

1 Models of the ventral stream that categorize and visualize 2 images

3 Elijah Christensen¹, Joel Zylberberg²

4
5 ¹Department of Physiology and Biophysics, University of Colorado Anschutz Medical Campus

6 ²Learning in Machines and Brains Program, CIFAR, Toronto, ON Canada

7 Abstract

8 A widely held view in visual systems neuroscience is that the ventral stream (VS) of mammalian
9 visual cortex maps visual inputs onto firing patterns that specify object identities. Recent
10 experiments, however, showed that information about object positions, sizes, etc. is encoded with
11 increasing explicitness along this pathway. Here, we show that computational models that identify
12 objects while *also* retaining enough information about the input image to allow its reconstruction,
13 can provide a better description of how primate VS represents this *category-orthogonal*
14 information, than do models that “just” identify objects. A thorough understanding of VS
15 computations might thus require considerations beyond object recognition.

16 Significance Statement

17 Our key finding is that ventral stream physiology is better described by a composite computational
18 objective of object recognition and reconstruction, rather than object recognition alone. Because this
19 finding potentially overturns the longstanding object recognition hypothesis of ventral stream function,
20 we expect it to have substantial impacts on visual systems neuroscience.

21 Introduction

22 The ventral stream (VS) of visual cortex begins in primary visual cortex (V1), ends in inferior temporal
23 cortex (IT), and is essential for object recognition. Accordingly, the long-standing belief in the field is
24 that the ventral stream could be understood as mapping visual scenes onto neuronal firing patterns that
25 represent object identity¹. Supporting that assertion, deep convolutional neural networks (DCNN's)
26 trained to categorize objects in natural images develop intermediate representations that resemble those
27 in primate VS²⁻⁵. However, several experimental findings appear at odds with the object recognition
28 hypothesis. VS and other visual areas are also engaged during visualization of both previously
29 encountered and novel scenes^{6,7}, suggesting that the VS can *generate* visual scenes in addition to
30 processing them as inputs. Furthermore, non-categorical information, about object positions⁸, sizes, etc.
31 is also represented with increasing explicitness in late VS areas V4 and IT⁹. This is not necessarily
32 expected in a “pure” object recognition system as the non-categorical information is not necessary for
33 the categorization task. Thus these recent findings challenge notion that ventral stream is purely an
34 object recognition system, and raise the question: What computational objective best explains VS
35 physiology¹⁰?

36 To address this question, we pursued a recently-popularized approach and trained deep neural networks
37 to perform different tasks: we then compared the trained neural networks' responses to image stimuli to
38 responses observed in neurophysiology experiments,^{3-5,9} to see which tasks yielded models that best
39 matched the neural data. We trained our networks to perform one of two visual tasks: a) recognize
40 objects; or b) recognize objects while *also* retaining enough information about the input image to allow

41 its reconstruction. We studied the evolution of categorical and non-categorical information
42 representations along the visual pathway within these models and compared that evolution with data
43 from monkey VS. Our main finding is that neural networks optimized for task (b) provide a better match
44 to the representation of non-categorical information in the monkey physiology data do those optimized
45 for task (a). This suggests that a full understanding of visual ventral stream computations might require
46 considerations other than object recognition.

47 **Materials and Methods**

48 **Dataset and Augmentation**

49 We constructed images of clothing items superimposed at random locations over natural image
50 backgrounds. To achieve this goal, we used all 70,000 images from the Fashion MNIST dataset, a
51 computer vision object recognition dataset comprised of images of clothing articles from 10 different
52 categories. We augmented this dataset by expanding the background of the image two-fold (from 28x28
53 pixels to 56x56 pixels) and drawing dx and dy linear pixel displacements from a uniform distribution
54 spanning 75% of the image field $\{-11,11\}$. Images were then shifted according the randomly drawn dx
55 and dy values. After applying positional shifts, the objects were superimposed over random patches
56 extracted from natural images from the BSDS500 natural image dataset to produce simplified natural
57 scenes which contain categorical (1 of 10 clothing categories) and non-categorical (position shifts)
58 variation. Random 56x56 pixel patches from the BSDS500 dataset were gray scaled before the shifted
59 object images were added to the background patch (Fig 1A). All augmentation was performed on-line
60 during training. That is, every position shift and natural image patch was drawn randomly every training
61 batch instead of pre-computing shifts and backgrounds. This allows every training batch to be composed
62 of unique examples from the dataset and prevents overfitting.

63 **Primate Electrophysiology**

64 Neural recordings were originally collected by the DiCarlo lab (Hong et al 2016) and shared with us for
65 this analysis. In brief, neural recordings were collected from the visual cortex of two awake and
66 behaving rhesus macaques using multi-electrode array electrophysiology recording systems (BlackRock
67 Microsystems). Animals were presented with a series of images showing 64 distinct objects from 8
68 classes rendered at varying eccentricity in the animal's visual field. After spike-sorting and quality
69 control this resulted in well-isolated single units from both IT (n=168) and V4 (n=128); higher-order
70 areas in primate visual cortex. A full description of the data and experimental methods is given by Hong
71 et al. (2016).

72 **Computational models**

73 Non-convolutional models were constructed by sequentially combining all-to-all (aka densely
74 connected) layers. Any given layer uses the previous layers' output as input, multiplying the inputs (x)
75 from by a weight matrix (w) and adds a bias to each unit in the output. Finally, this value is passed
76 through a nonlinear activation function. Each layer outputs an activation vector of its units (y) which is
77 function of its inputs (x).

$$78 \quad y = \sigma((x \cdot w) + b)$$

79 The size of each layer in the models used in our experiments were chosen to have layer sizes with
 80 roughly similar proportions to the number of output neurons in corresponding brain areas of ventral
 81 stream (Felleman and Van Essen 1991).

<i>Area</i>	<i># output neurons (10⁶)</i>	<i># layer outputs</i>	<i>Computational model</i>
		3136 (56 x 56)	Layer 0 (input image)
V1	37	3000	Layer 1
V2	29	2000	Layer 2
V4	15	2000	Layer 3
IT	10	70 (35+35)	Layer 4

82 Similar to the non-convolutional models, the convolutional models were constructed by sequentially
 83 combining convolutional layers. Each convolutional layer receives as input a spatially arranged map
 84 from the prior layer. A filter kernel is multiplied against the input at each spatial location in the input,
 85 and the resultant value is added to the bias and passed through the nonlinear activation function.

86 The convolutional models described in our paper were constructed according to the table below:

	<i>Output Size</i>	<i>Kernel Size</i>	<i>Activation Function</i>	<i>Dropout rate</i>	<i>Batch Normalization Momentum</i>
<i>Input</i>	56 x 56	N/A	N/A	N/A	N/A
<i>Layer 1</i>	28x28x16	3x3	LeakyReLU	25%	0.8
<i>Layer 2</i>	14x14x32	3x3	LeakyReLU	25%	0.8
<i>Layer 3</i>	7x7x64	3x3	LeakyReLU	25%	0.8
<i>Layer 4</i>	70(35+35)		Linear	0%	0.8

87 Models using the composite classify-reconstruct objective (see below) need an additional
 88 generator network to reconstruct the original stimulus input from the latent representation. The generator
 89 network (G) uses a residual convolutional neural network (ResNet) which has achieved state of the art
 90 performance in natural image generation.

91 The generator network uses is comprised of Deconvolutional layers and its architectural
 92 hyperparameters directly mirror those in the convolutional encoder.

93 **Objective functions and training parameters**

94 Models optimized for classification use categorical cross-entropy for the objective function. Categorical
 95 cross-entropy (XENT) is a commonly used objective function in machine learning to train neural
 96 network classifiers. Multilabel cross-entropy is calculated according to the equation below where M is
 97 the total number of classes

$$98 \quad XENT = - \sum_{c=1}^M y_c \cdot \ln(\hat{y}_c)$$

99 Here, \mathbf{y}_c is the true category label, represented as a one-hot vector, and $\hat{\mathbf{y}}_c$ is the network output
100 obtained from the linear readout of population V (see Fig. 1).

101 Models with an objective function term for reconstructing the original input scene use pixel-wise sum of
102 squared error (SSE) between the input and the generator’s output ($\hat{\mathbf{x}}$).

$$103 \quad SSE = \sum (\mathbf{x} - \hat{\mathbf{x}})^2$$

104

105 We trained each model in our experiment until classification accuracy plateaued on a validation dataset
106 of 512 objects from the 10,000 test images in the fashion MNIST dataset.

107 **Model Evaluation**

108 After training performance plateaus, 192 randomly chosen unit activations from Layers 1-3 in the
109 encoder model (Fig 1B) were used in comparisons with primate ventral stream electrophysiology. Unit
110 activations were generated using a random sample from held out test images (not used during training).
111 As in a (simulated) electrophysiology experiment, each image was input to the network, and the
112 corresponding unit activations were recorded. We then analyzed these unit activations in the same way
113 as we did the firing rates recorded in monkey visual cortex.

114 We measured selectivity of our artificial neurons in the same way as Hong et al 2016 (they call these
115 measures “performance” instead of selectivity). For continuous-valued scene attributes (e.g. horizontal
116 position) we measured selectivity as the absolute value of the Pearson correlation between the neuron’s
117 response and that attribute in the stimulus image. For categorical properties (e.g. object class) we
118 measure selectivity as the one-vs-all discriminability (d').

119 We quantified the similarity of each models’ layer-wise selectivity to corresponding layers in primate
120 ventral stream using Fisher’s Combined Probability Test (FCT). As discussed in the main paper, we first
121 used the Welch’s unpaired t-test to calculate p-values model-VS pairs for all selectivity metrics in the
122 corresponding layers, then used the FCT to combine those p-values into a single likelihood measure that
123 reflects the likelihood of observing the monkey physiology data, under the hypothesis that those data are
124 drawn from the same distribution as the units computational model: a larger p-value corresponds to a
125 model that more closely matches the monkey data.

126 **Results**

127 **Computational Models**

128 To identify the degree to which different computational objectives describe ventral stream physiology,
129 we optimized computational neural network models for different objectives, and compared them to
130 neural recordings from the primate ventral stream. Each computational model was constructed out of a
131 series of layers of artificial neurons, connected sequentially. The first layer takes as input an image \mathbf{x}
132 and at the final layer outputs a set of neuronal activities that represent the visual scene input (Fig 1B),
133 including object identity. We refer to this output as the *latent representation*. The input images, \mathbf{x} ,
134 consisted of images of clothing articles superimposed over natural image backgrounds (see Methods).
135 Each image used a single clothing article rendered in a randomly chosen position and placed over a
136 natural image background (Fig. 1A).

137 The models each had a total of three layers of processing (corresponding to cortical areas V1, V2, and
138 V4) between their inputs and these latent representations; the latent representations correspond to area
139 IT, for reasons we discuss below. The visual inputs to the model had normalized luminance values,
140 mimicking the normalization observed in thalamic inputs to V1¹¹. The connectivity between neurons in
141 each layer (and the artificial neurons' biases) were optimized within each model, to achieve the specified
142 objective (see Methods). We repeated this process for two different objectives, yielding two different
143 types of models.

144 The first type of model was optimized strictly for object recognition: the optimization maximized the
145 ability of a linear decoder to determine the identity of the clothing object in the visual scene from the
146 latent representation. (This mirrors the observation that neural activities in area IT can be linearly
147 decoded to recover object identity¹²). The second type of model was optimized for two tasks in parallel:
148 the ability of a linear decoder to determine object identity from latent representation, *and* the ability of a
149 decoder to reconstruct the object from the latent representation. (See Methods for details about the
150 optimization procedure). We repeated this procedure with both convolutional, and non-convolutional
151 neural network architectures, yielding a total of four models (Fig 1C).

152 In all cases, the models were optimized using sets of images containing randomly sampled objects, until
153 their object classification performance saturated on a set of held-out validation images. Good
154 performance on the categorization task was obtained in all models (Fig 1D). Having developed models
155 optimized for these different objectives, we could evaluate how well each model matched observations
156 from primate VS, and use that comparison to determine which computational objective provides the best
157 description of primate VS.

158 **Electrophysiology Comparisons**

159 To compare our neural network models to ventral stream physiology, we used the experimental data
160 from a previously-published study^{9,12} (see methods and ref^{9,12} for details). These data consisted of
161 electrode array recordings from areas V4 and IT of monkeys that were viewing images; many neurons in
162 each area were simultaneously observed. Within these data, we assessed each neuron's selectivity for
163 object identity, and for category-orthogonal image properties (e.g. horizontal object position), as in
164 Hong et al⁹ (see methods). We performed this analysis for the monkey data, and for the artificial neurons
165 in each layer of each of our computational models. We then compared the trends in image property
166 selectivity displayed by non-human primate VS neurons and units from each of our models along the
167 visual processing pathway.

168 In the primate VS, selectivity for both categorical and category-orthogonal scene attributes increased
169 along the ventral stream (Fig 2A), as reported by Hong et al⁹. This indicates that both types of attributes
170 are more explicitly represented in progressively deeper ventral stream areas.

171 Within our computational models, those models optimizing the composite objective showed the same
172 trends observed in primate ventral stream neurons (Fig 2C, 2E): both category and category-orthogonal
173 properties of the visual scene are represented more explicitly with each subsequent layers of the model.
174 This observation persisted for both the convolutional and the non-convolutional architectures. For
175 contrast, models optimized solely for object recognition (without the image reconstruction component of
176 the objective function) did not show consistent increases in position selectivity along the visual pathway
177 (Fig 2B, 2D). Again, this observation held for both convolutional and non-convolutional model
178 architectures.

179 Thus, models optimizing the composite objective function qualitatively recapitulate the trends in
180 neuronal selectivity along the visual pathways better than do models optimized strictly for object
181 recognition. This observation motivated us to quantify how well each model matched the primate VS
182 data. To achieve this goal, we performed the following analysis on each computational model. First, we
183 used unpaired t-tests to estimate the probability that there is no difference in object category selectivity
184 between the primate IT data and the model's latent representation. We then performed a t-test comparing
185 the primate V4 category selectivity to the corresponding layer of the computational model. Next, we
186 performed t-tests comparing the horizontal, and vertical, position selectivities in primate V4 and IT to
187 the corresponding layers of the computational model. This procedure yielded 6 p-values, describing the
188 probability that the model matched each of these attributes observed in the primate VS. Finally, we used
189 Fisher's method¹³ to combine those 6 p-values into a single number, that quantified the likelihood of
190 there being no difference between the computational model and the primate VS.

191 Comparing these likelihood values, we found that the convolutional models overall provided better
192 descriptions of the primate VS than did the non-convolutional ones (i.e., they had higher likelihood
193 values), and that the best model overall was the convolutional neural network optimized for the
194 composite classify-and-reconstruct objective (See Supplemental Fig. 1).

195 **Noise Robustness**

196 We found that the convolutional model, optimizing the composite objective (classify-and-reconstruct)
197 best matched the depth-dependent increase in position selectivity seen in single unit activities recorded
198 from primate ventral stream. This led us to ask whether there might be functional benefits for networks
199 optimizing this composite objective function, as compared with ones that are just trained to classify their
200 inputs.

201 Further motivating this question, we note that previous work has shown that convolutional neural
202 networks optimized for object recognition tend to perform poorly on object recognition tasks when the
203 images are corrupted by noise. Specifically, classification performance has been seen to decrease
204 significantly when networks are evaluated under noise conditions even marginally different from the
205 conditions under which it was trained¹⁴. This is different from the primate visual system, where object
206 recognition performance is more robust to image noise, leading us to speculate that the convolutional
207 networks trained for the composite classify-and-reconstruct task – which provide the best match to
208 primate VS data – might have classification performance that is more robust to image corruption than do
209 the networks trained purely for object recognition.

210 To test that hypothesis, we took each of our previously trained models, and measured their accuracy at
211 categorizing the clothing objects in test images corrupted by increasing levels of additive pixel noise
212 (see methods). Similar to previous work, the convolutional model trained purely for object recognition
213 showed a decrease in performance as the noise level increased. For the convolutional model trained on
214 the composite task, the decrease in performance with increasing noise level was less severe. This
215 suggests that, consistent with our hypothesis, there is a functional benefit to systems optimizing the
216 composite objective over “pure” object recognition systems: their object recognition performance is
217 more robust to noise.

218 The same finding also holds for the non-convolutional model architectures, and they are overall more
219 robust to image noise than are the convolutional ones. We repeated this analysis with multiplicative
220 (instead of additive) pixel noise (see Supplemental Fig. 2) and demonstrate that our findings can be
221 generalized across multiple noise types.

222 Discussion

223 Here we report evidence that convolutional neural networks (DCNNs) optimizing a two-part composite
224 objective (recognize and visualize) describe the depth-dependent evolution of categorical and non-
225 categorical information in primate VS better than do networks optimized for object recognition alone.
226 This is unexpected, as prior work posits that networks optimized strictly for object recognition should
227 form the best models of primate VS.^{2,4,9,10} Our results suggest that the evolution of category-orthogonal
228 information along the visual pathway could require a different functional explanation. Moreover,
229 consistent with previous work,^{2,4,9,10} our CNNs optimized for image classification resemble primate VS
230 more closely than do non-convolutional models optimizing the same objective.

231 Our findings may help reconcile discrepancies between the object recognition hypothesis of VS and
232 results which appear at odds with this interpretation,^{6,7,15,16} for example the finding that primate VS
233 explicitly retains information not useful for object recognition experiments tested previously.⁹ The
234 composite objective promotes retention of both category and category-orthogonal information because
235 both are necessary to reconstruct the stimulus.

236 Importantly, we used a different method to compare our neural networks to the primate VS than have
237 previous studies that compared the representational dissimilarity matrices (RDMs) for their models, with
238 those of the primate VS.^{2,4,9,10} While RDMs assay the (dis)similarity in how different images are
239 represented by the models, or primate VS, recent work suggests RDM analysis may be insufficient as a
240 universal metric of model similarity¹⁷; especially when the model cannot be trained using identical
241 image datasets (as in our case). Instead, our approach was to focus on the depth-dependent evolution of
242 neuronal selectivity to categorical and non-categorical variations in the input images. Our conclusion --
243 that an objective other than pure object categorization could best describe the computations in primate
244 VS -- differ from prior studies and further suggest that aspects of visual computation are not fully
245 captured by RDM analysis.

246 Furthermore, our findings suggest noise tolerance as another independent explanation for why the VS
247 might use a composite computational objective. VS classification accuracy measured in humans
248 tolerates noise corrupted images much better than DCNNs optimized for image classification alone¹⁴. In
249 contrast, convolutional models optimizing the composite objective demonstrate better noise tolerance
250 compared to identical models trained solely for classification (Fig 3). Importantly, improved noise
251 tolerance occurs without having to augment training images with noise. These findings complement the
252 expanding body of work to explain the neuronal computations in visual processing and have applications
253 in the computer vision models that emulate them.

254 Acknowledgements

255 We would like to thank the DiCarlo lab for sharing their primate electrophysiology recordings with us.
256 Special thanks to Alon Poley-Polsky for thoughtful discussion and direction and to Doug Crawford,
257 Shaiyan Keshvari, Martin Schrimpf, Rachel Sewell, and Heidi Sjoberg for helpful feedback on the
258 manuscript. EC acknowledges funding by an NDSEG fellowship through the US Department of
259 Defense. JZ acknowledges funding from CIFAR, the A.P. Sloan Foundation, Google, the Canada
260 Research Chairs Program, and the Natural Sciences and Engineering Research Council of Canada
261 (NSERC).

262 **References**

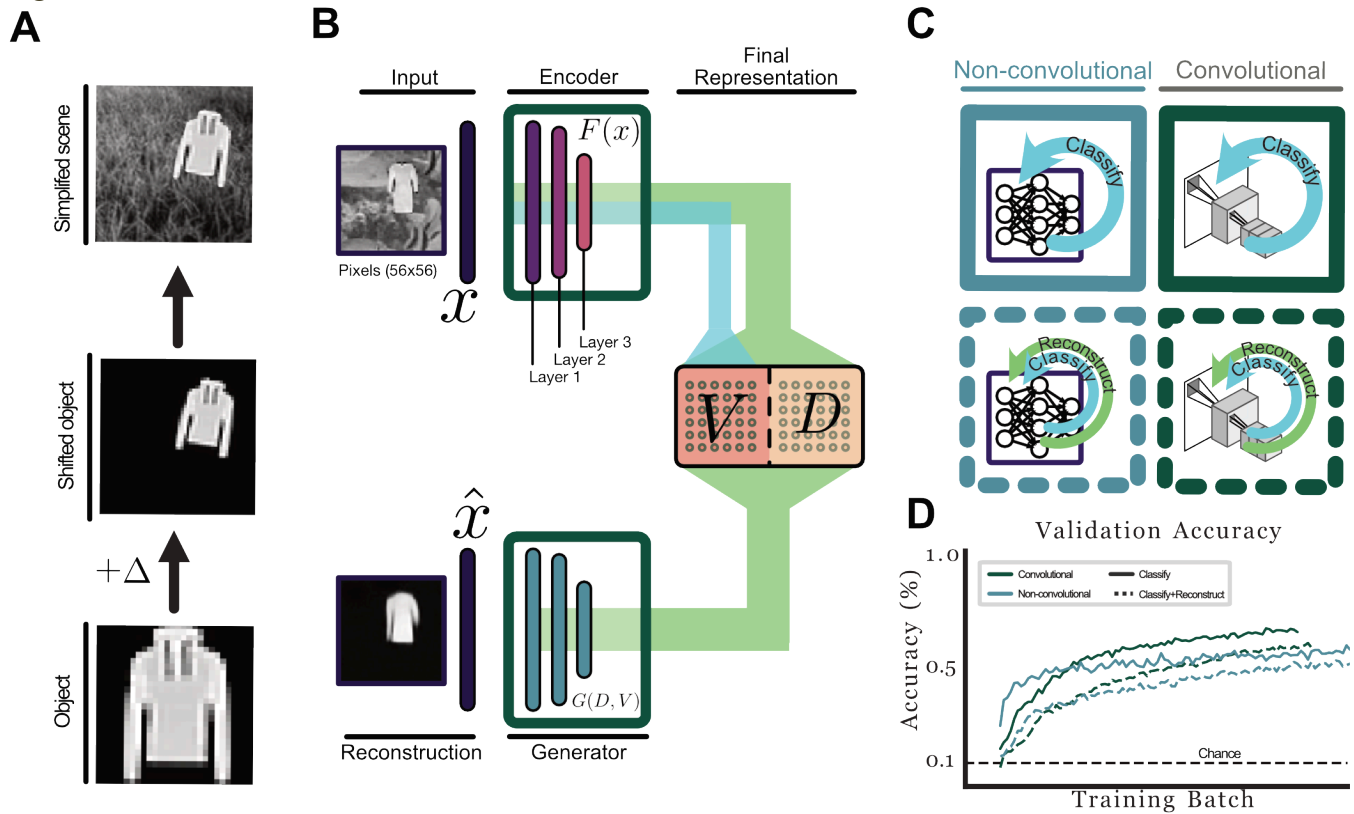
263

- 264 1. Felleman, D.J. & Van Essen, D.C. *Cereb. Cortex* **1**, 1–47 (1991).
- 265 2. Yamins, D.L.K. & DiCarlo, J.J. *Nat. Neurosci.* **19**, 356–365 (2016).
- 266 3. Yamins, D.L.K. et al. *Proc. Natl. Acad. Sci. U.S.A.* **111**, 8619–8624 (2014).
- 267 4. Cadieu, C.F. et al. *PLoS Comput. Biol.* **10**, e1003963 (2014).
- 268 5. Güçlü, U. & van Gerven, M.A.J. *J. Neurosci.* **35**, 10005–10014 (2015).
- 269 6. Stokes, M., Thompson, R., Cusack, R. & Duncan, J. *J. Neurosci.* **29**, 1565–1572 (2009).
- 270 7. O'Craven, K.M. & Kanwisher, N. *J Cogn Neurosci* **12**, 1013–1023 (2000).
- 271 8. Chen, Y. & Crawford, J.D. *Annals of the New York Academy of Sciences* **46**, 774 (2019).
- 272 9. Hong, H., Yamins, D.L.K., Majaj, N.J. & DiCarlo, J.J. *Nat. Neurosci.* **19**, 613–622 (2016).
- 273 10. Richards, B.A. et al. *Nat. Neurosci.* **22**, 1761–1770 (2019).
- 274 11. Carandini, M. & Heeger, D.J. *Nat Rev Neurosci* **13**, 51–62 (2011).
- 275 12. Majaj, N.J., Hong, H., Solomon, E.A. & DiCarlo, J.J. *J. Neurosci.* **35**, 13402–13418 (2015).
- 276 13. Li, Q., Hu, J., Ding, J. & Zheng, G. *Biostatistics* **15**, 284–295 (2014).
- 277 14. Geirhos, R. et al. (2018).
- 278 15. Freud, E., Plaut, D.C. & Behrmann, M. *Trends in Cognitive Sciences* **20**, 773–784 (2016).
- 279 16. Sereno, A.B. & Lehky, S.R. *Front Comput Neurosci* **4**, 159 (2011).

- 280 17. Rezai, O., Stoffl, L. & Tripp, B. *Neural Netw* **121**, 122–131 (2020).

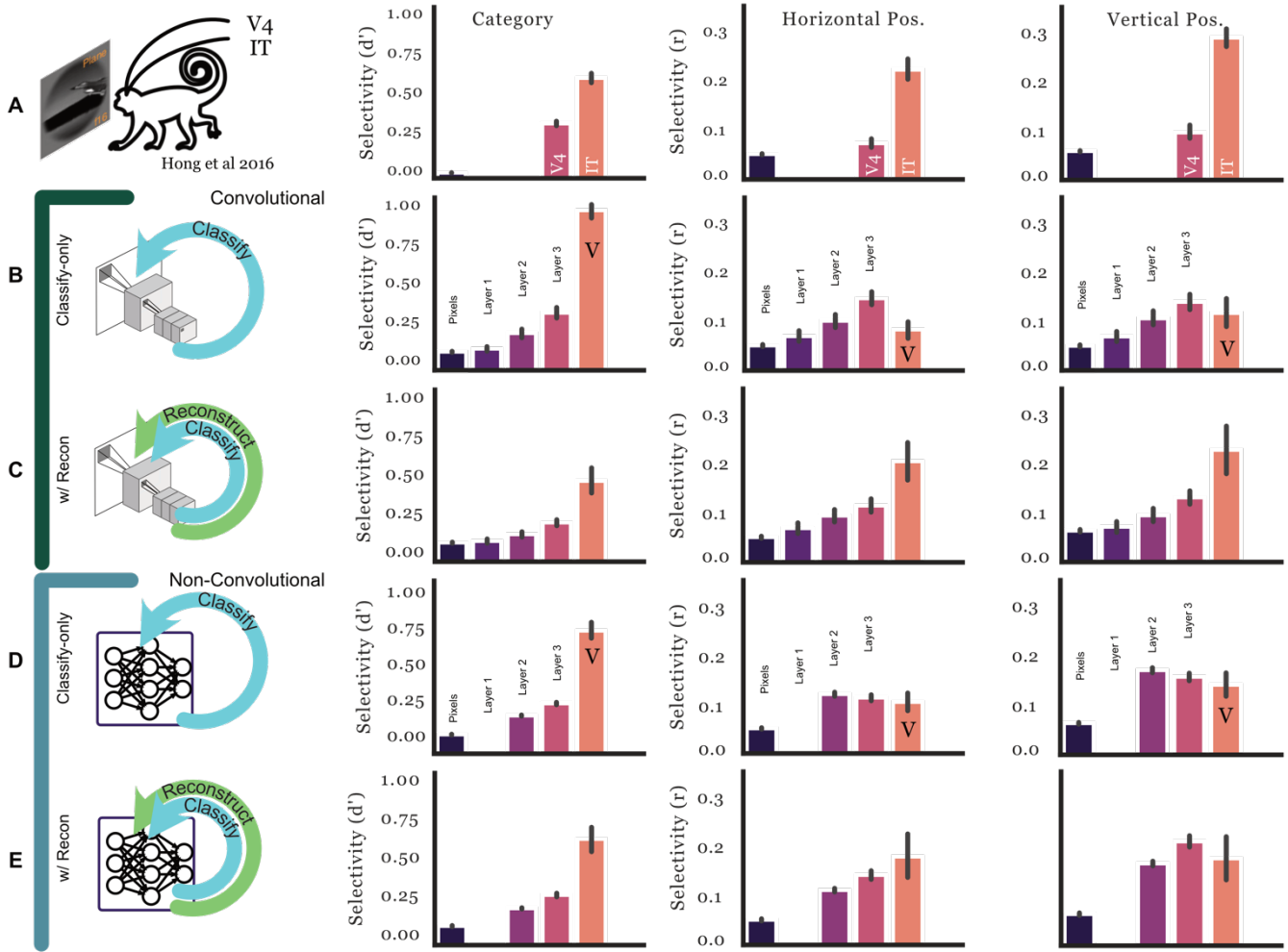
281

282



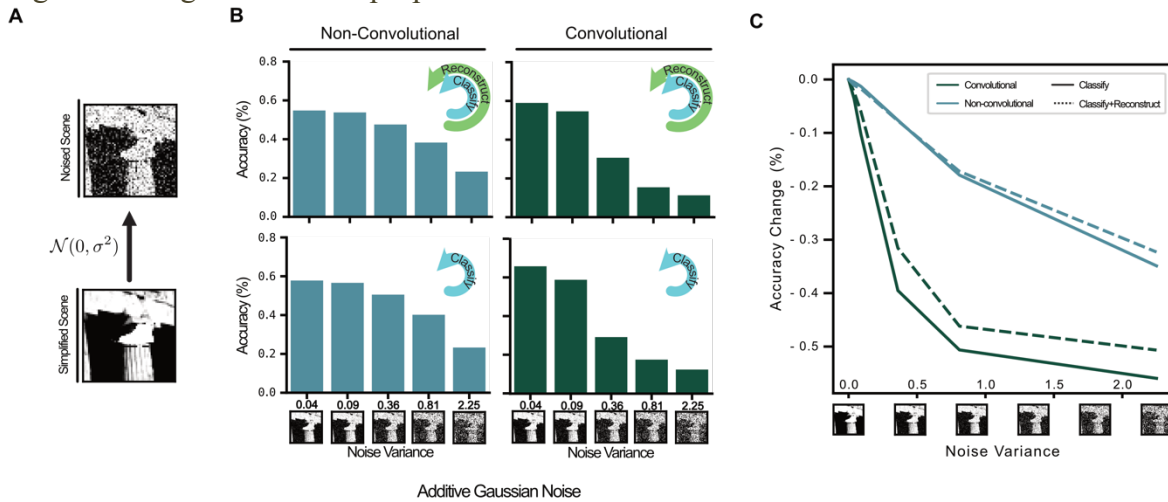
284

285 **A)** We constructed images of clothing items superimposed over natural image backgrounds at random
 286 eccentricities. **B)** We model the ventral stream as an encoder whose objective is to map input image (x)
 287 onto more abstract “latent” representations (D and V). In our models this entire latent space is
 288 represented by 70 artificial neurons (35 units in each of D and V). The generator network uses these
 289 latent representations (D and V) as input to reconstruct the object and its location within the scene. A
 290 separate linear decoder attempts to determine the object identity from the activities of the units in V . **C)**
 291 We trained both convolutional, and non-convolutional neural network architectures, on one of two tasks:
 292 object categorization (“classify”), or object categorization with concurrent image reconstruction. We
 293 note that, for the “pure” object recognition task, the generator network is superfluous. **D)** Neural
 294 networks with both architectures achieve comparable object recognition performance (accuracy) when
 295 using either classify-only and classify+reconstruct objective functions. This performance was assessed
 296 on held-out images, not used in training the networks.



298
 299 **A)** Category and position selectivity of single units recorded from macaque ventral stream (see Methods
 300 and Hong et al. 2016). **B&C)** Selectivity of units in the fully trained convolutional models optimized
 301 under classify-only objective (categorical cross-entropy) and the composite classify+reconstruct
 302 autoencoder objective. **D&E)** Non-convolutional or “all-to-all” models were also trained on both
 303 classify-only and classify+reconstruct. We measured property selectivity of both categorical and
 304 continuous valued category-orthogonal properties on units in the multi-electrode array data and each
 305 layer of the computational model encoders. We defined selectivity for categorical information on each
 306 unit in the dataset as the absolute value of that unit’s discriminability (one-vs-all d-prime). We defined
 307 selectivity for continuous valued attributes (horizontal and vertical position) on each unit as the absolute
 308 value of the Pearson correlation coefficient. Unit activities for models were sampled using 10000 held
 309 out test images to generate activations at each layer of the model. For layers containing more than 192
 310 units we randomly sampled 192 units for the analysis (to have a number of units similar to the number
 311 of IT units in the neural recordings).

312 Fig. 3: Noise generalization properties of models



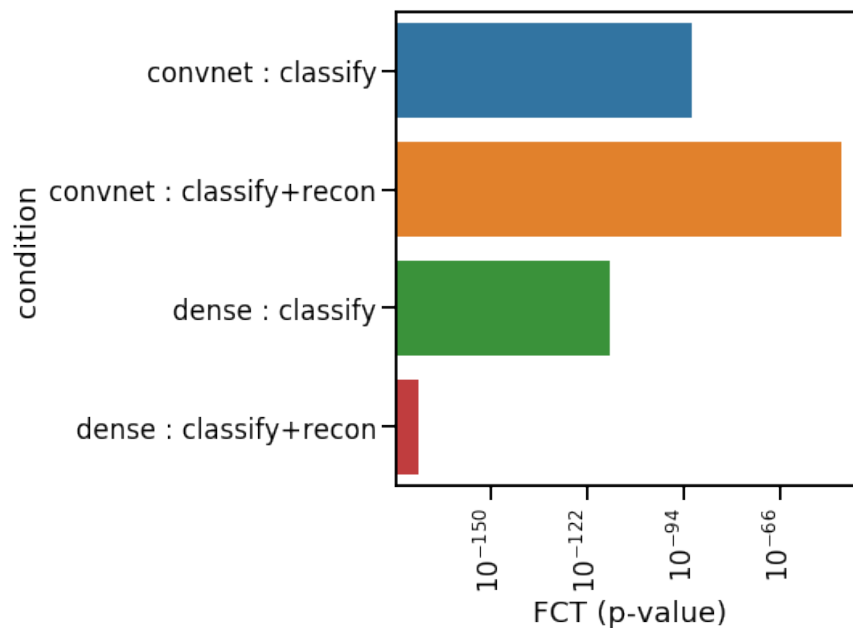
313 **A)** Additive gaussian noise (mean=0) was used to corrupt 10,000 testing images at increasing levels. **B)**
 314 Each model (defined its architecture – convolutional or non-convolutional -- and the objective on which
 315 it was trained) was evaluated on images corrupted with increasing levels of gaussian noise. We show the
 316 accuracy at categorizing the objects in the noise-corrupted images. These images were from a held-out
 317 dataset, not used in training the neural networks. **C)** Convolutional neural networks are more sensitive to
 318 noise than are non-convolutional ones; they show a larger decrease in accuracy with increasing noise
 319 variance. Adding a reconstruction component to the network objective reduces this sensitivity. Similar
 320 results were obtained with a multiplicative noise model (Fig. S2), indicating that this result is not
 321 sensitive to the specific type of noise that corrupts the images.
 322

323

324 **Supplemental**

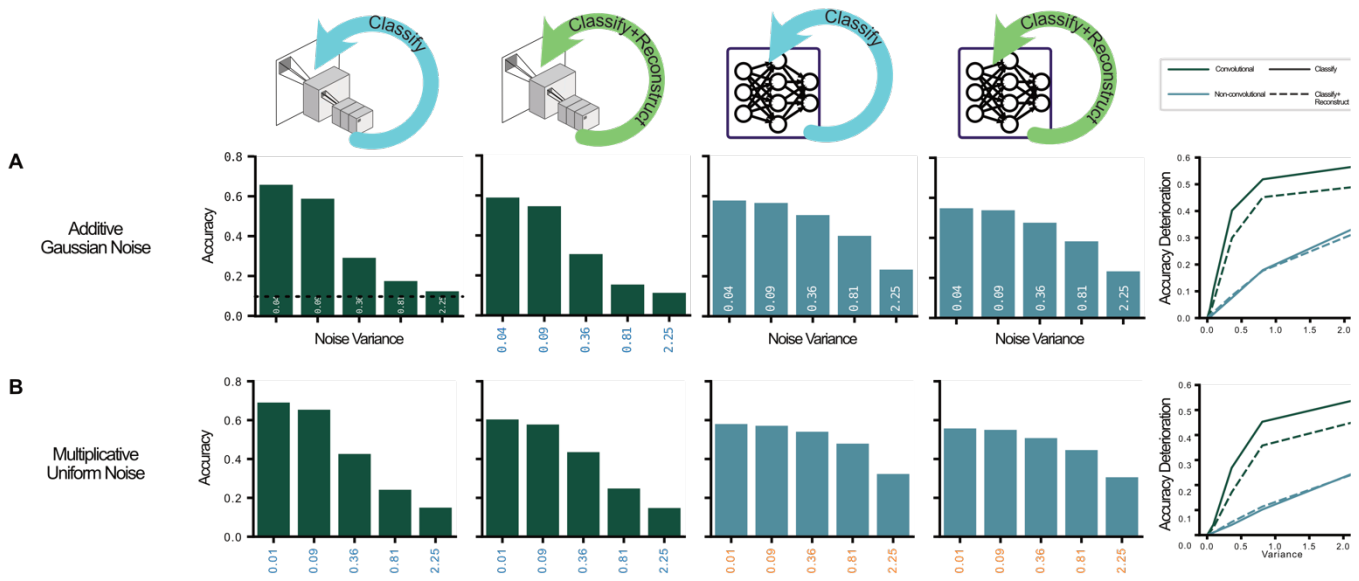
325

326 Supplemental Fig. 1: Fisher combined probability test (FCT).



327

328 We used the FCT to compute the likelihood of each model's category and position selectivity matching
329 the data observations made in monkey ventral stream recordings. Those likelihoods (p-values) are
330 shown for each model. Higher p-values (taller bars) correspond to models that more closely match the
331 neural data.



333 Each model (defined its architecture – convolutional or non-convolutional -- and the objective on which
 334 it was trained) was evaluated on images corrupted with increasing levels of noise. **A)** *Additive* gaussian
 335 noise (mean=0) was used to corrupt 10,000 testing images at increasing levels. **B)** *Multiplicative*
 336 noise (uniform noise) was used to corrupt 10,000 testing images at increasing levels. Bar plots show the
 337 accuracy of each neural network model at categorizing the objects in those noisy images. **C)** We show
 338 the deterioration in accuracy at each noise level, for each model. This comparison shows that the
 339 convolutional neural networks are more sensitive to noise but adding a reconstruction objective appears
 340 to improve this sensitivity.
 341