K Baruni, Brian Lau, and C Daniel Salzman. Reward expectation differ-
1251 entially modulates attentional behavior and activity in visual area v4. *Nature*
1252 *neuroscience*, 18(11):1656, 2015.
- 1253 [4] Narcisse P Bichot, Matthew T Heard, Ellen M DeGennaro, and Robert Desi-
1254 mone. A source for feature-based attention in the prefrontal cortex. *Neuron*, 88
1255 (4):832–844, 2015.
- 1256 [5] Ali Borji and Laurent Itti. Optimal attentional modulation of a neural popula-
1257 tion. *Frontiers in computational neuroscience*, 8, 2014.
- 1258 [6] Geoffrey M Boynton. A framework for describing the effects of attention on
1259 visual responses. *Vision research*, 49(10):1129–1143, 2009.
- 1260 [7] David A Bridwell and Ramesh Srinivasan. Distinct attention networks for feature
1261 enhancement and suppression in vision. *Psychological science*, 23(10):1151–1158,
1262 2012.
- 1263 [8] Elizabeth A Buffalo, Pascal Fries, Rogier Landman, Hualou Liang, and Robert
1264 Desimone. A backward progression of attentional effects in the ventral stream.
1265 *Proceedings of the National Academy of Sciences*, 107(1):361–365, 2010.
- 1266 [9] Claus Bundesen. A theory of visual attention. *Psychological review*, 97(4):523,
1267 1990.
- 1268 [10] Santiago A Cadena, George H Denfield, Edgar Y Walker, Leon A Gatys, An-
1269 dreas S Tolias, Matthias Bethge, and Alexander S Ecker. Deep convolutional
1270 models improve predictions of macaque v1 responses to natural images. *bioRxiv*,
1271 page 201764, 2017.
- 1272 [11] Marisa Carrasco. Visual attention: The past 25 years. *Vision research*, 51(13):
1273 1484–1525, 2011.
- 1274 [12] Kyle R Cave. The featuregate model of visual selection. *Psychological research*,
1275 62(2):182–194, 1999.
- 1276 [13] Leonardo Chelazzi, John Duncan, Earl K Miller, and Robert Desimone. Re-
1277 sponses of neurons in inferior temporal cortex during memory-guided visual
1278 search. *Journal of neurophysiology*, 80(6):2918–2940, 1998.
- 1279 [14] Sharat Chikkerur, Thomas Serre, Cheston Tan, and Tomaso Poggio. What and
1280 where: A bayesian inference theory of attention. *Vision research*, 50(22):2233–
1281 2247, 2010.
- 1282 [15] Marlene R Cohen and John HR Maunsell. Attention improves performance
1283 primarily by reducing interneuronal correlations. *Nature neuroscience*, 12(12):
1284 1594–1600, 2009.
- 1285 [16] Marlene R Cohen and John HR Maunsell. Using neuronal populations to study
1286 the mechanisms underlying spatial and feature attention. *Neuron*, 70(6):1192–
1287 1204, 2011.

- 1288 [17] Tolga Çukur, Shinji Nishimoto, Alexander G Huth, and Jack L Gallant. At-
1289 tention during natural vision warps semantic representation across the human
1290 brain. *Nature neuroscience*, 16(6):763–770, 2013.
- 1291 [18] Mohammad Reza Daliri, Vladislav Kozyrev, and Stefan Treue. Attention en-
1292 hances stimulus representations in macaque visual cortex without affecting their
1293 signal-to-noise level. *Scientific reports*, 6, 2016.
- 1294 [19] Gregory C DeAngelis, Bruce G Cumming, and William T Newsome. Cortical
1295 area mt and the perception of stereoscopic depth. *Nature*, 394(6694):677, 1998.
- 1296 [20] Rachel N Denison, William T Adler, Marisa Carrasco, and Wei Ji Ma. Humans
1297 flexibly incorporate attention-dependent uncertainty into perceptual decisions
1298 and confidence. *bioRxiv*, page 175075, 2017.
- 1299 [21] Cathryn J Downing. Expectancy and visual-spatial attention: effects on per-
1300 ceptual quality. *Journal of Experimental Psychology: Human perception and*
1301 *performance*, 14(2):188, 1988.
- 1302 [22] Miguel P Eckstein, Matthew F Peterson, Binh T Pham, and Jason A Droll.
1303 Statistical decision theory to relate neurons to behavior in the study of covert
1304 visual attention. *Vision research*, 49(10):1097–1128, 2009.
- 1305 [23] Michael Eickenberg, Alexandre Gramfort, Gaël Varoquaux, and Bertrand
1306 Thirion. Seeing it all: Convolutional network layers map the function of the
1307 human visual system. *NeuroImage*, 152:184–194, 2017.
- 1308 [24] Pascal Fries, John H Reynolds, Alan E Rorie, and Robert Desimone. Modulation
1309 of oscillatory neuronal synchronization by selective visual attention. *Science*, 291
1310 (5508):1560–1563, 2001.
- 1311 [25] Davi Frossard. *VGG in TensorFlow*. Accessed: 2017-03-01.
- 1312 [26] Kunihiko Fukushima. Neocognitron: A hierarchical neural network capable of
1313 visual pattern recognition. *Neural networks*, 1(2):119–130, 1988.
- 1314 [27] Robert Geirhos, David HJ Janssen, Heiko H Schütt, Jonas Rauber, Matthias
1315 Bethge, and Felix A Wichmann. Comparing deep neural networks against
1316 humans: object recognition when the signal gets weaker. *arXiv preprint*
1317 *arXiv:1706.06969*, 2017.
- 1318 [28] Ivan C Griffin and Anna C Nobre. Orienting attention to locations in internal
1319 representations. *Journal of cognitive neuroscience*, 15(8):1176–1194, 2003.
- 1320 [29] Umut Güçlü and Marcel AJ van Gerven. Deep neural networks reveal a gradient
1321 in the complexity of neural representations across the ventral stream. *Journal*
1322 *of Neuroscience*, 35(27):10005–10014, 2015.
- 1323 [30] FH Hamker. The role of feedback connections in task-driven visual search. In
1324 *Connectionist models in cognitive neuroscience*, pages 252–261. Springer, 1999.

- 1325 [31] Fred H Hamker and James Worcester. Object detection in natural scenes by
1326 feedback. In *International Workshop on Biologically Motivated Computer Vision*,
1327 pages 398–407. Springer, 2002.
- 1328 [32] Harold L Hawkins, Steven A Hillyard, Steven J Luck, Mustapha Mouloua,
1329 Cathryn J Downing, and Donald P Woodward. Visual attention modulates sig-
1330 nal detectability. *Journal of Experimental Psychology: Human Perception and*
1331 *Performance*, 16(4):802, 1990.
- 1332 [33] Benjamin Y Hayden and Jack L Gallant. Combined effects of spatial and feature-
1333 based attention on responses of v4 neurons. *Vision research*, 49(10):1182–1187,
1334 2009.
- 1335 [34] Hauke R Heekeren, Sean Marrett, Peter A Bandettini, and Leslie G Ungerleider.
1336 A general mechanism for perceptual decision-making in the human brain. *Nature*,
1337 431(7010):859–862, 2004.
- 1338 [35] Daniel Kaiser, Nikolaas N Oosterhof, and Marius V Peelen. The neural dynamics
1339 of attentional selection in natural scenes. *Journal of neuroscience*, 36(41):10522–
1340 10528, 2016.
- 1341 [36] Kohitij Kar, Jonas Kubilius, Elias Issa, Kailyn Schmidt, and James DiCarlo.
1342 Evidence that feedback is required for object identity inferences computed by
1343 the ventral stream. COSYNE, 2017.
- 1344 [37] Sabine Kastner and Mark A Pinsk. Visual attention as a multilevel selection
1345 process. *Cognitive, Affective, & Behavioral Neuroscience*, 4(4):483–500, 2004.
- 1346 [38] Seyed-Mahdi Khaligh-Razavi and Nikolaus Kriegeskorte. Deep supervised, but
1347 not unsupervised, models may explain it cortical representation. *PLoS compu-*
1348 *tational biology*, 10(11):e1003915, 2014.
- 1349 [39] Seyed-Mahdi Khaligh-Razavi, Linda Henriksson, Kendrick Kay, and Nikolaus
1350 Kriegeskorte. Fixed versus mixed rsa: Explaining visual representations by fixed
1351 and mixed feature sets from shallow and deep computational models. *Journal*
1352 *of Mathematical Psychology*, 76:184–197, 2017.
- 1353 [40] Saeed Reza Kheradpisheh, Masoud Ghodrati, Mohammad Ganjtabesh, and Tim-
1354 othée Masquelier. Deep networks can resemble human feed-forward vision in
1355 invariant object recognition. *Scientific reports*, 6:32672, 2016.
- 1356 [41] Mika Koivisto and Ella Kahila. Top-down preparation modulates visual cate-
1357 gorization but not subjective awareness of objects presented in natural back-
1358 grounds. *Vision Research*, 133:73–80, 2017.
- 1359 [42] Simon Kornblith and Doris Y Tsao. How thoughts arise from sights: inferotem-
1360 poral and prefrontal contributions to vision. *Current Opinion in Neurobiology*,
1361 46:208–218, 2017.
- 1362 [43] Richard J Krauzlis, Lee P Lovejoy, and Alexandre Zénon. Superior colliculus
1363 and visual spatial attention. *Annual review of neuroscience*, 36:165–182, 2013.

- 1364 [44] Jonas Kubilius, Stefania Bracci, and Hans P Op de Beeck. Deep neural networks
1365 as a computational model for human shape sensitivity. *PLoS computational*
1366 *biology*, 12(4):e1004896, 2016.
- 1367 [45] Brenden M Lake, Wojciech Zaremba, Rob Fergus, and Todd M Gureckis. Deep
1368 neural networks predict category typicality ratings for images. In *CogSci*, 2015.
- 1369 [46] Joonyeol Lee and John HR Maunsell. Attentional modulation of mt neurons
1370 with single or multiple stimuli in their receptive fields. *Journal of Neuroscience*,
1371 30(8):3058–3066, 2010.
- 1372 [47] Timothy P Lillicrap, Daniel Cownden, Douglas B Tweed, and Colin J Aker-
1373 man. Random synaptic feedback weights support error backpropagation for
1374 deep learning. *Nature communications*, 7, 2016.
- 1375 [48] Grace W Lindsay. Feature-based attention in convolutional neural networks.
1376 *arXiv preprint arXiv:1511.06408*, 2015.
- 1377 [49] Grace W Lindsay, Dan B Rubin, and Kenneth D Miller. The stabilized supralin-
1378 ear network replicates neural and performance correlates of attention. *COSYNE*,
1379 2017.
- 1380 [50] Drew Linsley, Sven Eberhardt, Tarun Sharma, Pankaj Gupta, and Thomas Serre.
1381 What are the visual features underlying human versus machine vision? In *Pro-*
1382 *ceedings of the IEEE Conference on Computer Vision and Pattern Recognition*,
1383 pages 2706–2714, 2017.
- 1384 [51] Bradley C Love, Olivia Guest, Piotr Slomka, Victor M Navarro, and Edward
1385 Wasserman. Deep networks as models of human and animal categorization. In
1386 *CogSci*, 2017.
- 1387 [52] Steven J Luck, Leonardo Chelazzi, Steven A Hillyard, and Robert Desimone.
1388 Neural mechanisms of spatial selective attention in areas v1, v2, and v4 of
1389 macaque visual cortex. *Journal of neurophysiology*, 77(1):24–42, 1997.
- 1390 [53] Thomas Zhihao Luo and John HR Maunsell. Neuronal modulations in visual
1391 cortex are associated with only one of multiple components of attention. *Neuron*,
1392 86(5):1182–1188, 2015.
- 1393 [54] Gary Lupyan and Michael J Spivey. Making the invisible visible: Verbal but not
1394 visual cues enhance visual detection. *PLoS One*, 5(7):e11452, 2010.
- 1395 [55] Gary Lupyan and Emily J Ward. Language can boost otherwise unseen objects
1396 into visual awareness. *Proceedings of the National Academy of Sciences*, 110(35):
1397 14196–14201, 2013.
- 1398 [56] Julio C Martinez-Trujillo and Stefan Treue. Feature-based attention increases
1399 the selectivity of population responses in primate visual cortex. *Current Biology*,
1400 14(9):744–751, 2004.
- 1401 [57] John HR Maunsell and Erik P Cook. The role of attention in visual processing.
1402 *Philosophical Transactions of the Royal Society of London B: Biological Sciences*,
1403 357(1424):1063–1072, 2002.

- 1404 [58] J Patrick Mayo and John HR Maunsell. Graded neuronal modulations related
1405 to visual spatial attention. *Journal of Neuroscience*, 36(19):5353–5361, 2016.
- 1406 [59] J Patrick Mayo, Marlene R Cohen, and John HR Maunsell. A refined neuronal
1407 population measure of visual attention. *PloS one*, 10(8):e0136570, 2015.
- 1408 [60] Carrie J McAdams and John HR Maunsell. Effects of attention on orientation-
1409 tuning functions of single neurons in macaque cortical area v4. *Journal of Neu-*
1410 *roscience*, 19(1):431–441, 1999.
- 1411 [61] Jude F Mitchell, Kristy A Sundberg, and John H Reynolds. Differential
1412 attention-dependent response modulation across cell classes in macaque visual
1413 area v4. *Neuron*, 55(1):131–141, 2007.
- 1414 [62] Sebastian Moeller, Trinity Crapse, Le Chang, and Doris Y Tsao. The effect of
1415 face patch microstimulation on perception of faces and objects. *Nature Neuro-*
1416 *science*, 20(5):743–752, 2017.
- 1417 [63] Ilya E Monosov, David L Sheinberg, and Kirk G Thompson. The effects of pre-
1418 frontal cortex inactivation on object responses of single neurons in the inferotem-
1419 poral cortex during visual search. *Journal of Neuroscience*, 31(44):15956–15961,
1420 2011.
- 1421 [64] Barbara Montagna, Franco Pestilli, and Marisa Carrasco. Attention trades off
1422 spatial acuity. *Vision research*, 49(7):735–745, 2009.
- 1423 [65] Tirin Moore and Katherine M Armstrong. Selective gating of visual signals by
1424 microstimulation of frontal cortex. *Nature*, 421(6921):370, 2003.
- 1425 [66] Sancho I Moro, Michiel Tolboom, Paul S Khayat, and Pieter R Roelfsema. Neu-
1426 ron activity in the visual cortex reveals the temporal order of cognitive opera-
1427 tions. *Journal of Neuroscience*, 30(48):16293–16303, 2010.
- 1428 [67] Brad C Motter. Neural correlates of feature selective memory and pop-out in
1429 extrastriate area v4. *Journal of Neuroscience*, 14(4):2190–2199, 1994.
- 1430 [68] Vidhya Navalpakkam and Laurent Itti. Search goal tunes visual features opti-
1431 mally. *Neuron*, 53(4):605–617, 2007.
- 1432 [69] Marino Pagan, Luke S Urban, Margot P Wohl, and Nicole C Rust. Signals
1433 in inferotemporal and perirhinal cortex suggest an untangling of visual target
1434 information. *Nature neuroscience*, 16(8):1132–1139, 2013.
- 1435 [70] William K Page and Charles J Duffy. Cortical neuronal responses to optic flow
1436 are shaped by visual strategies for steering. *Cerebral cortex*, 18(4):727–739, 2007.
- 1437 [71] Marius V Peelen and Sabine Kastner. A neural basis for real-world visual search
1438 in human occipitotemporal cortex. *Proceedings of the National Academy of Sci-*
1439 *ences*, 108(29):12125–12130, 2011.
- 1440 [72] Marius V Peelen, Li Fei-Fei, and Sabine Kastner. Neural mechanisms of rapid
1441 natural scene categorization in human visual cortex. *Nature*, 460(7251):94, 2009.

- 1442 [73] Joshua C Peterson, Joshua T Abbott, and Thomas L Griffiths. Adapting
1443 deep network features to capture psychological representations. *arXiv preprint*
1444 *arXiv:1608.02164*, 2016.
- 1445 [74] Dobromir Rahnev, Hakwan Lau, and Floris P de Lange. Prior expectation
1446 modulates the interaction between sensory and prefrontal regions in the human
1447 brain. *Journal of Neuroscience*, 31(29):10741–10748, 2011.
- 1448 [75] Waseem Rawat and Zenghui Wang. Deep convolutional neural networks for
1449 image classification: A comprehensive review. *Neural Computation*, 2017.
- 1450 [76] Maximilian Riesenhuber and Tomaso Poggio. Hierarchical models of object
1451 recognition in cortex. *Nature neuroscience*, 2(11), 1999.
- 1452 [77] Samuel Ritter, David GT Barrett, Adam Santoro, and Matt M Botvinick. Cog-
1453 nitive psychology for deep neural networks: A shape bias case study. *arXiv*
1454 *preprint arXiv:1706.08606*, 2017.
- 1455 [78] Mariel Roberts, Rachel Cymerman, R Theodore Smith, Lynne Kiorpes, and
1456 Marisa Carrasco. Covert spatial attention is functionally intact in amblyopic
1457 human adultsroberts et al. *Journal of vision*, 16(15):30–30, 2016.
- 1458 [79] Edmund T Rolls and Gustavo Deco. Attention in natural scenes: neurophysio-
1459 logical and computational bases. *Neural networks*, 19(9):1383–1394, 2006.
- 1460 [80] Douglas A Ruff and Richard T Born. Feature attention for binocular disparity
1461 in primate area mt depends on tuning strength. *Journal of neurophysiology*, 113
1462 (5):1545–1555, 2015.
- 1463 [81] Melissa Saenz, Giedrius T Buracas, and Geoffrey M Boynton. Global effects of
1464 feature-based attention in human visual cortex. *Nature neuroscience*, 5(7):631,
1465 2002.
- 1466 [82] C Daniel Salzman, Kenneth H Britten, and William T Newsome. Cortical mi-
1467 crostimulation influences perceptual judgements of motion direction. *Nature*,
1468 346(6280):174–177, 1990.
- 1469 [83] K Seeliger, M Fritsche, U Güçlü, S Schoenmakers, J-M Schoffelen, SE Bosch, and
1470 MAJ van Gerven. Cnn-based encoding and decoding of visual object recognition
1471 in space and time. *bioRxiv*, page 118091, 2017.
- 1472 [84] John T Serences, Jens Schwarzbach, Susan M Courtney, Xavier Golay, and
1473 Steven Yantis. Control of object-based attention in human cortex. *Cerebral*
1474 *Cortex*, 14(12):1346–1357, 2004.
- 1475 [85] Thomas Serre, Lior Wolf, Stanley Bileschi, Maximilian Riesenhuber, and Tomaso
1476 Poggio. Robust object recognition with cortex-like mechanisms. *IEEE transac-*
1477 *tions on pattern analysis and machine intelligence*, 29(3):411–426, 2007.
- 1478 [86] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for
1479 large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.

- 1480 [87] Hedva Spitzer, Robert Desimone, Jeffrey Moran, et al. Increased attention en-
1481 hances both behavioral and neuronal performance. *Science*, 240(4850):338–340,
1482 1988.
- 1483 [88] Timo Stein and Marius V Peelen. Object detection in natural scenes: Indepen-
1484 dent effects of spatial and category-based attention. *Attention, Perception, &*
1485 *Psychophysics*, 79(3):738–752, 2017.
- 1486 [89] Jan Theeuwes, Arthur F Kramer, and Paul Atchley. Attentional effects on preat-
1487 tentive vision: spatial precues affect the detection of simple features. *Journal of*
1488 *Experimental Psychology: Human Perception and Performance*, 25(2):341, 1999.
- 1489 [90] Anne M Treisman and Garry Gelade. A feature-integration theory of attention.
1490 *Cognitive psychology*, 12(1):97–136, 1980.
- 1491 [91] Stefan Treue. Neural correlates of attention in primate visual cortex. *Trends in*
1492 *neurosciences*, 24(5):295–300, 2001.
- 1493 [92] Stefan Treue and Julio C Martinez Trujillo. Feature-based attention influences
1494 motion processing gain in macaque visual cortex. *Nature*, 399(6736):575, 1999.
- 1495 [93] Bryan P Tripp. Similarities and differences between stimulus tuning in the
1496 inferotemporal visual cortex and convolutional networks. In *Neural Networks*
1497 *(IJCNN), 2017 International Joint Conference on*, pages 3551–3560. IEEE, 2017.
- 1498 [94] John K Tsotsos, Scan M Culhane, Winky Yan Kei Wai, Yuzhong Lai, Neal Davis,
1499 and Fernando Nufflo. Modeling visual attention via selective tuning. *Artificial*
1500 *intelligence*, 78(1-2):507–545, 1995.
- 1501 [95] Leslie G Ungerleider, Thelma W Galkin, Robert Desimone, and Ricardo Gattass.
1502 Cortical connections of area v4 in the macaque. *Cerebral Cortex*, 18(3):477–499,
1503 2007.
- 1504 [96] Preeti Verghese. Visual search and attention: A signal detection theory ap-
1505 proach. *Neuron*, 31(4):523–535, 2001.
- 1506 [97] Bram-Ernst Verhoef and John HR Maunsell. Attention-related changes in cor-
1507 related neuronal activity arise from normalization mechanisms. *Nature Neuro-*
1508 *science*, 20(7):969–977, 2017.
- 1509 [98] Aurel Wannig, Valia Rodríguez, and Winrich A Freiwald. Attention to surfaces
1510 modulates motion processing in extrastriate area mt. *Neuron*, 54(4):639–651,
1511 2007.
- 1512 [99] Louise Whiteley and Maneesh Sahani. Attention in a bayesian framework. *Fron-*
1513 *tiers in human neuroscience*, 6, 2012.
- 1514 [100] Jeremy M Wolfe. Guided search 2.0 a revised model of visual search. *Psycho-*
1515 *nomic bulletin & review*, 1(2):202–238, 1994.
- 1516 [101] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan
1517 Salakhudinov, Rich Zemel, and Yoshua Bengio. Show, attend and tell: Neural
1518 image caption generation with visual attention. In *International Conference on*
1519 *Machine Learning*, pages 2048–2057, 2015.

- 1520 [102] Daniel LK Yamins, Ha Hong, Charles F Cadieu, Ethan A Solomon, Darren Seib-
1521 ert, and James J DiCarlo. Performance-optimized hierarchical models predict
1522 neural responses in higher visual cortex. *Proceedings of the National Academy*
1523 *of Sciences*, 111(23):8619–8624, 2014.
- 1524 [103] Adam Zaidel, Gregory C DeAngelis, and Dora E Angelaki. Decoupled choice-
1525 driven and stimulus-related activity in parietal neurons may be misrepresented
1526 by choice probabilities. *Nature Communications*, 8, 2017.
- 1527 [104] Weiwei Zhang and Steven J Luck. Feature-based attention modulates feedfor-
1528 ward visual processing. *Nature neuroscience*, 12(1):24–25, 2009.
- 1529 [105] Ying Zhang, Ethan M Meyers, Narcisse P Bichot, Thomas Serre, Tomaso A Pog-
1530 gio, and Robert Desimone. Object decoding with attention in inferior temporal
1531 cortex. *Proceedings of the National Academy of Sciences*, 108(21):8850–8855,
1532 2011.
- 1533 [106] Huihui Zhou and Robert Desimone. Feature-based attention in the frontal eye
1534 field and area v4 during visual search. *Neuron*, 70(6):1205–1217, 2011.