

How biological attention mechanisms improve task performance in a large-scale visual system model

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

^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, Columbia University, New York, New York, 10032, USA*

^c *Mortimer B. Zuckerman Mind Brain Behavior Institute, Department of Neuroscience, Swartz Program in Theoretical Neuroscience, Kavli Institute for Brain Science, College of Physicians and Surgeons, Columbia University, New York, New York, USA*

Abstract

How does attentional modulation of neural activity enhance performance? Here we use a deep convolutional neural network as a large-scale model of the visual system to address this question. We model the feature similarity gain model of attention, in which attentional modulation is applied according to neural stimulus tuning. Using a variety of visual tasks, we show that neural modulations of the kind and magnitude observed experimentally lead to performance changes of the kind and magnitude observed experimentally. We find that, at earlier layers, attention applied according to tuning does not successfully propagate through the network, and has a weaker impact on performance than attention applied according to values computed for optimally modulating higher areas. This raises the question of whether biological attention might be applied at least in part to optimize function rather than strictly according to tuning. We suggest a simple experiment to distinguish these alternatives.

1. Introduction

1 Covert visual attention—applied according to spatial location or visual features—
2 has been shown repeatedly to enhance performance on challenging visual tasks [13].
3 To explore the neural mechanisms behind this enhancement, neural responses to the
4 same 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, 17, 25, 56]. But how do these
7 neural activity changes impact performance? Previous theoretical studies have offered
8 helpful insights on how attention may work to enhance performance [67, 77, 94, 14, 30,
9 99, 29, 23, 7, 98, 11, 90, 97, 16]. However, much of this work is either based on small,
10 hand-designed models or lacks direct mechanistic interpretability. Here, we utilize a
11 large-scale model of the ventral visual stream to explore the extent to which neural
12 changes like those observed experimentally can lead to performance enhancements on
13 realistic visual tasks. Specifically, we use a deep convolutional neural network trained
14 to perform object classification to test effects of the feature similarity gain model of
15 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 to study the
21 biological underpinnings of object recognition for decades [27, 76, 84]. Recently it has
22 been shown that when these networks are trained to successfully perform object classi-
23 fication on real-world images, the intermediate representations learned are remarkably
24 similar to those of the primate visual system, making CNNs state-of-the-art models
25 of the ventral stream [101, 41, 40, 42, 37, 12, 93, 50, 46]. A key finding has been the
26 correspondence between different areas in the ventral stream and layers in the deep
27 CNNs, with early convolutional layers best able to capture the representation of V1
28 and middle and higher layers best able to capture V4 and IT, respectively [28, 24, 82].
29 The generalizability of these networks is limited, however, and the models are not able
30 to match all elements of visual behavior [95, 2, 3]. But given that CNNs can reach
31 near-human performance on some visual tasks and have architectural and represen-
32 tational similarities to the visual system, they are well-positioned for exploring how
33 neural correlates of attention can impact behavior.

34 One popular framework to describe attention’s effects on firing rates is the feature
35 similarity gain model (FSGM). This model, introduced by Treue & Martinez-Trujillo,
36 claims that a neuron’s activity is multiplicatively scaled up (or down) according to
37 how much it prefers (or doesn’t prefer) the properties of the attended stimulus [92,
38 55]. Attention to a certain visual attribute, such as a specific orientation or color,
39 is generally referred to as feature-based attention (FBA). FBA effects are spatially
40 global: if a task performed at one location in the visual field activates attention to a
41 particular feature, neurons that represent that feature across the visual field will be
42 affected [103, 79]. Overall, this leads to a general shift in the representation of the
43 neural population towards that of the attended stimulus [20, 36, 71]. Spatial attention
44 implies that a particular portion of the visual field is being attended. According to the
45 FSGM, spatial location is treated as an attribute like any other. Therefore, a neuron’s
46 modulation due to attention can be predicted by how well its preferred features and
47 spatial receptive field align with the features and location of the attended stimulus.
48 The effects of combined feature and spatial attention have been found to be additive
49 [32].

50 A debated issue in the attention literature is where in the visual stream attention
51 effects can be seen. Many studies of attention focus on V4 and MT/MST [91], as
52 these areas have reliable attentional effects. Some studies do find effects at earlier
53 areas [65], though they tend to be weaker and occur later in the visual response [38].
54 Therefore, a leading hypothesis is that attention signals, coming from prefrontal areas
55 [63, 62, 6, 44], target later visual areas, and the feedback connections that those areas
56 send to earlier ones cause the weaker effects seen there later [10, 51].

57 In this study, we define the FSGM of attention mathematically and implement
58 it in a deep CNN. By applying attention at different layers in the network and for
59 different tasks, we see how neural changes at one area propagate through the network
60 and change performance.

61 **2. Results**

62 The network used in this study—VGG-16, [85]—is shown in Figure 1A and ex-
63 plained in Methods 4.2. Briefly, at each convolutional layer, the application of a given

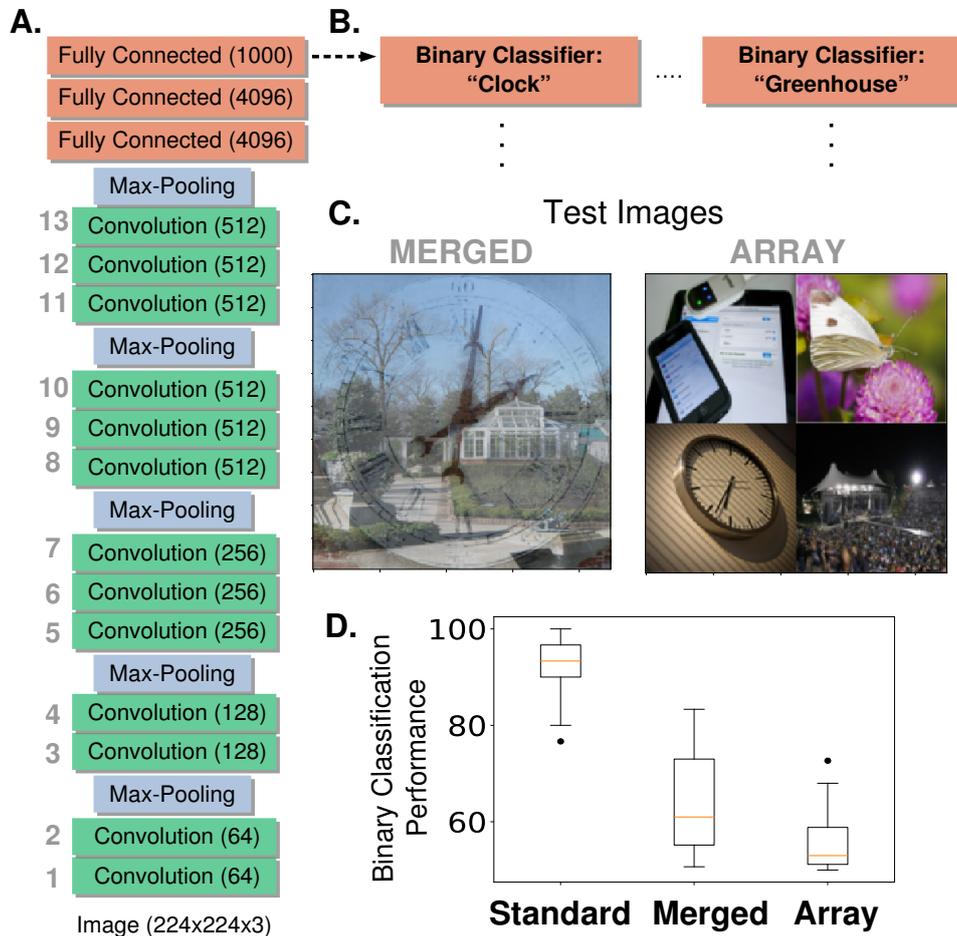


Figure 1: Network Architecture and Feature-Based Attention Task Setup. A.) The model used is a pre-trained deep neural network (VGG-16) that contains 13 convolutional layers (labeled in gray, number of feature maps given in parenthesis) and is trained on the ImageNet dataset to do 1000-way object classification. 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, 20 total). 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 the exception of outliers marked as dots). Standard images are regular ImageNet images not used in the binary classifier training set.

convolutional filter results in a feature map, which is a 2-D grid of artificial neurons that represent how well the bottom-up input at each location aligns with the filter. Therefore a "retinotopic" layout is built into the structure of the network, and the same visual features are represented across that retinotopy (akin to how cells that prefer a given orientation exist at all locations across the V1 retinotopy). This network was explored in [28], where it was shown that early convolutional layers of this CNN are best at predicting activity of voxels in V1, while late convolutional layers are best at predicting activity of voxels in the object-selective lateral occipital area (LO).

2.1. The Relationship between Tuning and Classification

The feature similarity gain model of attention posits that neural activity is modulated by attention in proportion to how strongly a neuron prefers the attended features, as assessed by its tuning. However, the relationship between a neuron's tuning and its ability to influence downstream readouts remains a difficult one to investigate biologically. We use our hierarchical model to explore this question. We do so by using backpropagation to calculate "gradient values", which we compare to tuning curves (see Methods 4.4 and 4.6.1 for details). Gradient values indicate the ways in which feature map activities should change in order to make the network more likely to classify an image as being of a certain object category. Tuning values represent the degree to which the feature map responds preferentially to images of a given category. If there is a correspondence between tuning and classification, a feature map that prefers a given object category (that is, responds strongly to it) should also have a high positive gradient value for that category. In Figure 2A we show gradient values and tuning curves for three example feature maps. In Figure 2C, we show the average correlation coefficients between tuning values and gradient values for all feature maps at each of the 13 convolutional layers. As can be seen, tuning curves in all layers show higher correlation with gradient values than expected by chance (as assayed by shuffled controls), but this correlation is relatively low, increasing across layers from about .2 to .5. Overall tuning quality also increases with layer depth (Figure 2B), but less strongly.

Even at the highest layers, there can be serious discrepancies between tuning and gradient values. In Figure 2D, we show the gradient values of feature maps at the final four convolutional layers, segregated according to tuning value. In red are gradient values that correspond to tuning values greater than one (for example, category 12 for the feature map in the middle pane of Figure 2A). As these distributions show, strong tuning values can be associated with weak or even negative gradient values. Negative gradient values indicate that increasing the activity of that feature map makes the network less likely to categorize the image as the given category. Therefore, even feature maps that strongly prefer a category (and are only a few layers from the classifier) still may not be involved in its classification, or even be inversely related to it. This is aligned with a recent neural network ablation study that shows category selectivity does not predict impact on classification [64].

2.2. Feature-based Attention Improves Performance on Challenging Object Classification Tasks

To determine if manipulation according to tuning values can enhance performance, we created challenging visual images composed of multiple objects for the network to classify. These test images are of two types: merged (two object images transparently overlaid, such as in [83]) or array (four object images arranged on a grid) (see Figure

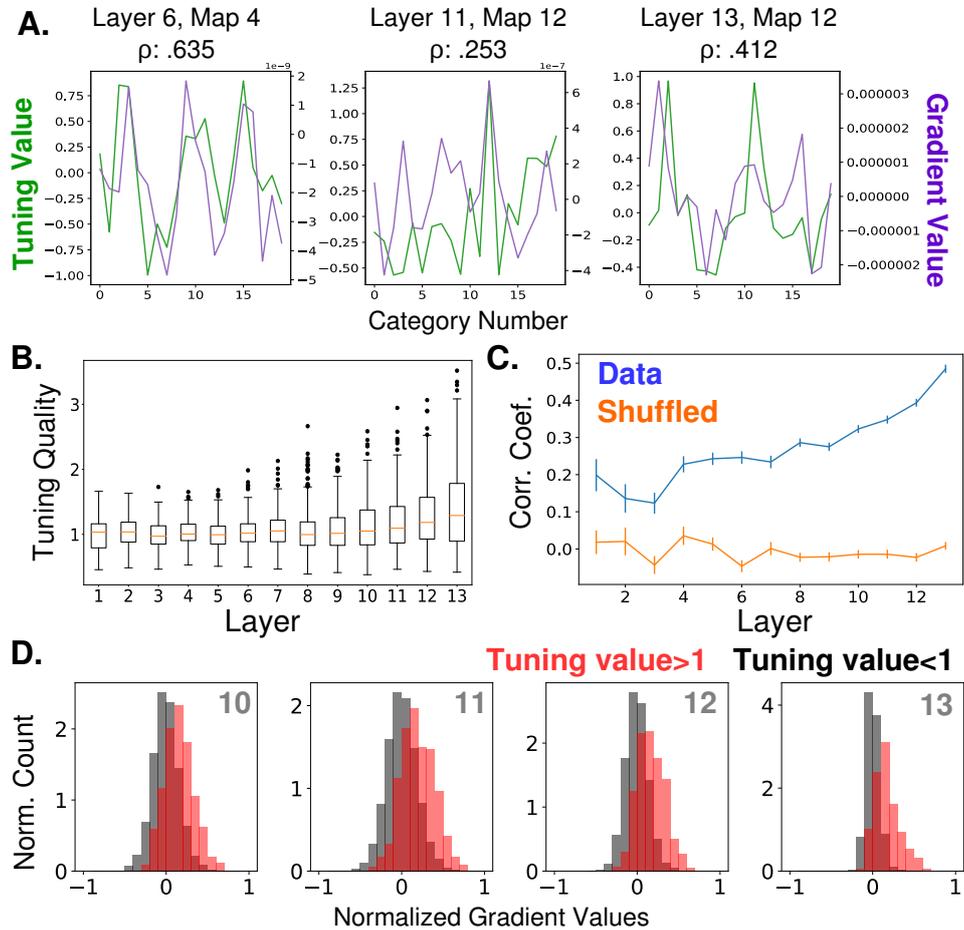


Figure 2: Relationship Between Feature Map Tuning and Gradient Values. A.) Example tuning values (green, left axis) and gradient values (purple, right axis) of three different feature maps from three different layers (identified in titles, layers as labeled in Figure 1A) over the 20 tested object categories. Tuning values indicate how the response to a category differs from the mean response; gradient values indicate how activity should change in order to classify input as from the category. Correlation coefficients between tuning curves and gradient values given in titles. All gradient and tuning values available in Figure 2 Source Data file B.) 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) C.) Correlation coefficients between tuning curves and gradient value curves averaged over feature maps and plotted across layers (errorbars \pm S.E.M., data values in blue and shuffled controls in orange). D.) Distributions of gradient values when tuning is strong. In red, histogram of gradient values associated with tuning values larger than one (i.e. for feature maps that strongly prefer the category), across all feature maps in layers 10, 11, 12, and 13. 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).

110 1C examples). The task for the network is to detect the presence of a given object
111 category in these images. It does so using a series of binary classifiers trained on
112 standard images of these objects, which replace the last layer of the network (Figure
113 1B). The performance of these classifiers on the test images indicates that this is a
114 challenging task for the network (64.4% on merged images and 55.6% on array, Figure
115 1D. Chance is 50%), and thus a good opportunity to see the effects of attention.

116 We implement feature-based attention in this network by modulating the activity
117 of units in each feature map according to how strongly the feature map prefers the
118 attended object category (see Methods 4.6.1 and 4.6). A schematic of this is shown
119 in Figure 3A. The slope of the activation function of units in a given feature map is
120 scaled according to the tuning value of that feature map for the attended category
121 (positive tuning values increase the slope while negative tuning values decrease it).
122 Thus the impact of attention on activity is multiplicative and bi-directional.

123 The effects of attention are measured when attention is applied in this way at each
124 layer individually (Figure 3B; solid lines) or all layers simultaneously (Figure 3: Figure
125 Supplement 1A, red). For both image types (merged and array), attention enhances
126 performance and there is a clear increase in performance enhancement as attention is
127 applied at later layers in the network (numbering is as in Figure 1A). In particular,
128 attention applied at the final convolutional layer performs best, leading to an 18.8%
129 percentage point increase in binary classification on the merged images task and 22.8%
130 increase on the array images task. Thus, FSGM-like effects can have large beneficial
131 impacts on performance.

132 Attention applied at all layers simultaneously does not lead to better performance
133 than attention applied at any individual layer (Figure 3: Figure Supplement 1A). We
134 also performed a control experiment to ensure that nonspecific scaling of activity does
135 not alone enhance performance (Figure 3: Figure Supplement 1C).

136 Some components of the FSGM are debated, e.g. whether attention impacts re-
137 sponses multiplicatively or additively [8, 5, 51, 59], and whether the activity of cells
138 that do not prefer the attended stimulus is actually suppressed [9, 67]. Comparisons
139 of different variants of the FSGM can be seen in Figure 3: Figure Supplement 2. In
140 general, multiplicative and bidirectional effects work best.

141 We also measure performance when attention is applied using gradient values rather
142 than tuning values (these gradient values are calculated to maximize performance
143 on the binary classification task, rather than classify the image as a given category;
144 therefore technically they differ from those shown in Figure 2, however in practice
145 they are strongly correlated. See Methods 4.4 and 4.6.2 for details). Attention applied
146 using gradient values shows the same layer-wise trend as when using tuning values.
147 It also reaches the same performance enhancement peak when attention is applied at
148 the final layers. The major difference, however, comes when attention is applied at
149 middle layers of the network. Here, attention applied according to gradient values
150 outperforms that of tuning values.

151 *2.3. Attention Strength and the Tradeoff between Increasing True and False Positives*

152 In the previous section, we examined the best possible effects of attention by choos-
153 ing the strength for each layer and category that optimized performance. Here, we
154 look at how performance changes as we vary the overall strength (β) of attention.

155 In Figure 4A we break the binary classification performance into true and false
156 positive rates. Here, each colored line indicates a different category and increasing dot

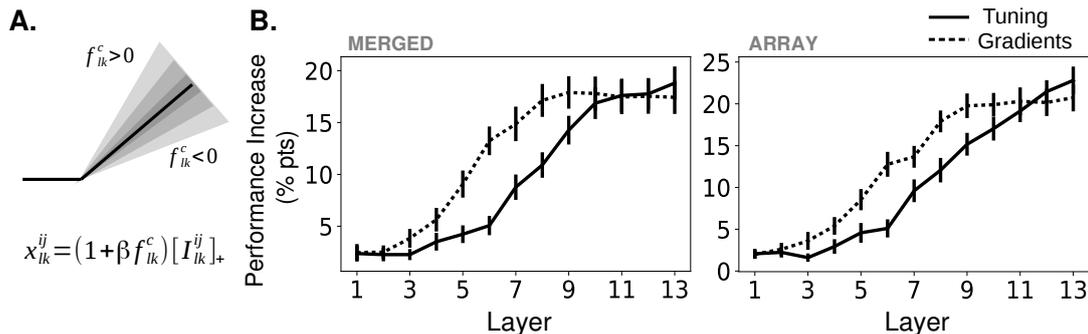


Figure 3: Effects of Applying Feature-Based Attention on Object Category Tasks. A.) Schematic of how attention modulates the activity function. All units in a feature map are modulated the same way. The slope of the activation function is altered based on the tuning (or gradient) value, f_{ik}^c , of a given feature map (here, the k^{th} feature map in the l^{th} layer) for the attended category, c , along with an overall strength parameter β . I_{ik}^{ij} is the input to this unit from the previous layer. For more information, see Methods 4.6. B.) Average increase in binary classification performance as a function of layer attention is applied at (solid line represents using tuning values, dashed line using gradient values, errorbars \pm S.E.M.). In all cases, best performing strength from the range tested is used for each instance. Performance shown separately for merged (left) and array (right) images. Gradients perform significantly ($p < .05$, $N = 20$) better than tuning at layers 5-8 ($p = 4.6e-3, 2.6e-5, 6.5e-3, 4.4e-3$) for merged images and 5-9 ($p = 3.1e-2, 2.3e-4, 4.2e-2, 6.1e-3, 3.1e-2$) for array images. Raw performance values in Figure 3 Source Data file

157 size represents increasing strength of attention. Ideally, true positives would increase
 158 without an equivalent increase (and possibly with a decrease) in false positive rates.
 159 If they increase in tandem, attention does not have a net beneficial effect. Looking at
 160 the effects of applying attention at different layers, we can see that attention at lower
 161 layers is less effective at moving the performance in this space and that movement is in
 162 somewhat random directions, although there is an average increase in performance with
 163 moderate attentional strength. With attention applied at later layers, true positive
 164 rates are more likely to increase for moderate attentional strengths, while substantial
 165 false positive rate increases occur only with higher strengths. Thus, when attention
 166 is applied with modest strength at layer 13, most categories see a substantial increase
 167 in true positives with only modest increases in false positives. As strength continues
 168 to increase however, false positives increase substantially and eventually lead to a net
 169 decrease in overall classifier performance (representing as crossing the dotted line in
 170 Figure 4A).

171 Applying attention according to negated tuning values leads to a decrease in true
 172 and false positive values with increasing attention strength, which decreases overall
 173 performance (Figure 4: Figure Supplement 1A). This verifies that the effects of atten-
 174 tion are not from non-specific changes in activity.

175 Experimentally, when switching from no or neutral attention, neurons in MT
 176 showed an average increase in activity of 7% when attending their preferred motion
 177 direction (and similar decrease when attending the non-preferred) [55]. In our model,
 178 when $\beta = .75$ (roughly the value at which performance peaks at later layers; Figure 4:
 179 Figure Supplement 1B), given the magnitude of the tuning values (average magnitude:
 180 .38), attention scales activity by an average of 28.5%. This value refers to how much
 181 activity is modulated in comparison to the $\beta = 0$ condition, which is probably more
 182 comparable to passive or anesthetized viewing, as task engagement has been shown

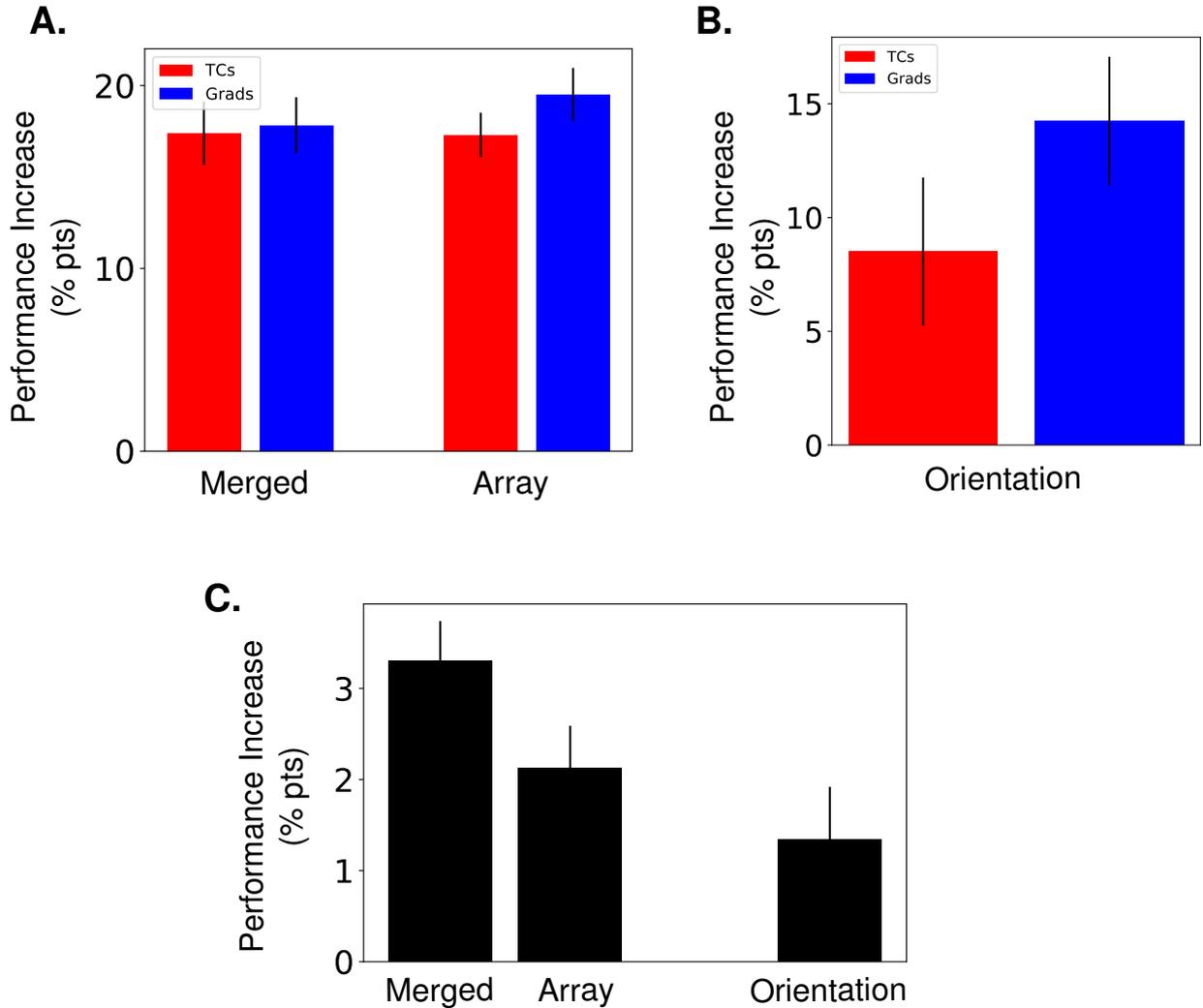


Figure 3: Figure Supplement 1. A.) Effect of applying attention at all layers simultaneously for the category detection task. Performance increase in merged (left) and array (right) image tasks when attention is applied with tuning curves (blue) or gradients (red). Range of strengths tested is one-tenth that of the range tested when applying attention at only one layer and best-performing strength for each category is used. Errorbars are \pm S.E.M. B.) Same as (A) but for orientation detection task (Figure 5A) C.) Control experiment. Instead of using tuning values or gradient values to determine how activity modulates feature maps, all feature maps are scaled by the same amount. Best-performing strengths are used for each category. These results show that merely scaling activity is insufficient to create the performance gains seen when attention is applied in a specific manner. Note: these results are independent of the layer at which the modulation takes place because $[(1+\beta) * I_{ik}^{ij}]_+ = (1+\beta)[I_{ik}^{ij}]_+$ if $(1+\beta) > 0$.

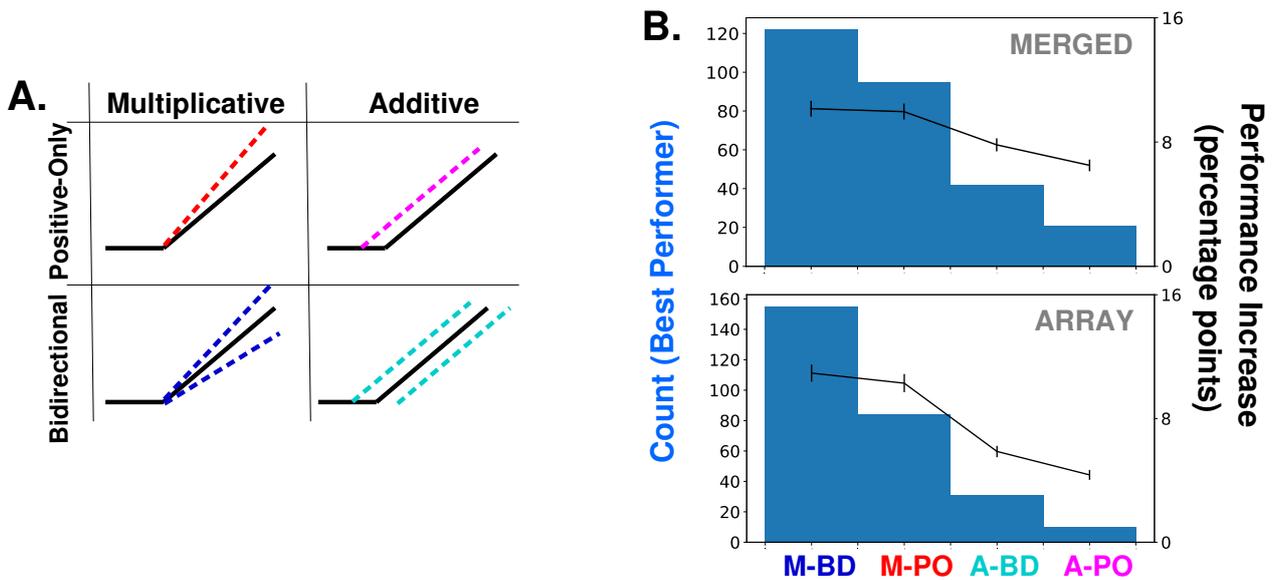


Figure 3: Figure Supplement 2. A.) Schematics of how attention can modulate the activity function. Feature-based attention modulates feature maps according to their tuning values but this modulation can scale the activity multiplicatively or additively, and can either only enhance feature maps that prefer the attended category (positive-only) or also decrease the activity of feature maps that do not prefer it (bidirectional). See Methods 4.6.4 for details of these implementations. The main body of this paper only uses multiplicative bi-directional. B.) Comparison of binary classification performance when attention is applied in each of the four ways described in (A). Considering the combination of attention applied to a given category at a given layer/layers 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 \pm S.E.M.).

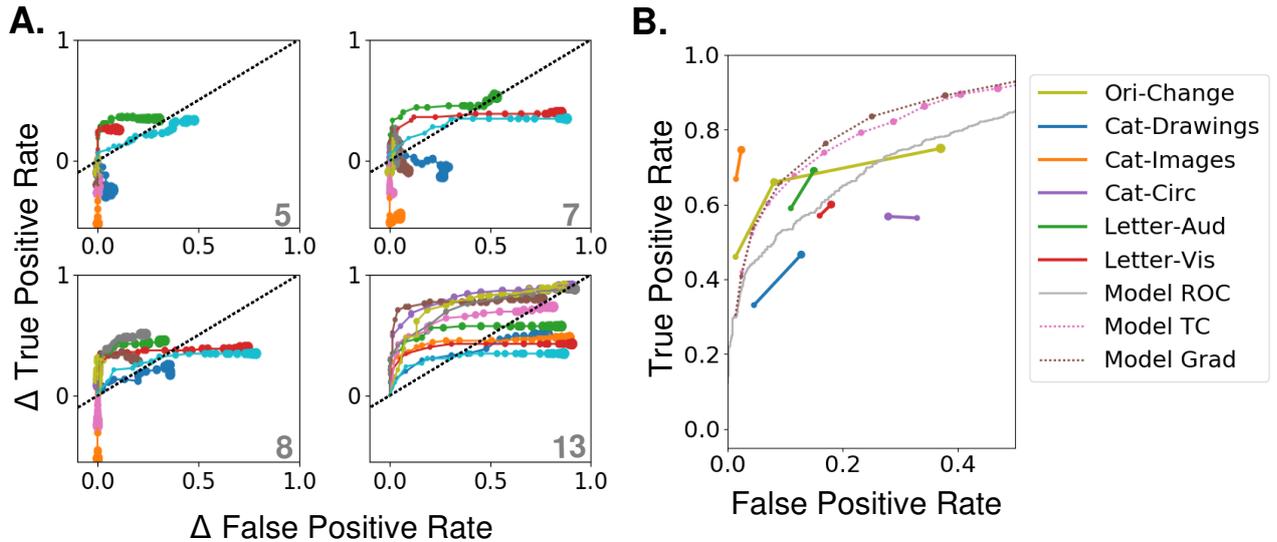


Figure 4: Effects of Varying Attention Strength A.) Effect of increasing attention strength (β) in true and false positive rate space for attention applied at each of four layers (layer indicated in bottom right of each panel, attention applied using tuning values). Each line represents performance for an individual category (only 10 categories shown for visibility), with each increase in dot size representing a .15 increase in β . Baseline (no attention) values are subtracted for each category such that all start at (0,0). The black dotted line represents equal changes in true and false positive rates. B.) Comparisons from experimental data. The true and false positive rates from six experiments in four previously published studies are shown for conditions of increasing attentional strength (solid lines). Cat-Drawings=[54], Exp. 1; Cat-Images=[54],Exp. 2; Objects=[43], Letter-Aud.=[53], Exp. 1; Letter-Vis.=[53], Exp. 2. Ori-Change=[57]. See Methods 4.10 for details of experiments. Dotted lines show model results for merged images, averaged over all 20 categories, when attention is applied using either tuning (TC) or gradient (Grad) values at layer 13. Model results are shown for attention applied with increasing strengths (starting at 0, with each increasing dot size representing a .15 increase in β). Receiver operating curve (ROC) for the model using merged images, which corresponds to the effect of changing the threshold in the final, readout layer, is shown in gray. Raw performance values in Figure 3 Source Data file

183 to scale neural responses generally [70]. This complicates the relationship between
 184 modulation strength in our model and the values reported in the data.

185 To allow for a more direct comparison, in Figure 4B, we collected the true and
 186 false positive rates obtained experimentally during different object detection tasks
 187 (explained in Methods 4.10), and plotted them in comparison to the model results
 188 when attention is applied at layer 13 using tuning values (pink line) or gradient value
 189 (brown line). Five experiments (second through sixth studies) are human studies.
 190 In all of these, uncued trials are those in which no information about the upcoming
 191 visual stimulus is given, and therefore attention strength is assumed to be low. In
 192 cued trials, the to-be-detected category is cued before the presentation of a challenging
 193 visual stimulus, allowing attention to be applied to that object or category.

194 The majority of these experiments show a concurrent increase in both true and false
 195 positive rates as attention strength is increased. The rates in the uncued conditions
 196 (smaller dots) are generally higher than the rates produced by the $\beta = 0$ condition in
 197 our model, consistent with neutrally cued conditions corresponding to $\beta > 0$. We find
 198 (see Methods 4.10), that the average corresponding β value for the neutral conditions
 199 is .37 and for the attended conditions .51. Because attention scales activity by $1 + \beta f_c^{lk}$

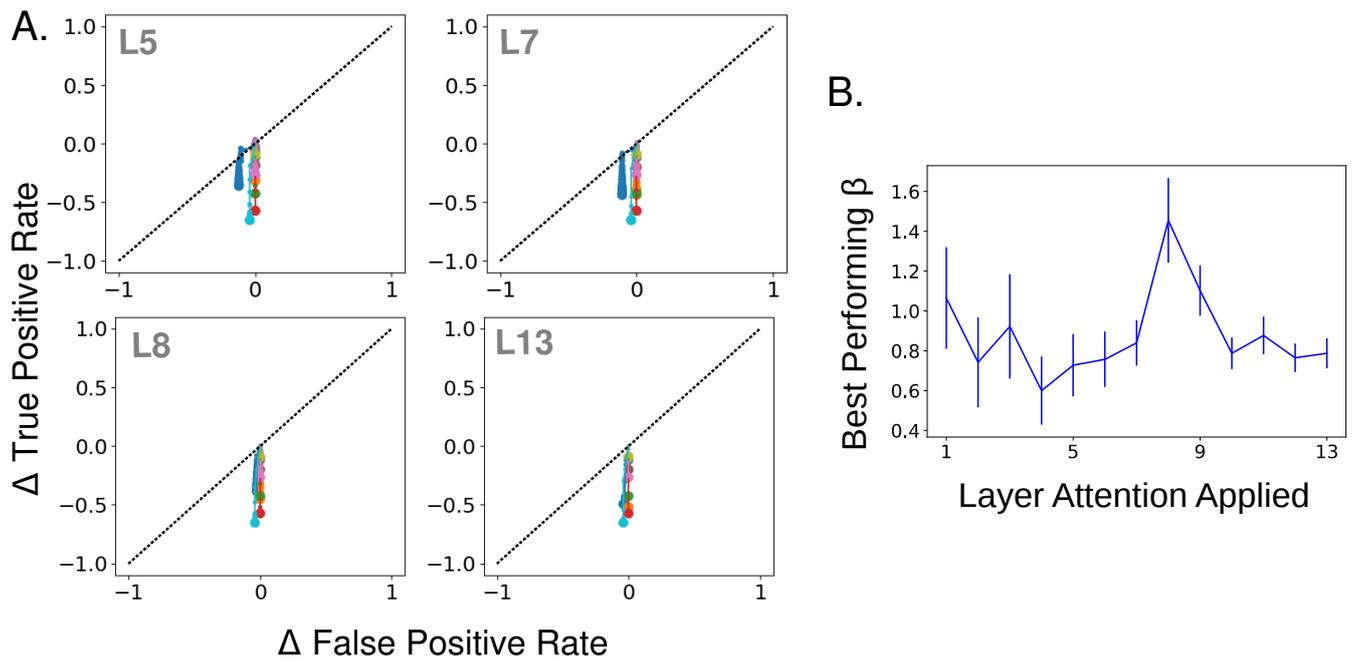


Figure 4: Figure Supplement 1. A.) Effect of strength increase in true and false positive rate space when tuning values are negated. Negated tuning values have the same overall level of positive and negative modulation but in the opposite direction of tuning for a given category. Plot same as in Figure 4A. Layer attention applied at indicated in gray. Attention applied in this way decreases true positives, and to a lesser extent false positives (the initial false positive rate when no attention is applied is very low). B. Mean best performing strength (β value, using regular non-negated attention) across categories as a function of the layer attention is applied at, according to merged images task. Errorbars \pm S.E.M.

200 (where f_c^{lk} is the tuning value), these changes correspond to a $\approx 5\%$ change in activity.

201 The first dataset included in the plot (Ori-Change; yellow line in Figure 4B) comes
202 from a macaque change detection study (see Methods 4.10 for details). Because the
203 attention cue was only 80% valid, attention strength could be of three levels: low
204 (for the uncued stimuli on cued trials), medium (for both stimuli on neutrally-cued
205 trials), or high (for the cued stimuli on cued trials). Like the other studies, this study
206 shows a concurrent increase in both true positive (correct change detection) and false
207 positive (premature response) rates with increasing attention strength. For the model
208 to achieve the performance changes observed between low and medium attention a
209 roughly 12% activity change is needed, but average V4 firing rates recorded during
210 this task show an increase of only 3.6%. This discrepancy may suggest that changes
211 in correlations [17] or firing rate changes in areas aside from V4 also make important
212 contributions to observed performance changes.

213 Thus, according to our model, the size of experimentally observed performance
214 changes is broadly consistent with the size of experimentally observed neural changes.
215 While other factors are likely also relevant for performance changes, this rough align-
216 ment between the magnitude of firing rate changes and magnitude of performance
217 changes supports the idea that the former could be a major causal factor for the lat-
218 ter. In addition, the fact that the model can capture this relationship provides further
219 support for its usefulness as a model of the biology.

220 Finally, we show the change in true and false positive rates when the threshold of
221 the final layer binary classifier is varied (a "receiver operating characteristic" analy-
222 sis, Figure 4B, gray line; no attention was applied during this analysis). Comparing
223 this to the pink line, it is clear that varying the strength of attention applied at the
224 final convolutional layer has more favorable performance effects than altering the clas-
225 sifier threshold (which corresponds to an additive effect of attention at the classifier
226 layer). This points to the limitations that could come from attention targeting only
227 downstream readout areas.

228 Overall, the model roughly matches experiments in the amount of neural modula-
229 tion needed to create the observed changes in true and false positive rates. However,
230 it is clear that the details of the experimental setup are relevant, and changes aside
231 from firing rate and/or outside the ventral stream also likely play a role [67].

232 *2.4. Feature-based Attention Enhances Performance on Orientation Detection Task*

233 Some of the results presented above, particularly those related to the layer at which
234 attention is applied, may be influenced by the fact that we are using an object catego-
235 rization task. To see if results are comparable using the simpler stimuli frequently used
236 in macaque studies, we created an orientation detection task (Figure 5A). Here, binary
237 classifiers trained on full-field oriented gratings are tested using images that contain
238 two gratings of different orientation and color. The performance of these binary clas-
239 sifiers without attention is above chance (distribution across orientations shown in
240 inset of Figure 5A). The performance of the binary classifier associated with vertical
241 orientation (0 degrees) was abnormally high (92% correct without attention, other ori-
242 entations average 60.25%. This likely reflects the over-representation of vertical lines
243 in the training images) and this orientation was excluded from further performance
244 analysis.

245 Attention is applied according to orientation tuning values of the feature maps
246 (tuning quality by layer is shown in Figure 5B) and tested across layers. We find

247 (Figure 5D, solid line and Figure 3: Figure Supplement 1B, red) that the trend in
248 this task is similar to that of the object task: applying attention at later layers leads
249 to larger performance increases (14.4% percentage point increase at layer 10). This is
250 despite the fact that orientation tuning quality peaks in the middle layers.

251 We also calculate the gradient values for this orientation detection task. While
252 overall the correlations between gradient values and tuning values are lower (and even
253 negative for early layers), the average correlation still increases with layer (Figure
254 5C), as with the category detection task. Importantly, while this trend in correlation
255 exists in both detection tasks tested here, it is not a universal feature of the network
256 or an artifact of how these values are calculated. Indeed, an opposite pattern in
257 the correlation between orientation tuning and gradient values is shown when using
258 attention to orientation to classify the color of a stimulus with the attended orientation
259 (Figure 7B, Methods 4.5 and 4.6.2).

260 The results of applying attention according to gradient values is shown in Figure
261 5D (dashed line). Here again, using gradient value creates similar trends as using
262 tuning values, with gradient values performing better in the middle layers.

263 *2.5. Feature-based Attention Primarily Influences Criteria and Spatial Attention Pri-* 264 *marily Influences Sensitivity*

265 Signal detection theory is frequently used to characterize the effects of attention
266 on performance [97]. Here, we use a joint feature-spatial attention task to explore
267 effects of attention in the model. The task uses the same two-grating stimuli described
268 above. The same binary orientation classifiers are used and the task of the model
269 is to determine if a given orientation is present in a given quadrant of the image.
270 Performance is then measured when attention is applied to an orientation, a quadrant,
271 or both an orientation and a quadrant (effects are combined additively, for more, see
272 Methods 4.6). Two key signal detection measurements are computed: criteria and
273 sensitivity. Criteria is a measure of the threshold that’s used to mark an input as
274 positive, with a higher criteria leading to fewer positives; sensitivity is a measure of
275 the separation between the two populations (positives and negatives), with higher
276 sensitivity indicating a greater separation.

277 Figure 5E shows that both spatial and feature-based attention influence sensitivity
278 and criteria. However, feature-based attention decreases criteria more than spatial
279 attention does. Intuitively, feature-based attention shifts the representations of all
280 stimuli in the direction of the attended category, implicitly lowering the detection
281 threshold. Starting from a high threshold, this can lead to the observed behavioral
282 pattern wherein true positives increase before false positives do. Sensitivity increases
283 more for spatial attention alone than for feature-based attention alone, indicating that
284 spatial attention amplifies differences in the representation of whichever features are
285 present. These general trends hold regardless of the layer at which attention is ap-
286 plied and whether feature-based attention is applied using tuning curves or gradients.
287 Changes in true and false positive rates for this task can be seen explicitly in Figure
288 5: Figure Supplement 1.

289 In line with our results, spatial attention was found experimentally to increase sensi-
290 tivity and (less reliably) decrease criteria [31, 22]. Furthermore, feature-based attention
291 is known to decrease criteria, with lesser effects on sensitivity ([74, 4], though see [87]).
292 A study that looked explicitly at the different effects of spatial and category-based at-
293 tention [88] found that spatial attention increases sensitivity more than category-based

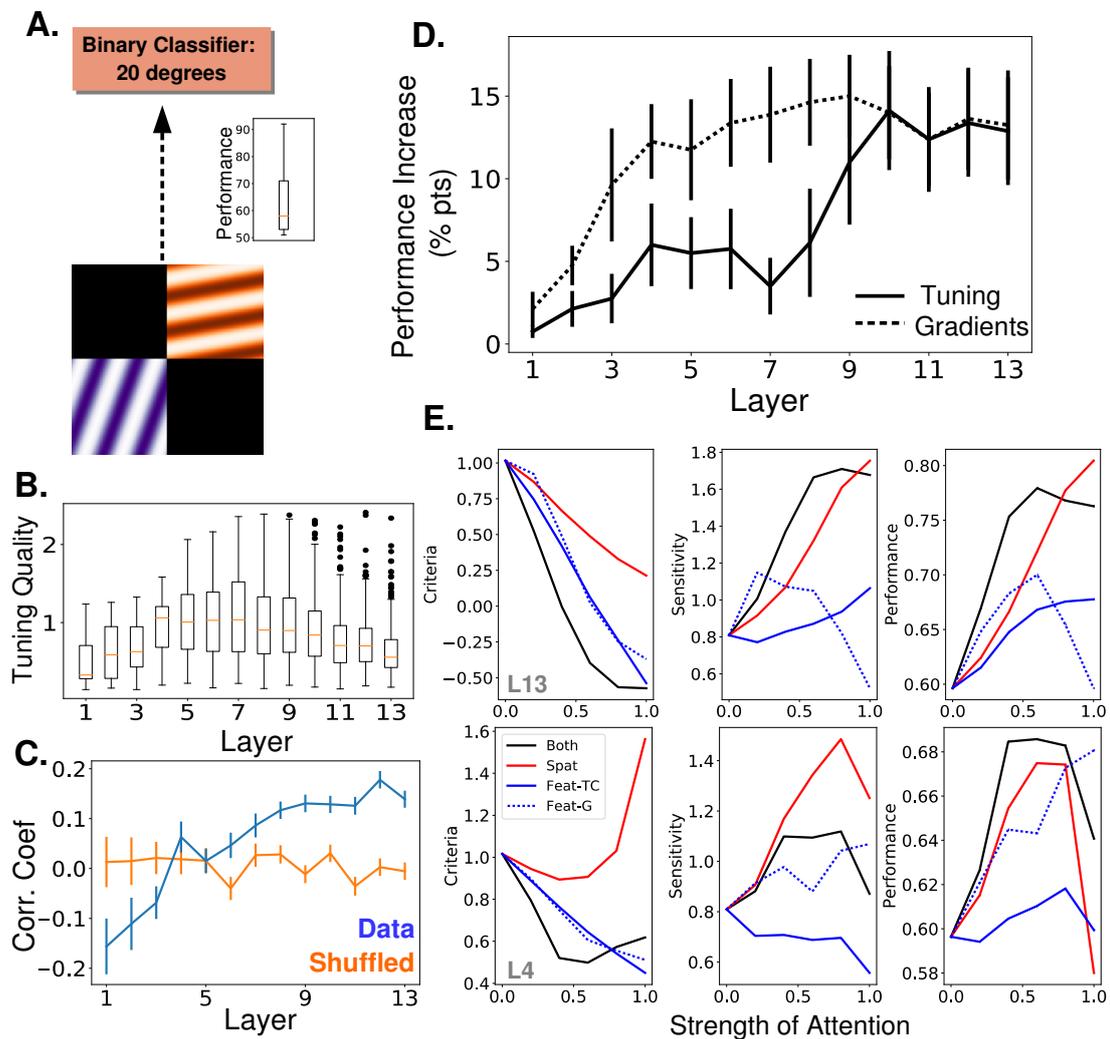


Figure 5: Attention Task and Results Using Oriented Gratings. A.) Orientation detection task. Like with the object category detection tasks, separate binary classifiers trained to detect each of 9 different orientations replaced the final layer of the network. Test images included 2 oriented gratings of different color and orientation located at 2 of 4 quadrants. Inset shows performance over 9 orientations without attention B.) Orientation tuning quality as a function of layer. C.) Average correlation coefficient between orientation tuning curves and gradient curves across layers (blue). Shuffled correlation values in orange. Errorbars are \pm S.E.M. D.) Comparison of performance on orientation detection task when attention is determined by tuning values (solid line) or gradient values (dashed line) and applied at different layers. As in Figure 3B, best performing strength is used in all cases. Errorbars are \pm S.E.M. Gradients perform significantly ($p = 1.9e - 2$) better than tuning at layer 7. Raw performance values available in Figure 5 Source Data-1 file. E.) Change in signal detection values and performance (percent correct) when attention is applied in different ways—spatial (red), feature according to tuning (solid blue), feature according to gradients (dashed blue), and both spatial and feature (according to tuning, black)—for the task of detecting a given orientation in a given quadrant. Top row is when attention is applied at layer 13 and bottom when applied at layer 4. Raw performance values available in Figure 5 Source Data-2 file.

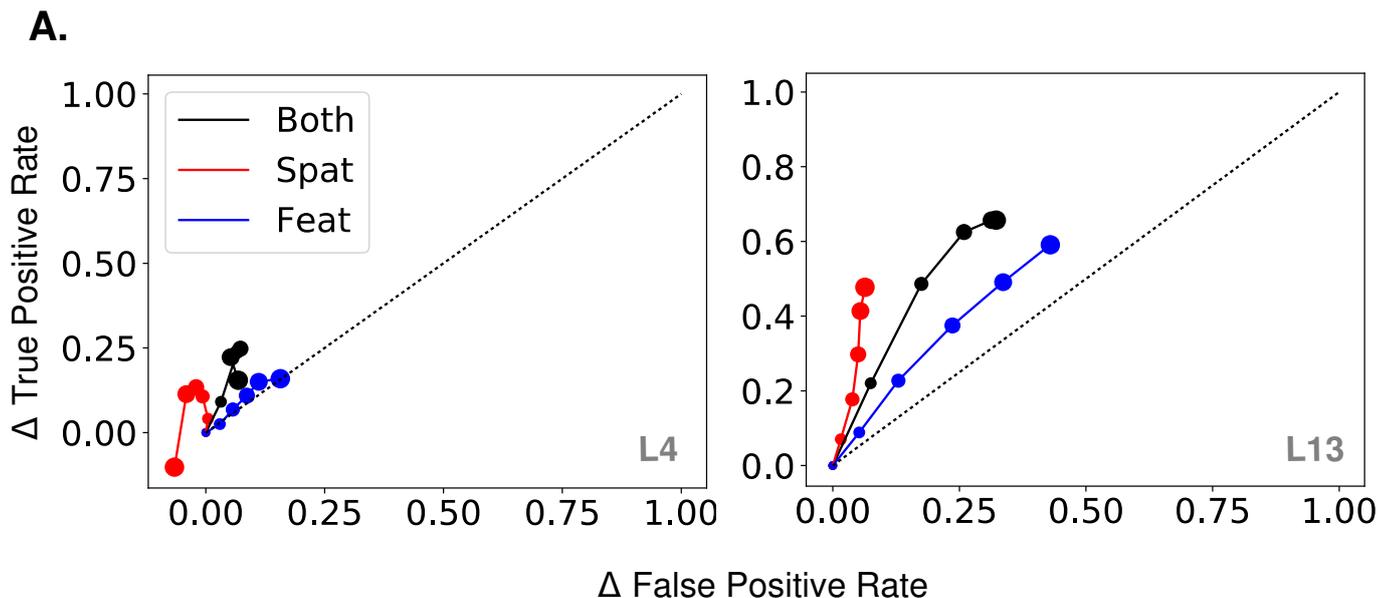


Figure 5: Figure Supplement 1. A.) Effect of strength increase in true and false positive rate space when attention is applied according to quadrant, orientation, or both in the orientation detection task. Rates averaged over orientations/locations. Increasing dot size corresponds to .2 increase in β each. No-attention rates are subtracted and the black dotted line indicates equal increase in true and false positives. Layer attention applied at indicated in gray.

294 attention (most visible in their Experiment 3c, which uses natural images), and the
 295 effects of the two are additive.

296 Attention and priming are known to impact neural activity beyond pure sensory
 297 areas [45, 19]. This idea is borne out by a study that aimed to isolate the neural
 298 changes associated with sensitivity and criteria changes [52]. In this study, the authors
 299 designed behavioral tasks that encouraged changes in behavioral sensitivity or criteria
 300 exclusively: high sensitivity was encouraged by associating a given stimulus location
 301 with higher overall reward, while high criteria was encouraged by rewarding correct
 302 rejects more than hits (and vice versa for low sensitivity/criteria). Differences in V4
 303 neural activity were observed between trials using high versus low sensitivity stimuli.
 304 No differences were observed between trials using high versus low criteria stimuli.
 305 This indicates that areas outside of the ventral stream (or at least outside V4) are
 306 capable of impacting criteria [86]. Importantly, it does not mean that changes in V4
 307 don't impact criteria, but merely that those changes can be countered by the impact
 308 of changes in other areas. Indeed, to create sessions wherein sensitivity was varied
 309 without any change in criteria, the authors had to increase the relative correct reject
 310 reward (i.e., increase the criteria) at locations of high absolute reward, which may have
 311 been needed to counter a decrease in criteria induced by attention-related changes in
 312 V4 (similarly, they had to decrease the correct reject reward at low reward locations).
 313 Our model demonstrates clearly how such effects from sensory areas alone can impact
 314 detection performance, which, in turn highlights the role downstream areas may play
 315 in determining the final behavioral outcome.

2.6. Recordings Show How Feature Similarity Gain Effects Propagate

To explore how attention applied at one location in the network impacts activity later on, we apply attention at various layers and "record" activity at others (Figure 6A, in response to full field oriented gratings). In particular, we record activity of feature maps at all layers while applying attention at layers 2, 6, 8, 10, or 12 individually.

To understand the activity changes occurring at each layer, we use an analysis from [55] that was designed to test for FSGM-like effects and is explained in Figure 6B. Here, the activity of a feature map in response to a given orientation when attention is applied is divided by the activity in response to the same orientation without attention. These ratios are organized according to the feature map's orientation preference (most to least) and a line is fit to them. According to the FSGM of attention, this ratio should be greater than one for more preferred orientations and less than one for less preferred, creating a line with an intercept greater than one and negative slope.

In Figure 6C, we plot the median value of the slopes and intercepts across all feature maps at a layer, when attention is applied at different layers (indicated by color). When attention is applied directly at a layer according to its tuning values (left), FSGM effects are seen by default (intercept values are plotted in terms of how they differ from one; comparable average values from [55] are intercept: .06 and slope: 0.0166, but we are using $\beta = 0$ for the no-attention condition in the model which, as mentioned earlier, is not necessarily the best analogue for no-attention conditions experimentally. Therefore we use these measures to show qualitative effects). As these activity changes propagate through the network, however, the FSGM effects wear off, suggesting that activating units tuned for a stimulus at one layer does not necessarily activate cells tuned for that stimulus at the next. This misalignment between tuning at one layer and the next explains why attention applied at all layers simultaneously isn't more effective (Figure 3: Figure Supplement 1). In fact, applying attention to a category at one layer can actually have effects that counteract attention at a later layer (see Figure 6: Figure Supplement 1).

In Figure 6C (right), we show the same analysis, but while applying attention according to gradient values. The effects at the layer at which attention is applied do not look strongly like FSGM, however FSGM properties evolve as the activity changes propagate through the network, leading to clear FSGM-like effects at the final layer. Finding FSGM-like behavior in neural data could thus be a result of FSGM effects at that area or non-FSGM effects at an earlier area (here, attention applied according to gradients which, especially at earlier layers, are not aligned with tuning).

An alternative model of the neural effects of attention—the feature matching (FM) model—suggests that the effect of attention is to amplify the activity of a neuron whenever the stimulus in its receptive field matches the attended stimulus. In Figure 6D, we calculate the fraction of feature maps at a given layer that show feature matching behavior (defined as having activity ratios greater than one when the stimulus orientation matches the attended orientation for both preferred and anti-preferred orientations). As early as one layer post-attention, some feature maps start showing feature matching behavior. The fact that the attention literature contains conflicting findings regarding the feature similarity gain model versus the feature matching model [66, 78] may result from this finding that FSGM effects can turn into FM effects as they propagate through the network. In particular, this mechanism can explain the observations that feature matching behavior is observed more in FEF than V4 [105] and that match information is more easily read out from perirhinal cortex than IT

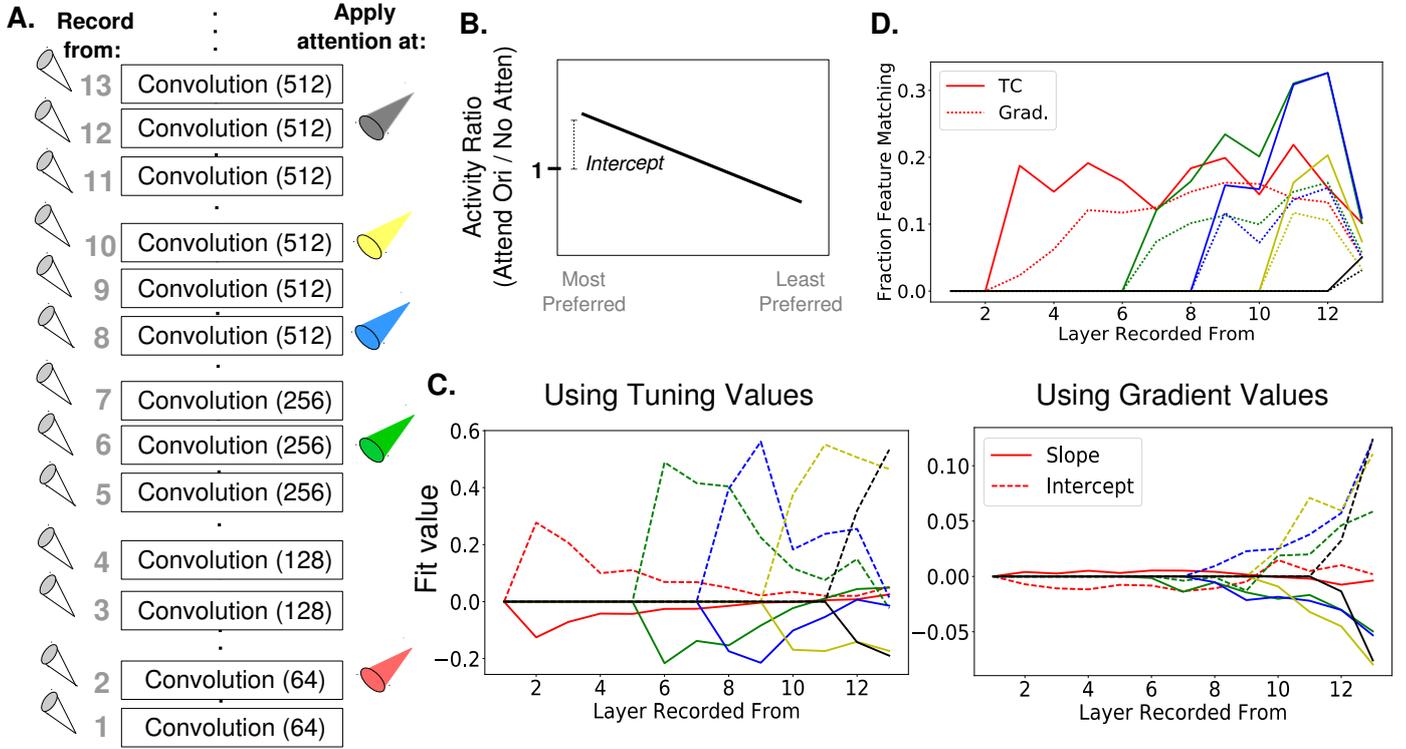


Figure 6: How Attention-Induced Activity Changes Propagate through the Network. A.) Recording setup. The spatially averaged activity of feature maps at each layer was recorded (left) while attention was applied at layers 2, 6, 8, 10, or 12 individually. Activity was in response to a full field oriented grating. B.) Schematic of metric used to test for the feature similarity gain model. Activity when a given orientation is present and attended is divided by the activity when no attention is applied, giving 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 in (C). Intercept values are plotted in terms of how they differ from 1, so positive values are an intercept greater than 1. (FSGM predicts negative slope and positive intercept) C.) The median slope (solid line) and intercept (dashed line) values as described in (B) plotted for each layer when attention is applied to the layer indicated by the line color as labeled in (A). On the left, attention applied according to tuning values and on the right, attention applied according to gradient values. Raw slope and intercept values when using tuning curves available in Figure 6 Source Data-1 file and for gradients in Figure 6 Source Data-2 file. D.) Fraction of feature maps displaying feature matching behavior at each layer when attention is applied at the layer indicated by line color. Shown for attention applied according to tuning (solid lines) and gradient values (dashed line).

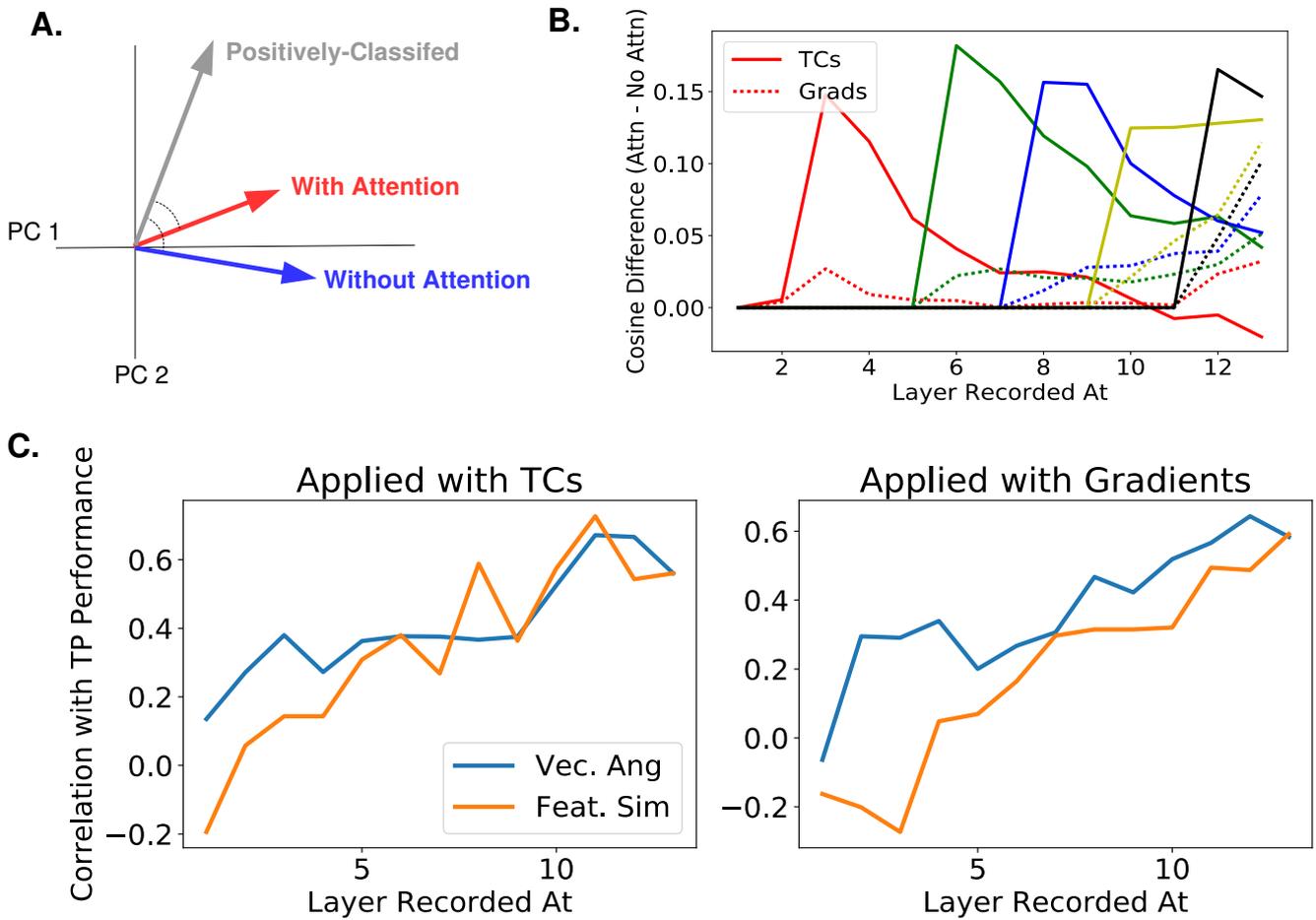


Figure 6: Figure Supplement 2. A.) A new measure of activity changes inspired by gradient values. The gray vector represents the average pattern of neural activity in response to images the classifier indicates as containing the given orientation (i.e., positively-classified in the absence of attention, whether or not the orientation was present in the image). The blue vector (activity without attention) and red vector (activity when attention is applied) are then made using images that do contain the given orientation. Assuming that attention makes activity look more like activity during positive classification, this measure compares the cosine of angle between the positively-classified and with-attention vectors to the cosine of the angle between the positively-classified and without-attention vectors. We use $\cos(\theta)$ as the measure, but results are similar using θ . B.) Using the same color scheme as Figure 6, this plot shows how attention applied at different layers causes activity changes throughout the network, as measured by the vector method introduced in (A). Specifically, the cosine of the angle between the positively-classified and without-attention vectors is subtracted from the cosine of the angle between the positively-classified and with-attention vectors. Solid lines indicate median value of this difference (across images) when attention is applied with tuning curves and dashed line when applied with gradients. C.) How activity changes correlate with performance changes. The correlation coefficient between the change in true positive rate with attention and activity changes as measured by: difference in cosines of angles (blue line) or feature similarity gain model-like behavior (orange line, see Methods 4.9 for how this is calculated). Activity and performance changes are collected when attention is applied at different layers individually (using a range of strengths) according to tuning curves (left) or gradient values (right). Activity is recorded at and after the layer at which attention is applied. For a given layer L , the correlation coefficient is thus computed across data points, where there is one data point for each combination of orientation, strength of attention applied, and layer ($1 \leq L$) at which attention is applied. A bootstrap analysis determined that at layers 1, 2, 3, 4, and 5 the vector angle method had significantly ($p < .05$) higher correlation with performance for both application options than the FSGM measure.

364 [69].

365 We also investigated the extent to which measures of attention’s neural effects
366 correlate with changes in performance (see Methods 4.9). For this we developed a
367 new, experimentally-feasible way of calculating attention’s effects on neural activity
368 that is inspired by the gradient-based approach to attention (that is, it focuses on
369 classification rather than tuning). We show (Figure 6: Figure Supplement 2) that this
370 new measure better correlates with performance changes than the FSGM measure of
371 activity changes, particularly at earlier layers.

372 There is a simple experiment that would distinguish whether factors beyond tuning,
373 such as gradients, play a role in guiding attention. It requires using two tasks with
374 very different objectives (which should produce different gradients) but with the same
375 attentional cue. An example is described in Figure 7. Here, the two tasks used would
376 be an orientation-based color classification task (two gratings each with their own
377 color and orientation are simultaneously shown, and the task is to report the color of
378 the grating with the attended orientation) and an orientation detection task (report if
379 the attended orientation is present or absent in the image). In both cases, attention
380 is cued according to orientation. But gradient-based attention will produce different
381 neural modulations for the two tasks, while the FSGM predicts identical modulations
382 (Figure 7C). Thus, an experiment that recorded from the same neurons during both
383 tasks could distinguish between tuning-based and gradient-based attention.

384 3. Discussion

385 In this work, we utilized a deep convolutional neural network (CNN) as a model of
386 the visual system to probe the relationship between modulation of neural activity, as in
387 attention, and performance. Specifically, we formally define the feature similarity gain
388 model (FSGM) of attention (the basic tenets of which have been described in several
389 experimental studies) as a multiplicative modulation of neuronal activity proportional
390 to the neuron’s mean-subtracted feature tuning. This formalization allows us to in-
391 vestigate the FSGM’s ability to enhance a CNN’s performance on challenging visual
392 tasks. We found that, across a variety of tasks, neural activity changes matching the
393 type and magnitude of those observed experimentally can indeed lead to performance
394 changes of the kind and magnitude observed experimentally.

395 We used the full observability of the model to investigate the relationship between
396 tuning and function. We compared attention applied according to feature tuning
397 (the FSGM) with attention designed to optimally modulate activity to improve per-
398 formance (as determined by the gradient of performance with respect to the neural
399 activity). Attention applied according to tuning does not successfully propagate from
400 lower or middle to higher layers; that is, enhancing the activity of neurons that most
401 prefer a given category at lower layers need not selectively enhance the activity of neu-
402 rons preferring that category at higher layers. As a result, attention applied according
403 to the FSGM performs poorly when applied at early to middle layers, while attention
404 applied according to gradients at these layers performs better.

405 Attention is most effective applied at later layers (e.g., layers 9-13), where tuning
406 and gradient values are better correlated. According to [28], these layers correspond
407 most to areas V4 and LO. Such areas are known and studied for reliably showing
408 attentional effects, whereas earlier areas such as V1 are generally not [51, 1]. In
409 a study involving detection of objects in natural scenes, the strength of category-
410 specific preparatory activity in object selective cortex was correlated with performance,

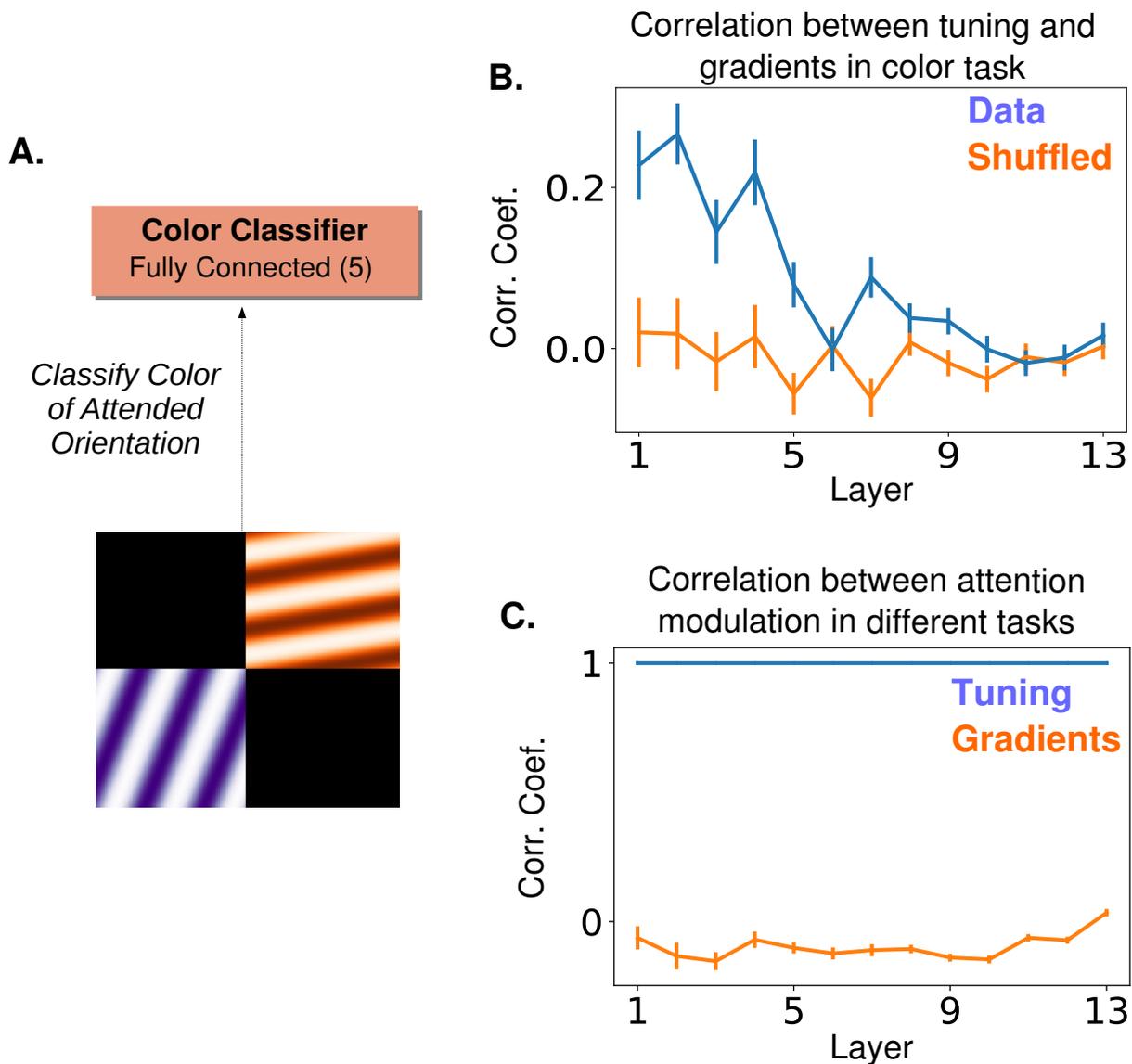


Figure 7: A Proposed Experiment to Distinguish between Tuning-based and Gradient-based Attention A.) "Cross-featural" attention task. Here, the final layer of the network is replaced with a color classifier and the task is to classify the color of the attended orientation in a two-orientation stimulus. Importantly, in both this and the orientation detection task (Figure 5A), a subject performing the task would be cued to attend to an orientation. B.) The correlation coefficient between the gradient values calculated for this task and orientation tuning values (as in Figure 5C). Correlation peaks at lower layers for this task. C.) Correlation between tuning values for the two tasks (blue) and between gradient values for the two tasks (orange). If attention does target cells based on tuning, the modulation would be the same in both the color classification task and the orientation detection task. If a gradient-based targeting is used, no (or even a slight anti-) correlation is expected. Tuning and gradient values available in Figure 7 Source Data file.

411 whereas such preparatory activity in V1 was anti-correlated with performance [71].
412 This is in line with our finding that feature-based attention effects at earlier areas
413 can counter the beneficial effects of that attention at later areas (Figure 6: Figure
414 Supplement 1).

415 Our work raises the question: is attention applied simply according to tuning or
416 is it targeted to best optimize function on a given task? We suggested a simple ex-
417 periment (Figure 7) that would reveal whether non-tuning factors, such as gradients,
418 guide attentional modulation. In [15] the correlation coefficient between an index of
419 tuning and an index of attentional modulation was .52 for a population of V4 neurons,
420 suggesting factors other than selectivity influence attention. Furthermore, many at-
421 tention studies, including that one, use only preferred and anti-preferred stimuli and
422 therefore don't include a thorough investigation of the relationship between tuning and
423 attentional modulation. [55] uses multiple stimuli to provide support for the FSGM,
424 however the interpretation is limited by the fact that they only report population av-
425 erages. [78] investigated the relationship between tuning strength and the strength
426 of attentional modulation on a cell-by-cell basis. While they did find a correlation
427 (particularly for binocular disparity tuning), it was relatively weak, which leaves room
428 for the possibility that tuning is not the primary factor that determines attentional
429 modulation. Local connectivity is also likely to play a role, as a correlation between
430 normalization and attentional modulation has been shown [68].

431 A major challenge for understanding the biological implementation of selective
432 attention is determining how such a precise attentional signal is carried by feedback
433 connections. We believe that it is plausible that the visual system can learn the
434 connections needed to carry out gradient-based attention. For example, if a high-level
435 neuron related to the classification of an image sends a feedback connection to lower
436 areas, an anti-Hebbian post-pre spike timing-dependent learning rule would strengthen
437 the connection from the high level neuron to the low level one, if the lower level one
438 causes the firing of the higher. In this way, neurons in later areas can learn to target
439 the cells in earlier areas that caused them to fire. In contrast, it is actually more
440 difficult to imagine how higher areas could learn the connections needed to target
441 neurons according to their tuning, as in the FSGM. The machine learning literature
442 on attention and learning may inspire other useful hypotheses on underlying brain
443 mechanisms [100, 47].

444 The concept of attention has been introduced in these models previously in the
445 machine learning literature [60]. Generally, this kind of attention relates to what
446 would be called overt spatial attention in the neuroscience literature. That is, the
447 attention mechanism serially selects areas of the input image for further processing,
448 rather than modulating the activity of neurons representing those areas (as in our
449 model of spatial attention). Other work has been done using attention to selectively
450 process image features [89] and it would be interesting to compare the workings of
451 that model to the feature-based attention used in our study.

452 While CNNs have representations that are similar to the ventral stream, they lack
453 many biological details including recurrent connections, dynamics, cell types, and noisy
454 responses. Preliminary work has shown that these elements can be incorporated into
455 a CNN structure, and attention can enhance performance in this more biologically-
456 realistic architecture [49]. Furthermore, while the current work does not include neural
457 noise independent of the stimulus, the fact that a given image is presented in many
458 contexts (different merged images or different array images) can be thought of as a

459 form of highly structured noise that does produce variable responses to the same image.

460 Another biological detail that this model lacks is "skip connections," where one
461 layer feeds into both the layer directly after it and deeper layers after that [33, 35]
462 as in connections from V2 to V4 or V4 to parietal areas [96]. Our results regarding
463 propagation of changes through the network suggest that synaptic distance from the
464 classifier is a relevant feature—one that is less straight forward to determine in a
465 network with skip connections.

466 Because experimenters can easily control the image, defining a cells function in
467 terms of how it responds to stimuli makes practical sense. However, it may be that
468 thinking about visual areas in terms of their synaptic distance from decision-making
469 areas such as prefrontal cortex [34] can be more useful for the study of attention
470 than thinking in terms of their distance from the retina. Thus far, coarse stimulation
471 protocols have found a relationship between the tuning of neural populations and their
472 impact on perception [61, 21, 81]. However, studies of the relationship between tuning
473 and choice probabilities suggest that a neurons preferred stimulus is not always an
474 indication of its causal role in classification ([102, 73], though see [39]). Targeted
475 stimulation protocols and a more fine-grained ability to determine both upstream
476 drivers of, and downstream responses driven by, stimulated neurons will be needed to
477 better address these issues.

478 4. Methods

479 4.1. Key Resources

480 The weights for the model ("VGG-16") came from [26] (RRID SCR_016494).

481 4.2. Network Model

482 This work uses a deep convolutional neural network (CNN) as a model of the
483 ventral visual stream. Convolutional neural networks are feedforward artificial neural
484 networks that consist of a few basic operations repeated in sequence, key among them
485 being the convolution. The specific CNN architecture used in the study comes from
486 [85] (VGG-16D) and is shown in Figure 1A (a previous variant of this work used
487 a smaller network [48]). For this study, all the layers of the CNN except the final
488 classifier layer were pre-trained using backpropagation on the ImageNet classification
489 task, which involves doing 1000-way object categorization (weights provided by [26]).
490 The training of the top layer is described in subsequent sections. Here we describe the
491 basic workings of the CNN model we use, with details available in [85].

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

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

495 where \star indicates convolution, and $[x]_+ = x$ if $x > 0$, otherwise $x = 0$. W^{lk} is the
496 k^{th} convolutional filter at the l^{th} layer. The application of each filter results in a 2-D
497 feature map (the number of filters used varies across layers and is given in parenthesis
498 in Figure 1A). x_{ij}^{lk} is the activity of the unit at the i, j^{th} spatial location in the k^{th}
499 feature map at the l^{th} layer. X^{l-1} is thus the activity of all units at the layer below
500 the l^{th} layer. The input to the network is a 224 by 224 pixel RGB image, and thus the
501 first convolution is applied to these pixel values. Convolutional filters are 3x3. For the

502 purposes of this study the convolutional layers are most relevant, and will be referred
503 to according to their numbering in Figure 1A (numbers in parentheses indicate number
504 of feature maps per layer).

505 Max pooling layers reduce the size of the feature maps by taking the maximum
506 activity value of units in a given feature map in non-overlapping 2x2 windows. Through
507 this, the size of the feature maps decreases after each max pooling (layers 1 and 2: 224
508 x 224; 3 and 4: 112 x 112; 5, 6, and 7: 56 x 56. 8, 9, and 10: 28 x 28; 11, 12, and 13:
509 14 x 14).

510 The final two layers before the classifier are each fully-connected to the layer below
511 them, with the number of units per layer given in parenthesis in Figure 1A. Therefore,
512 connections exist from all units from all feature maps in the last convolutional layer
513 (layer 13) to all 4096 units of the next layer, and so on. The top readout layer of
514 the network in [85] contained 1000 units upon which a softmax classifier was used to
515 output a ranked list of category labels for a given image. Looking at the top-5 error
516 rate (wherein an image is correctly labeled if the true category appears in the top five
517 categories given by the network), this network achieved 92.7% accuracy. With the
518 exception of the gradient calculations described below, we did not use this 1000-way
519 classifier, but rather replaced it with a series of binary classifiers.

520 4.3. Object Category Attention Tasks

521 The tasks we use to probe the effects of feature-based attention in this network
522 involve determining if a given object category is present in an image or not, similar to
523 tasks used in [88, 72, 43]. To have the network perform this specific task, we replaced
524 the final layer in the network with a series of binary classifiers, one for each category
525 tested (Figure 1B). We tested a total of 20 categories: paintbrush, wall clock, seashore,
526 paddlewheel, padlock, garden spider, long-horned beetle, cabbage butterfly, toaster,
527 greenhouse, bakery, stone wall, artichoke, modem, football helmet, stage, mortar,
528 consomme, dough, bathtub. Binary classifiers were trained using ImageNet images
529 taken from the 2014 validation set (and were therefore not used in the training of
530 the original model). A total of 35 unique true positive images were used for training
531 for each category, and each training batch was balanced with 35 true negative images
532 taken from the remaining 19 categories. The results shown here come from using
533 logistic regression as the binary classifier, though trends in performance are similar if
534 support vector machines are used.

535 Once these binary classifiers are trained, they are then used to classify more chal-
536 lenging test images. Experimental results suggest that classifiers trained on unattended
537 and isolated object images are appropriate for reading out attended objects in clut-
538 tered images [104]. These test images are composed of multiple individual images
539 (drawn from the 20 categories) and are of two types: "merged" and "array". Merged
540 images are generated by transparently overlaying two images, each from a different
541 category (specifically, pixel values from each are divided by two and then summed).
542 Array images are composed of four separate images (all from different categories) that
543 are scaled down to 112 by 112 pixels and placed on a two by two grid. The images that
544 comprise these test images also come from the 2014 validation set, but are separate
545 from those used to train the binary classifiers. See examples of each in Figure 1C. Test
546 image sets are balanced (50% do contain the given category and 50% do not, 150 total
547 test images per category). Both true positive and true negative rates are recorded and
548 overall performance is the average of these rates.

549 *4.4. Object Category Gradient Calculations*

550 When neural networks are trained via backpropagation, gradients are calculated
 551 that indicate how a given weight in the network impacts the final classification. We
 552 use this same method to determine how a given unit’s activity impacts the final clas-
 553 sification. Specifically, we input a ”merged” image (wherein one of the images belongs
 554 to the category of interest) to the network. We then use gradient calculations to deter-
 555 mine the changes in activity that would move the 1000-way classifier toward classifying
 556 that image as belonging to the category of interest (i.e. rank that category highest).
 557 We average these activity changes over images and over all units in a feature map.
 558 This gives a single value per feature map:

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

559 where H and W are the spatial dimensions of layer l and N_c is the total number of
 560 images from the category (here $N_c = 35$, and the merged images used were generated
 561 from the same images used to generate tuning curves, described below). $E(n)$ is
 562 the error of the 1000-way classifier in response to image n , which is defined as the
 563 difference between the activity vector of the final layer (after the soft-max operation)
 564 and a one-hot vector, wherein the correct label is the only non-zero entry. Because
 565 we are interested in activity changes that would decrease the error value, we negate
 566 this term. The gradient value we end up with thus indicates how the feature map’s
 567 activity would need to change to make the network more likely to classify an image as
 568 the desired category. Repeating this procedure for each category, we obtain a set of
 569 gradient values (one for each category, akin to a tuning curve), for each feature map:
 570 \mathbf{g}^{lk} . Note that, as these values result from applying the chain rule through layers of
 571 the network, they can be very small, especially for the earliest layers. For this study,
 572 the sign and relative magnitudes are of more interest than the absolute values.

573 *4.5. Oriented Grating Attention Tasks*

574 In addition to attending to object categories, we also test attention on simpler
 575 stimuli. In the orientation detection task, the network detects the presence of a given
 576 orientation in an image. Again, the final layer of the network is replaced by a series
 577 of binary classifiers, one for each of 9 orientations (0, 20, 40, 60, 80, 100, 120, 140,
 578 and 160 degrees. Gratings had a frequency of .025 cycles/pixel). The training sets
 579 for each were balanced (50% had only the given orientation and 50% had one of 8
 580 other orientations) and composed of full field (224 by 224 pixel) oriented gratings in
 581 red, blue, green, orange, or purple (to increase the diversity of the training images,
 582 they were randomly degraded by setting blocks of pixels ranging uniformly from 0%
 583 to 70% of the image to 0 at random). Test images were each composed of two oriented
 584 gratings of different orientation and color (same options as training images). Each
 585 of these gratings were of size 112 by 112 pixels and placed randomly in a quadrant
 586 while the remaining two quadrants were black (Figure 5A). Again, the test sets were
 587 balanced and performance was measured as the average of the true positive and true
 588 negative rates (100 test images per orientation).

589 These same test images were used for a task wherein the network had to classify the
 590 color of the grating that had the attended orientation (cross-featural task paradigms
 591 like this are commonly used in attention studies, such as [80]). For this, the final layer

592 of the network was replaced with a 5-way softmax color classifier. This color classifier
 593 was trained using the same full field oriented gratings used to train the binary classifiers
 594 (therefore, the network saw each color at all orientation values).

595 For another analysis, a joint feature and spatial attention task was used. This
 596 task is almost identical to the setup of the orientation detection task, except that the
 597 searched-for orientation would only appear in one of the four quadrants. Therefore,
 598 performance could be measured when applying feature-based attention to the searched-
 599 for orientation, spatial attention to the quadrant in which it could appear, or both.

600 4.6. How Attention is Applied

601 This study aims to test variations of the feature similarity gain model of attention,
 602 wherein neural activity is modulated by attention according to how much the neuron
 603 prefers the attended stimulus. To replicate this in our model, we therefore must first
 604 determine the extent to which units in the network prefer different stimuli ("tuning
 605 values"). When attention is applied to a given category, for example, units' activities
 606 are modulated according to these values.

607 4.6.1. Tuning Values

608 To determine tuning to the 20 object categories used, we presented the network with
 609 images of each object category (the same images on which the binary classifiers were
 610 trained) and measured the relative activity levels. Because feature-based attention is
 611 a spatially global phenomenon [103, 79], we treat all units in a feature map identically,
 612 and calculate tuning by averaging over them.

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

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

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

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

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

629 We also define the category with the highest tuning value as that feature map's
 630 most preferred, and the category with the lowest (most negative) value as the least or
 631 anti-preferred.

632 We apply the same procedure to generate tuning curves for orientation by using
633 the full field gratings used to train the orientation detection classifiers. The orientation
634 tuning values were used when applying attention in these tasks.

635 When measuring how correlated tuning values are with gradient values, shuffled
636 comparisons are used. To do this shuffling, correlation coefficients are calculated from
637 pairing each feature map’s tuning values with a random other feature map’s gradient
638 values.

639 *4.6.2. Gradient Values*

640 In addition to applying attention according to tuning, we also attempt to generate
641 the ”best possible” attentional modulation by utilizing gradient values. These gradient
642 values are calculated slightly differently from those described above (4.4), because they
643 are meant to represent how feature map activity should change in order to increase
644 binary classification performance, rather than just increase the chance of classifying
645 an image as a certain object.

646 The error functions used to calculate gradient values for the category and orienta-
647 tion detection tasks were for the binary classifiers associated with each object/orientation.
648 A balanced set of test images was used. Therefore a feature map’s gradient value for
649 a given object/orientation is the averaged activity change that would increase binary
650 classification performance for that object/orientation. Note that on images that the
651 network already classifies correctly, gradients are zero. Therefore, the gradient values
652 are driven by the errors: false negatives (classifying an image as not containing the
653 category when it does) and false positives (classifying an image as containing the cate-
654 gory when it does not). In our detection tasks, the former error is more prevalent than
655 the latter, and thus is the dominant impact on the gradient values. Because of this,
656 gradient values calculated this way end up very similar to those described in Methods
657 4.4, as they are driven by a push to positively classify the input as the given category.

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

665 In the orientation detection task, the test images used for gradient calculations (50
666 images per orientation) differed from those used to assess performance. For the object
667 detection task, images used for gradient calculations (45 per category; preliminary
668 tests for some categories using 90 images gave similar results) were drawn from the
669 same pool as, but different from, those used to test detection performance. Gradient
670 values were calculated separately for merged and array images.

671 *4.6.3. Spatial Attention*

672 In the feature similarity gain model of attention, attention is applied according to
673 how much a cell prefers the attended feature, and location is considered a feature like
674 any other. In CNNs, each feature map results from applying the same filter at different
675 spatial locations. Therefore, the 2-D position of a unit in a feature map represents
676 more or less the spatial location to which that unit responds. Via the max-pooling
677 layers, the size of each feature map shrinks deeper in the network, and each unit

678 responds to a larger area of image space, but the "retinotopy" is still preserved. Thus,
 679 when we apply spatial attention to a given area of the image, we enhance the activity
 680 of units in that area of the feature maps and decrease the activity of units in other
 681 areas. In this study, spatial attention is applied to a given quadrant of the image.

682 4.6.4. Implementation Options

683 The values discussed above determine how strongly different feature maps or units
 684 should be modulated under different attentional conditions. We will now lay out the
 685 different implementation options for that modulation. The multiplicative bidirectional
 686 form of attention is used throughout this paper (with the exception of Figure 3 where
 687 it is compared to the others). Other implementations are only used for the Supple-
 688 mentary Results.

689 First, the modulation can be multiplicative or additive. That is, when attending
 690 to category c , the slope of the rectified linear units can be multiplied the tuning value
 691 for category c weighted by the strength parameter, β :

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

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

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

694 Strength of attention is varied via the strength parameter, β . For the additive effect,
 695 manipulations are multiplied by μ_l , the average activity level across all units of layer
 696 l in response to all images (for each of the 13 layers respectively: 20, 100, 150, 150,
 697 240, 240, 150, 150, 80, 20, 20, 10, 1). When gradient values are used in place of tuning
 698 values, we normalize them by the maximum value at a layer, to be the same order of
 699 magnitude as the tuning values: $\mathbf{g}^l / \max(|\mathbf{g}^l|)$.

700 Recall that for feature-based attention all units in a feature map are modulated the
 701 same way, as feature-based attention has been found to be spatially global. In the case
 702 of spatial attention, however, tuning values are not used and a unit's modulation is
 703 dependent on its location in the feature map. Specifically, the tuning value term is set
 704 to +1 if the i, j position of the unit is in the attended quadrant and to -1 otherwise.
 705 For feature-based attention tasks, β ranged from 0 to a maximum of 11.85 (object
 706 attention) and 0 to 4.8 (orientation attention). For spatial attention tasks, it ranged
 707 from 0 to 1.

708 Next, we chose whether attention only enhances units that prefer the attended
 709 feature, or also decreases activity of those that don't prefer it. For the latter, the
 710 tuning values are used as-is. For the former, the tuning values are positively-rectified:
 711 $[\mathbf{f}^{lk}]_+$.

712 Combining these two factors, there are four implementation options: additive
 713 positive-only, multiplicative positive-only, additive bidirectional, and multiplicative
 714 bidirectional.

715 The final option is the layer in the network at which attention is applied. We try
 716 attention at all convolutional layers individually and, in figure 3, simultaneously (when
 717 applying simultaneously the strength range tested is a tenth of that when applying to
 718 a single layer).

719 *4.7. Signal Detection Calculations*

720 For the joint spatial-feature attention task (Figure 5), we calculated criteria (c ,
721 "threshold") and sensitivity (d') using true (TP) and false (FP) positive rates as follows
722 [52] :

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

723 where Φ^{-1} is the inverse cumulative normal distribution function. c is a measure of
724 the distance from a neutral threshold situated between the mean of the true negative
725 and true positive distributions. Thus, a positive c indicates a stricter threshold (fewer
726 inputs classified as positive) and a negative c indicates a more lenient threshold (more
727 inputs classified as positive). The sensitivity was calculated as:

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

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

730 To prevent the individual terms in these expressions from going to $\pm\infty$, false
731 positive rates of $< .01$ were set to $.01$ and true positive rates of $> .99$ were set to $.99$.

732 *4.8. Assessment of Feature Similarity Gain Model and Feature Matching Behavior*

733 In Figure 6, we examined the effects that applying attention at certain layers in the
734 network (specifically 2, 6, 8, 10, and 12) has on activity of units at other layers. Atten-
735 tion was applied with $\beta = .5$. The recording setup is designed to mimic the analysis
736 of [55]. Here, the images presented to the network are full-field oriented gratings of
737 all orientation-color combinations. Feature map activity is measured as the spatially
738 averaged activity of all units in a feature map in response to an image. Activity in
739 response to a given orientation is further averaged over all colors. We calculate the
740 ratio of activity when attention is applied to a given orientation (and the orientation
741 is present in the image) over activity in response to the same image when no attention
742 is applied. These ratios are then organized according to orientation preference: the
743 most preferred is at location 0, then the average of next two most preferred at location
744 1, and so on with the average of the two least preferred orientations at location 4 (the
745 reason for averaging of pairs is to match [55] as closely as possible). Fitting a line to
746 these points gives a slope and intercept for each feature map (lines are fit using the
747 least squares method). FSGM predicts a negative slope and an intercept greater than
748 one.

749 To test for signs of feature matching behavior, each feature map's preferred (most
750 positive tuning value) and anti-preferred (most negative tuning value) orientations
751 are determined. Activity is recorded when attention is applied to the preferred or
752 anti-preferred orientation and activity ratios are calculated. According to the FSGM,
753 activity when the preferred orientation is attended should be greater than when the
754 anti-preferred is attended, regardless of whether the image is of the preferred or anti-
755 preferred orientation. According to the feature matching (FM) model, however, ac-
756 tivity when attending the presented orientation should be greater than activity when
757 attending an absent orientation, regardless of whether the orientation is preferred or
758 not. Therefore, we say that a feature map is displaying feature matching behavior
759 if (1) activity is greater when attending the preferred orientation when the preferred
760 is present versus when the anti-preferred is present, and (2) activity is greater when
761 attending the anti-preferred orientation when the anti-preferred is present versus when

762 the preferred is present. The second criteria distinguishes feature matching behavior
763 from FSGM.

764 *4.9. Correlating Activity Changes with Performance*

765 In Figure 6: Figure Supplement 2, we use two different measures of attention-
766 induced activity changes in order to probe the relationship between activity and clas-
767 sification performance. In both cases, the network is performing the orientation de-
768 tection task described in Figure 5A and performance is measured only in terms of
769 true positive rates. Because we know attention to increase both true and false posi-
770 tive rates, we would expect a positive correlation between activity changes and true
771 positive performance, but a negative correlation between activity changes and true
772 negative rates. This predicts that activity changes will have a monotonic relation-
773 ship with true positive performance, but an inverted U-shaped relationship with total
774 performance. Since we are calculating correlation coefficients of activity with perfor-
775 mance, which measure a linear relationship, we use the rate of true positives as our
776 measure of performance.

777 The first measure is meant to capture feature similarity gain model-like behavior
778 in a way similar to the metric described in Figure 6B. The main difference is that
779 that measure is calculated over a population of images of different stimuli, whereas
780 the variant introduced here can be calculated on an image-by-image basis. Images
781 that contain a given orientation are shown to the network and the spatially-averaged
782 activity of feature maps is recorded when attention is applied to that orientation
783 and when it is not. The ratio of these activities is then plotted against each feature
784 map’s tuning value for the orientation. According to the FSGM, this ratio should be
785 above 1 for feature maps with positive tuning values and less than one for those with
786 negative tuning values. Therefore, we use the slope of the line fitted to these ratios
787 plotted as a function of tuning values as an indication of the extent to which activity
788 is FSGM-like (with positive slopes more FSGM-like). The median slope over a set of
789 images of a given orientation is paired with the change in performance on those images
790 with attention. This gives one pair for each combination of orientation, strength ($\beta =$
791 $.15, .30, .45, .60, .75, .90$), and layer at which attention was applied (activity changes are
792 only recorded if attention was applied at or before the recorded layer). The correlation
793 coefficient between these value pairs is plotted as the orange line in Figure 6: Figure
794 Supplement 2C.

795 The second measure aims to characterize activity in terms of its downstream effects,
796 rather than the contents of the input (“Vector Angle” measure, see Figure 6: Figure
797 Supplement 2A for a visualization). It is therefore more aligned with the gradient-
798 based approach to attention rather than tuning, and is thus related to “choice proba-
799 bility” measures [102, 73]. First, for a particular orientation, images that both do and
800 do not contain that orientation are shown to the network. Activity (spatially-averaged
801 over each feature map) in response to images classified as containing the orientation
802 (i.e., both true and false positives) is averaged in order to construct a vector in activity
803 space that represents positive classification for a given layer. To reduce complications
804 of working with vectors in high dimensions, principal components are found that cap-
805 ture at least 90% of the variance of the activity in response to all images, and all
806 computations are done in this lower dimensional space. The next step is to determine
807 if attention moves activity in a given layer closer to this direction of positive classi-
808 fication. For this, only images that contain the given orientation are used. For each

809 image, the cosine of the angle between the positive-classification vector and the activ-
810 ity in response to the image is calculated. The median of these angles over a set of
811 images is calculated separately for when attention is applied and when it is not. The
812 difference between these medians (with-attention minus without-attention) is paired
813 with the change in performance that comes with attention on those images. Then the
814 same correlation calculation is done with these pairs as described above.

815 The outcome of these analyses is a correlation coefficient between the measure of
816 activity changes and performance changes. This gives two values per layer: one for
817 the FSGM-like measure and one for the vector angle measure. To determine if these
818 two values are significantly different, we performed a bootstrap analysis. For this,
819 correlation coefficients were recalculated using simulated data made by sampling with
820 replacement from the true data. We do this 100 times and perform a two-sided t-test
821 to test for differences between the two measures.

822 *4.10. Experimental Data*

823 Model results were compared to previously published data coming from several
824 studies. In [54], a category detection task was performed using stereogram stimuli
825 (on object present trials, the object image was presented to one eye and a noise mask
826 to another). The presentation of the visual stimuli was preceded by a verbal cue
827 that indicated the object category that would later be queried (cued trials) or by
828 meaningless noise (uncued trials). After visual stimulus presentation, subjects were
829 asked if an object was present and, if so, if the object was from the cued category
830 (categories were randomized for uncued trials). In Experiment 1 ('Cat-Drawings' in
831 Figure 4B), the object images were line drawings (one per category) and the stimuli
832 were presented for 1.5 sec. In Experiment 2 ('Cat-Images'), the object images were
833 grayscale photographs (multiple per category) and presented for 6 sec (of note: this
834 presumably allows for several rounds of feedback processing, in contrast to our purely
835 feedforward model). True positives were counted as trials wherein a given object
836 category was present and the subject correctly indicated its presence when queried.
837 False positives were trials wherein no category was present and subjects indicated that
838 the queried category was present.

839 In [53], a similar detection task was used. Here, subjects detected the presence of
840 an uppercase letter that (on target present trials) was presented rapidly and followed
841 by a mask. Prior to the visual stimulus, a visual ('Letter-Vis') or audio ('Letter-Aud')
842 cue indicated a target letter. After the visual stimulus, the subjects were required to
843 indicate whether any letter was present. True positives were trials in which a letter was
844 present and the subject indicated it (only uncued trials or validly cued trials—where
845 the cued letter was the letter shown—were considered here). False positives were trials
846 where no letter was present and the subject indicated that one was.

847 The task in [43] was also an object category detection task ('Objects'). Here, an
848 array of several images was flashed on the screen with one image marked as the target.
849 All images were color photographs of objects in natural scenes. In certain blocks,
850 the subjects knew in advance which category they would later be queried about (cued
851 trials). On other trials, the queried category was only revealed after the visual stimulus
852 (uncued). True positives were trials in which the subject indicated the presence of the
853 queried category when it did exist in the target image. False positives were trials in
854 which the subject indicated the presence of the cued category when it was not in the
855 target image. Data from trials using basic category levels with masks were used for

856 this study.

857 Finally, we include one study using macaques ('Ori-Change') wherein both neural
858 and performance changes were measured [57]. In this task, subjects had to report a
859 change in orientation that could occur in one of two stimuli. On cued trials, the change
860 occurred in the cued stimulus in 80% of trials and the uncued stimulus in 20% of tri-
861 als. On neutrally-cued trials, subjects were not given prior information about where
862 the change was likely to occur (50% at each stimulus). Therefore performance could
863 be compared under conditions of low (uncued stimuli), medium (neutrally cued stim-
864 ulti), and high (cued stimuli) attention strength. Correct detection of an orientation
865 change in a given stimulus (indicated by a saccade) is considered a true positive and a
866 saccade to the stimulus prior to any orientation change is considered a false positive.
867 True negatives are defined as correct detection of a change in the uncued stimulus
868 (as this means the subject correctly did not perceive a change in the stimulus under
869 consideration) and false negatives correspond to a lack of response to an orientation
870 change. While this task includes a spatial attention component, it is still useful as a
871 test of feature-based attention effects. Previous work has demonstrated that, during a
872 change detection task, feature-based attention is deployed to the pre-change features
873 of a stimulus [18, 58]. Therefore, because the pre-change stimuli are of differing orien-
874 tations, the cueing paradigm used here controls the strength of attention to orientation
875 as well.

876 In cases where the true and false positive rates were not published, they were ob-
877 tained via personal communications with the authors. Not all changes in performance
878 were statistically significant, but we plot them to show general trends.

879 We calculate the activity changes required in the model to achieve the behavioral
880 changes observed experimentally by using the data plotted in Figure 4B. We determine
881 the average β value for the neutral and cued conditions by finding the β value of the
882 point on the model line nearest to the given data point. Specifically, we average the β
883 values found for the four datasets whose experiments are most similar to our merged
884 image task (Cat-Drawings, Cat-Images, Letter-Aud, and Letter-Vis).

885 5. Acknowledgements

886 We are very grateful to the authors who so readily shared details of their behavioral
887 data upon request: J. Patrick Mayo, Gary Lupyan, and Mika Koivisto. We further
888 thank J. Patrick Mayo for helpful comments on the manuscript.

889 6. References

- 890 [1] Mohamed Abdelhack and Yukiyasu Kamitani. Sharpening of hierarchical visual
891 feature representations of blurred images. *eNeuro*, pages ENEURO-0443, 2018.
- 892 [2] Aharon Azulay and Yair Weiss. Why do deep convolutional networks generalize
893 so poorly to small image transformations? *arXiv preprint arXiv:1805.12177*,
894 2018.
- 895 [3] Nicholas Baker, Hongjing Lu Lu, Gennady Erlikhman, and Philip Kellman. Deep
896 convolutional networks do not make classifications based on global object shape.
897 *Vision Sciences Society Annual Meeting*, 2018.

- 898 [4] Ji Won Bang and Dobromir Rahnev. Stimulus expectation alters decision crite-
899 rion but not sensory signal in perceptual decision making. *Scientific reports*, 7
900 (1):17072, 2017.
- 901 [5] Jalal K Baruni, Brian Lau, and C Daniel Salzman. Reward expectation differ-
902 entially modulates attentional behavior and activity in visual area v4. *Nature*
903 *neuroscience*, 18(11):1656, 2015.
- 904 [6] Narcisse P Bichot, Matthew T Heard, Ellen M DeGennaro, and Robert Desi-
905 monne. A source for feature-based attention in the prefrontal cortex. *Neuron*, 88
906 (4):832–844, 2015.
- 907 [7] Ali Borji and Laurent Itti. Optimal attentional modulation of a neural popula-
908 tion. *Frontiers in computational neuroscience*, 8, 2014.
- 909 [8] Geoffrey M Boynton. A framework for describing the effects of attention on
910 visual responses. *Vision research*, 49(10):1129–1143, 2009.
- 911 [9] David A Bridwell and Ramesh Srinivasan. Distinct attention networks for feature
912 enhancement and suppression in vision. *Psychological science*, 23(10):1151–1158,
913 2012.
- 914 [10] Elizabeth A Buffalo, Pascal Fries, Rogier Landman, Hualou Liang, and Robert
915 Desimone. A backward progression of attentional effects in the ventral stream.
916 *Proceedings of the National Academy of Sciences*, 107(1):361–365, 2010.
- 917 [11] Claus Bundesen. A theory of visual attention. *Psychological review*, 97(4):523,
918 1990.
- 919 [12] Santiago A Cadena, George H Denfield, Edgar Y Walker, Leon A Gatys, An-
920 dreas S Tolias, Matthias Bethge, and Alexander S Ecker. Deep convolutional
921 models improve predictions of macaque v1 responses to natural images. *bioRxiv*,
922 page 201764, 2017.
- 923 [13] Marisa Carrasco. Visual attention: The past 25 years. *Vision research*, 51(13):
924 1484–1525, 2011.
- 925 [14] Kyle R Cave. The featuregate model of visual selection. *Psychological research*,
926 62(2):182–194, 1999.
- 927 [15] Leonardo Chelazzi, John Duncan, Earl K Miller, and Robert Desimone. Re-
928 sponses of neurons in inferior temporal cortex during memory-guided visual
929 search. *Journal of neurophysiology*, 80(6):2918–2940, 1998.
- 930 [16] Sharat Chikkerur, Thomas Serre, Cheston Tan, and Tomaso Poggio. What and
931 where: A bayesian inference theory of attention. *Vision research*, 50(22):2233–
932 2247, 2010.
- 933 [17] Marlene R Cohen and John HR Maunsell. Attention improves performance
934 primarily by reducing interneuronal correlations. *Nature neuroscience*, 12(12):
935 1594–1600, 2009.

- 936 [18] Marlene R Cohen and John HR Maunsell. Using neuronal populations to study
937 the mechanisms underlying spatial and feature attention. *Neuron*, 70(6):1192–
938 1204, 2011.
- 939 [19] Trinity B Crapse, Hakwan Lau, and Michele A Basso. A role for the superior
940 colliculus in decision criteria. *Neuron*, 97(1):181–194, 2018.
- 941 [20] Tolga Çukur, Shinji Nishimoto, Alexander G Huth, and Jack L Gallant. At-
942 tention during natural vision warps semantic representation across the human
943 brain. *Nature neuroscience*, 16(6):763–770, 2013.
- 944 [21] Gregory C DeAngelis, Bruce G Cumming, and William T Newsome. Cortical
945 area mt and the perception of stereoscopic depth. *Nature*, 394(6694):677, 1998.
- 946 [22] Cathryn J Downing. Expectancy and visual-spatial attention: effects on per-
947 ceptual quality. *Journal of Experimental Psychology: Human perception and*
948 *performance*, 14(2):188, 1988.
- 949 [23] Miguel P Eckstein, Matthew F Peterson, Binh T Pham, and Jason A Droll.
950 Statistical decision theory to relate neurons to behavior in the study of covert
951 visual attention. *Vision research*, 49(10):1097–1128, 2009.
- 952 [24] Michael Eickenberg, Alexandre Gramfort, Gaël Varoquaux, and Bertrand
953 Thirion. Seeing it all: Convolutional network layers map the function of the
954 human visual system. *NeuroImage*, 152:184–194, 2017.
- 955 [25] Pascal Fries, John H Reynolds, Alan E Rorie, and Robert Desimone. Modulation
956 of oscillatory neuronal synchronization by selective visual attention. *Science*, 291
957 (5508):1560–1563, 2001.
- 958 [26] Davi Frossard. *VGG in TensorFlow*. <https://www.cs.toronto.edu/~frossard/post/vgg16>
959 Accessed: 2017-03-01.
- 960 [27] Kunihiko Fukushima. Neocognitron: A hierarchical neural network capable of
961 visual pattern recognition. *Neural networks*, 1(2):119–130, 1988.
- 962 [28] Umut Güçlü and Marcel AJ van Gerven. Deep neural networks reveal a gradient
963 in the complexity of neural representations across the ventral stream. *Journal*
964 *of Neuroscience*, 35(27):10005–10014, 2015.
- 965 [29] FH Hamker. The role of feedback connections in task-driven visual search. In
966 *Connectionist models in cognitive neuroscience*, pages 252–261. Springer, 1999.
- 967 [30] Fred H Hamker and James Worcester. Object detection in natural scenes by
968 feedback. In *International Workshop on Biologically Motivated Computer Vision*,
969 pages 398–407. Springer, 2002.
- 970 [31] Harold L Hawkins, Steven A Hillyard, Steven J Luck, Mustapha Mouloua,
971 Cathryn J Downing, and Donald P Woodward. Visual attention modulates sig-
972 nal detectability. *Journal of Experimental Psychology: Human Perception and*
973 *Performance*, 16(4):802, 1990.

- 974 [32] Benjamin Y Hayden and Jack L Gallant. Combined effects of spatial and feature-
975 based attention on responses of v4 neurons. *Vision research*, 49(10):1182–1187,
976 2009.
- 977 [33] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learn-
978 ing for image recognition. In *Proceedings of the IEEE conference on computer*
979 *vision and pattern recognition*, pages 770–778, 2016.
- 980 [34] Hauke R Heekeren, Sean Marrett, Peter A Bandettini, and Leslie G Ungerleider.
981 A general mechanism for perceptual decision-making in the human brain. *Nature*,
982 431(7010):859–862, 2004.
- 983 [35] Gao Huang, Zhuang Liu, Kilian Q Weinberger, and Laurens van der Maaten.
984 Densely connected convolutional networks. In *Proceedings of the IEEE confer-*
985 *ence on computer vision and pattern recognition*, volume 1, page 3, 2017.
- 986 [36] Daniel Kaiser, Nikolaas N Oosterhof, and Marius V Peelen. The neural dynamics
987 of attentional selection in natural scenes. *Journal of neuroscience*, 36(41):10522–
988 10528, 2016.
- 989 [37] Kohitij Kar, Jonas Kubilius, Elias Issa, Kailyn Schmidt, and James DiCarlo.
990 Evidence that feedback is required for object identity inferences computed by
991 the ventral stream. COSYNE, 2017.
- 992 [38] Sabine Kastner and Mark A Pinsk. Visual attention as a multilevel selection
993 process. *Cognitive, Affective, & Behavioral Neuroscience*, 4(4):483–500, 2004.
- 994 [39] Leor N Katz, Jacob L Yates, Jonathan W Pillow, and Alexander C Huk. Dis-
995 sociated functional significance of decision-related activity in the primate dorsal
996 stream. *Nature*, 535(7611):285, 2016.
- 997 [40] Seyed-Mahdi Khaligh-Razavi and Nikolaus Kriegeskorte. Deep supervised, but
998 not unsupervised, models may explain it cortical representation. *PLoS compu-*
999 *tational biology*, 10(11):e1003915, 2014.
- 1000 [41] Seyed-Mahdi Khaligh-Razavi, Linda Henriksson, Kendrick Kay, and Nikolaus
1001 Kriegeskorte. Fixed versus mixed rsa: Explaining visual representations by fixed
1002 and mixed feature sets from shallow and deep computational models. *Journal*
1003 *of Mathematical Psychology*, 76:184–197, 2017.
- 1004 [42] Saeed Reza Kheradpisheh, Masoud Ghodrati, Mohammad Ganjtabesh, and Tim-
1005 othée Masquelier. Deep networks can resemble human feed-forward vision in
1006 invariant object recognition. *Scientific reports*, 6:32672, 2016.
- 1007 [43] Mika Koivisto and Ella Kahila. Top-down preparation modulates visual cate-
1008 gorization but not subjective awareness of objects presented in natural back-
1009 grounds. *Vision Research*, 133:73–80, 2017.
- 1010 [44] Simon Kornblith and Doris Y Tsao. How thoughts arise from sights: inferotem-
1011 poral and prefrontal contributions to vision. *Current Opinion in Neurobiology*,
1012 46:208–218, 2017.

- 1013 [45] Richard J Krauzlis, Lee P Lovejoy, and Alexandre Zénon. Superior colliculus
1014 and visual spatial attention. *Annual review of neuroscience*, 36:165–182, 2013.
- 1015 [46] Jonas Kubilius, Stefania Bracci, and Hans P Op de Beeck. Deep neural networks
1016 as a computational model for human shape sensitivity. *PLoS computational
1017 biology*, 12(4):e1004896, 2016.
- 1018 [47] Timothy P Lillicrap, Daniel Cownden, Douglas B Tweed, and Colin J Aker-
1019 man. Random synaptic feedback weights support error backpropagation for
1020 deep learning. *Nature communications*, 7, 2016.
- 1021 [48] Grace W Lindsay. Feature-based attention in convolutional neural networks.
1022 *arXiv preprint arXiv:1511.06408*, 2015.
- 1023 [49] Grace W Lindsay, Dan B Rubin, and Kenneth D Miller. The stabilized supralin-
1024 ear network replicates neural and performance correlates of attention. COSYNE,
1025 2017.
- 1026 [50] Bradley C Love, Olivia Guest, Piotr Slomka, Victor M Navarro, and Edward
1027 Wasserman. Deep networks as models of human and animal categorization. In
1028 *CogSci*, 2017.
- 1029 [51] Steven J Luck, Leonardo Chelazzi, Steven A Hillyard, and Robert Desimone.
1030 Neural mechanisms of spatial selective attention in areas v1, v2, and v4 of
1031 macaque visual cortex. *Journal of neurophysiology*, 77(1):24–42, 1997.
- 1032 [52] Thomas Zhihao Luo and John HR Maunsell. Neuronal modulations in visual
1033 cortex are associated with only one of multiple components of attention. *Neuron*,
1034 86(5):1182–1188, 2015.
- 1035 [53] Gary Lupyan and Michael J Spivey. Making the invisible visible: Verbal but not
1036 visual cues enhance visual detection. *PLoS One*, 5(7):e11452, 2010.
- 1037 [54] Gary Lupyan and Emily J Ward. Language can boost otherwise unseen objects
1038 into visual awareness. *Proceedings of the National Academy of Sciences*, 110(35):
1039 14196–14201, 2013.
- 1040 [55] Julio C Martinez-Trujillo and Stefan Treue. Feature-based attention increases
1041 the selectivity of population responses in primate visual cortex. *Current Biology*,
1042 14(9):744–751, 2004.
- 1043 [56] John HR Maunsell and Erik P Cook. The role of attention in visual processing.
1044 *Philosophical Transactions of the Royal Society of London B: Biological Sciences*,
1045 357(1424):1063–1072, 2002.
- 1046 [57] J Patrick Mayo and John HR Maunsell. Graded neuronal modulations related
1047 to visual spatial attention. *Journal of Neuroscience*, 36(19):5353–5361, 2016.
- 1048 [58] J Patrick Mayo, Marlene R Cohen, and John HR Maunsell. A refined neuronal
1049 population measure of visual attention. *PloS one*, 10(8):e0136570, 2015.

- 1050 [59] Carrie J McAdams and John HR Maunsell. Effects of attention on orientation-
1051 tuning functions of single neurons in macaque cortical area v4. *Journal of Neu-*
1052 *roscience*, 19(1):431–441, 1999.
- 1053 [60] Volodymyr Mnih, Nicolas Heess, Alex Graves, et al. Recurrent models of visual
1054 attention. In *Advances in neural information processing systems*, pages 2204–
1055 2212, 2014.
- 1056 [61] Sebastian Moeller, Trinity Crapse, Le Chang, and Doris Y Tsao. The effect of
1057 face patch microstimulation on perception of faces and objects. *Nature Neuro-*
1058 *science*, 20(5):743–752, 2017.
- 1059 [62] Ilya E Monosov, David L Sheinberg, and Kirk G Thompson. The effects of pre-
1060 frontal cortex inactivation on object responses of single neurons in the inferotem-
1061 poral cortex during visual search. *Journal of Neuroscience*, 31(44):15956–15961,
1062 2011.
- 1063 [63] Tirin Moore and Katherine M Armstrong. Selective gating of visual signals by
1064 microstimulation of frontal cortex. *Nature*, 421(6921):370, 2003.
- 1065 [64] Ari S Morcos, David GT Barrett, Neil C Rabinowitz, and Matthew Botvinick.
1066 On the importance of single directions for generalization. *arXiv preprint*
1067 *arXiv:1803.06959*, 2018.
- 1068 [65] Sancho I Moro, Michiel Tolboom, Paul S Khayat, and Pieter R Roelfsema. Neu-
1069 ronal activity in the visual cortex reveals the temporal order of cognitive opera-
1070 tions. *Journal of Neuroscience*, 30(48):16293–16303, 2010.
- 1071 [66] Brad C Motter. Neural correlates of feature selective memory and pop-out in
1072 extrastriate area v4. *Journal of Neuroscience*, 14(4):2190–2199, 1994.
- 1073 [67] Vidhya Navalpakkam and Laurent Itti. Search goal tunes visual features opti-
1074 mally. *Neuron*, 53(4):605–617, 2007.
- 1075 [68] Amy M Ni, Supratim Ray, and John HR Maunsell. Tuned normalization explains
1076 the size of attention modulations. *Neuron*, 73(4):803–813, 2012.
- 1077 [69] Marino Pagan, Luke S Urban, Margot P Wohl, and Nicole C Rust. Signals
1078 in inferotemporal and perirhinal cortex suggest an untangling of visual target
1079 information. *Nature neuroscience*, 16(8):1132–1139, 2013.
- 1080 [70] William K Page and Charles J Duffy. Cortical neuronal responses to optic flow
1081 are shaped by visual strategies for steering. *Cerebral cortex*, 18(4):727–739, 2007.
- 1082 [71] Marius V Peelen and Sabine Kastner. A neural basis for real-world visual search
1083 in human occipitotemporal cortex. *Proceedings of the National Academy of Sci-*
1084 *ences*, 108(29):12125–12130, 2011.
- 1085 [72] Marius V Peelen, Li Fei-Fei, and Sabine Kastner. Neural mechanisms of rapid
1086 natural scene categorization in human visual cortex. *Nature*, 460(7251):94, 2009.
- 1087 [73] Gopathy Purushothaman and David C Bradley. Neural population code for fine
1088 perceptual decisions in area mt. *Nature neuroscience*, 8(1):99, 2005.

- 1089 [74] Dobromir Rahnev, Hakwan Lau, and Floris P de Lange. Prior expectation
1090 modulates the interaction between sensory and prefrontal regions in the human
1091 brain. *Journal of Neuroscience*, 31(29):10741–10748, 2011.
- 1092 [75] Waseem Rawat and Zenghui Wang. Deep convolutional neural networks for
1093 image classification: A comprehensive review. *Neural Computation*, 2017.
- 1094 [76] Maximilian Riesenhuber and Tomaso Poggio. Hierarchical models of object
1095 recognition in cortex. *Nature neuroscience*, 2(11), 1999.
- 1096 [77] Edmund T Rolls and Gustavo Deco. Attention in natural scenes: neurophysio-
1097 logical and computational bases. *Neural networks*, 19(9):1383–1394, 2006.
- 1098 [78] Douglas A Ruff and Richard T Born. Feature attention for binocular disparity
1099 in primate area mt depends on tuning strength. *Journal of neurophysiology*, 113
1100 (5):1545–1555, 2015.
- 1101 [79] Melissa Saenz, Giedrius T Buracas, and Geoffrey M Boynton. Global effects of
1102 feature-based attention in human visual cortex. *Nature neuroscience*, 5(7):631,
1103 2002.
- 1104 [80] Melissa Saenz, Giedrius T Buraças, and Geoffrey M Boynton. Global feature-
1105 based attention for motion and color. *Vision research*, 43(6):629–637, 2003.
- 1106 [81] C Daniel Salzman, Kenneth H Britten, and William T Newsome. Cortical mi-
1107 crostimulation influences perceptual judgements of motion direction. *Nature*,
1108 346(6280):174–177, 1990.
- 1109 [82] K Seeliger, M Fritsche, U Güçlü, S Schoenmakers, J-M Schoffelen, SE Bosch, and
1110 MAJ van Gerven. Cnn-based encoding and decoding of visual object recognition
1111 in space and time. *bioRxiv*, page 118091, 2017.
- 1112 [83] John T Serences, Jens Schwarzbach, Susan M Courtney, Xavier Golay, and
1113 Steven Yantis. Control of object-based attention in human cortex. *Cerebral*
1114 *Cortex*, 14(12):1346–1357, 2004.
- 1115 [84] Thomas Serre, Lior Wolf, Stanley Bileschi, Maximilian Riesenhuber, and Tomaso
1116 Poggio. Robust object recognition with cortex-like mechanisms. *IEEE transac-*
1117 *tions on pattern analysis and machine intelligence*, 29(3):411–426, 2007.
- 1118 [85] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for
1119 large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- 1120 [86] Devarajan Sridharan, Nicholas A Steinmetz, Tirin Moore, and Eric I Knudsen.
1121 Does the superior colliculus control perceptual sensitivity or choice bias during
1122 attention? evidence from a multialternative decision framework. *Journal of*
1123 *Neuroscience*, 37(3):480–511, 2017.
- 1124 [87] Timo Stein and Marius V Peelen. Content-specific expectations enhance stim-
1125 ulus detectability by increasing perceptual sensitivity. *Journal of Experimental*
1126 *Psychology: General*, 144(6):1089, 2015.

- 1127 [88] Timo Stein and Marius V Peelen. Object detection in natural scenes: Independent
1128 effects of spatial and category-based attention. *Attention, Perception, &*
1129 *Psychophysics*, 79(3):738–752, 2017.
- 1130 [89] Marijn F Stollenga, Jonathan Masci, Faustino Gomez, and Jürgen Schmidhuber.
1131 Deep networks with internal selective attention through feedback connections.
1132 In *Advances in neural information processing systems*, pages 3545–3553, 2014.
- 1133 [90] Anne M Treisman and Garry Gelade. A feature-integration theory of attention.
1134 *Cognitive psychology*, 12(1):97–136, 1980.
- 1135 [91] Stefan Treue. Neural correlates of attention in primate visual cortex. *Trends in*
1136 *neurosciences*, 24(5):295–300, 2001.
- 1137 [92] Stefan Treue and Julio C Martinez Trujillo. Feature-based attention influences
1138 motion processing gain in macaque visual cortex. *Nature*, 399(6736):575, 1999.
- 1139 [93] Bryan P Tripp. Similarities and differences between stimulus tuning in the
1140 inferotemporal visual cortex and convolutional networks. In *Neural Networks*
1141 *(IJCNN), 2017 International Joint Conference on*, pages 3551–3560. IEEE, 2017.
- 1142 [94] John K Tsotsos, Scan M Culhane, Winky Yan Kei Wai, Yuzhong Lai, Neal Davis,
1143 and Fernando Nuflo. Modeling visual attention via selective tuning. *Artificial*
1144 *intelligence*, 78(1-2):507–545, 1995.
- 1145 [95] Shimon Ullman, Liav Assif, Ethan Fetaya, and Daniel Harari. Atoms of recog-
1146 nition in human and computer vision. *Proceedings of the National Academy of*
1147 *Sciences*, 113(10):2744–2749, 2016.
- 1148 [96] Leslie G Ungerleider, Thelma W Galkin, Robert Desimone, and Ricardo Gattass.
1149 Cortical connections of area v4 in the macaque. *Cerebral Cortex*, 18(3):477–499,
1150 2007.
- 1151 [97] Preeti Verghese. Visual search and attention: A signal detection theory ap-
1152 proach. *Neuron*, 31(4):523–535, 2001.
- 1153 [98] Louise Whiteley and Maneesh Sahani. Attention in a bayesian framework. *Fron-*
1154 *tiers in human neuroscience*, 6, 2012.
- 1155 [99] Jeremy M Wolfe. Guided search 2.0 a revised model of visual search. *Psycho-*
1156 *nomics bulletin & review*, 1(2):202–238, 1994.
- 1157 [100] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan
1158 Salakhudinov, Rich Zemel, and Yoshua Bengio. Show, attend and tell: Neural
1159 image caption generation with visual attention. In *International Conference on*
1160 *Machine Learning*, pages 2048–2057, 2015.
- 1161 [101] Daniel LK Yamins, Ha Hong, Charles F Cadieu, Ethan A Solomon, Darren Seib-
1162 ert, and James J DiCarlo. Performance-optimized hierarchical models predict
1163 neural responses in higher visual cortex. *Proceedings of the National Academy*
1164 *of Sciences*, 111(23):8619–8624, 2014.

- 1165 [102] Adam Zaidel, Gregory C DeAngelis, and Dora E Angelaki. Decoupled choice-
1166 driven and stimulus-related activity in parietal neurons may be misrepresented
1167 by choice probabilities. *Nature Communications*, 8, 2017.
- 1168 [103] Weiwei Zhang and Steven J Luck. Feature-based attention modulates feedfor-
1169 ward visual processing. *Nature neuroscience*, 12(1):24–25, 2009.
- 1170 [104] Ying Zhang, Ethan M Meyers, Narcisse P Bichot, Thomas Serre, Tomaso A Pog-
1171 gio, and Robert Desimone. Object decoding with attention in inferior temporal
1172 cortex. *Proceedings of the National Academy of Sciences*, 108(21):8850–8855,
1173 2011.
- 1174 [105] Huihui Zhou and Robert Desimone. Feature-based attention in the frontal eye
1175 field and area v4 during visual search. *Neuron*, 70(6):1205–1217, 2011.