

Zero-shot neural decoding of visual categories without prior exemplars

Thomas P. O’Connell ^{1,2,3,*}, Marvin M. Chun ^{4,5}, Gabriel Kreiman ^{3,6}

1 Department of Brain and Cognitive Sciences, MIT

2 McGovern Institute for Brain Research, MIT

3 Center for Brains, Minds, and Machines, MIT

4 Department of Psychology, Yale University

5 Department of Neuroscience, Yale University

6 Children’s Hospital, Harvard Medical School

* tpo@mit.edu

Abstract

Decoding information from neural responses in visual cortex demonstrates interpolation across repetitions or exemplars. Is it possible to decode novel categories from neural activity without any prior training on activity from those categories? We built zero-shot neural decoders by mapping responses from macaque inferior temporal cortex onto a deep neural network. The resulting models correctly interpreted responses to novel categories, even extrapolating from a single category.

1 Neural decoding approaches typically train machine learning classifiers on
2 responses to a set of stimuli and subsequently test the classifier using either
3 different repetitions of the same training stimuli or responses to different exem-
4 plars from the same training categories. These approaches have been extremely
5 successful in a wide variety of domains [1], but show limited generalization.

6 Zero-shot neural decoding, or interpreting neural activity without prior
7 exposure to any similar information [2–6], holds great promise to improve
8 the generalizability of neural information processing models. While standard
9 decoders predict information directly from patterns of neural activity, zero-shot
10 decoders map neural activity to an intermediate representation that constitutes
11 a computational hypothesis for the neural code [2]. The intermediate represen-
12 tation is selected such that it has a known or easily learned relationship to a

13 wide variety of to-be-predicted outputs. In an impressive recent demonstration
14 of zero-shot decoding, Anumanchipalli and colleagues [6] reconstructed recog-
15 nizable human speech from electrophysiological recordings in human motor
16 cortex via a computational model of articulatory movement. Even though the
17 decoding model was only trained to map neural activity to the articulatory
18 model, and not representations of words or semantics, the models could recon-
19 struct intelligible human speech. Here, we demonstrate such zero-shot decoding
20 from electrophysiological responses for visual objects.

21 Beyond a feat of engineering, the degree of generalization has important
22 consequences for the conclusions that can be drawn from a model of neural
23 information processing. The greater the generalization, the stronger the evi-
24 dence that a model captures generic processing beyond any particular set or
25 class of stimuli. As an example, consider a standard linear decoder trained to
26 distinguish whether responses along the ventral stream were evoked by images
27 of airplanes or chairs. The decoder could interpolate within its training space
28 to label neural responses to new images of airplanes or chairs, but it would
29 not be able to accurately label neural responses to cars or tables. A zero-shot
30 model can capture generic visual information and extrapolate to new categories
31 on which it was not trained.

32 Constructing generic zero-shot decoders for visual objects necessitates a
33 model for visual processing in the primate brain. How well do we understand
34 the neural code for visual object processing? Deep convolutional neural net-
35 works (DCNNs) constitute a promising initial approximation to the cascade
36 of computations along the ventral stream that support visual object recogni-
37 tion [7–11]. DCNNs are goal-directed, hierarchical, image-computable models
38 capable of recognizing complex, natural objects and scenes [12], and repre-
39 sentations in DCNNs predict object-evoked neural activity in *rhesus macaque*
40 inferior temporal cortex (IT) [13–15], which is at the top of the ventral visual
41 stream hierarchy and plays a central role in visual object recognition [16, 17].

42 While DCNNs are powerful pattern extractors, it remains possible that their
43 performance predicting IT responses is driven by generalization within stimuli
44 (e.g., different views of the same chair) or within categories (e.g., one type
45 of chair to another). To test whether DCNNs capture the type of flexible
46 visual processing accomplished by biological vision, the mapping from DCNNs
47 to IT should generalize across object categories (e.g. chairs to cars). While
48 some studies have shown extrapolation across categories [3, 13], the degree
49 of generalization remains unclear. In an extreme case, can IT to DCNN
50 mappings learned from neural activity evoked by a single object category
51 extrapolate to novel categories? If the mappings generalize to new images from
52 the same category, we can conclude that DCNN responses capture category-
53 level information within IT. If the mappings generalize to new images from
54 novel categories, this suggests that DCNNs capture generic visual information
55 in IT beyond any one category.

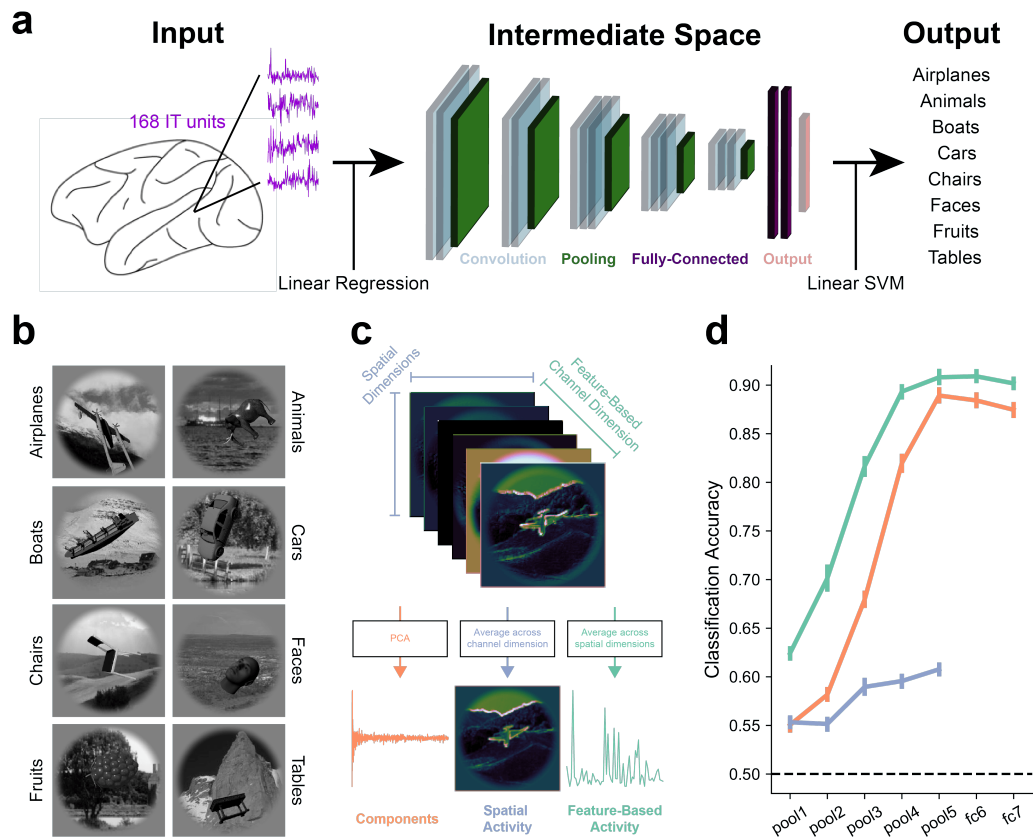


Figure 1. Overview of zero-shot approach and deep convolutional neural network (DCNN) architecture. **a.** Overview of zero-shot decoding pipeline. IT recordings were mapped to an intermediate space defined as unit activity in a deep convolutional neural network trained for object categorization. Pre-learned mappings from DCNN unit activity to object categories were used to generate predictions from DCNN-aligned IT recordings. The decoders are zero-shot if neural recordings from the test categories are withheld when learning the IT to DCNN. **b.** Example images from the eight object categories: Airplanes, Animals, Boats, Cars, Chairs, Faces, Fruits, Tables [18]. **c.** VGG-16 trained for object categorization on the ImageNet dataset [19] was used as the intermediate basis-space to which IT recordings were aligned. In the pooling layers, DCNN activity is organized along two spatial dimensions and a feature-based channel dimension. To summarize the full spatio-featural activity-space, we encoded DCNN activity into 1000 components using principle components analysis (left, orange). To isolate spatial activity, we averaged across the channel dimension in DCNN activity (center, blue). To isolate spatially invariant feature-based activity, we averaged across the spatial dimensions (right, green). **d.** Object category (binary prediction between all 28 combinations of the eight categories, chance = 50%) could be predicted from all types of DCNN activity, with components and feature-based activity showing a sharp rise in decodability across layers. This plot does not include any neural data, it only shows decodability using DCNN activity.

56 To determine whether representations in DCNNs capture generic visual
57 processing in the primate brain, we built zero-shot neural decoders for object
58 category from multi-electrode array recordings in *rhesus macaque* IT (**Fig. 1a**).
59 IT responses were evoked by images of computer-generated objects on natural
60 scene backgrounds with high variation in position, size, and orientation (**Fig.**
61 **1b**). We tested whether zero-shot decoders trained on neural responses to a
62 set of categories (e.g., airplane and chair images) can accurately label neural
63 activity evoked by novel categories (e.g., cars and tables). In the most extreme
64 instance, we tested whether zero-shot decoders trained on neural responses
65 from a single category can generalize to label neural responses evoked by seven
66 novel categories.

67 First, we defined an intermediate space based on unit activity in a DCNN
68 trained for object categorization (VGG-16 pre-trained on ImageNet, **Supple-**
69 **mental Fig. 1**) [19]. For the pooling layers, we analyzed three features
70 computed from DCNN activity: principal component scores, feature-based
71 channel activity, and spatial activity (**Fig. 1c**). For the fully-connected layers,
72 we analyzed two DCNN features: principal components and feature-based
73 channel activity. Within each layer, each DCNN feature for each layer was
74 mapped to object category using a series of 8-choose-2 linear support vector
75 machine (SVM) classifiers to make binary predictions of object category. Im-
76 portantly, the training of these SVMs did not incorporate any neural data.
77 Object category could be predicted from each of the seven VGG-16 layers
78 for all three DCNN features (**Fig. 1d**). For components and feature-based
79 activity, and to a lesser extent for spatial activity, performance improved for
80 progressively deeper layers relative to earlier layers.

81 Next, we mapped IT activity to DCNN activity using linear regression (20
82 category-matched 75% train, 25% test folds). The resultant IT to DCNN
83 transformation matrices were multiplied by IT activity vectors from the test
84 set to transform IT activity into the same space as DCNN component scores.
85 The transformed IT activity was multiplied by the transpose of the PCA
86 transformation matrix to reconstruct the full DCNN activity space for a given
87 layer. As with true DCNN activity, feature-based and spatial reconstructions
88 were computed for the five pooling layers, and the full 4096 reconstructed
89 channel activations were analyzed for fully-connected layers. This procedure
90 was repeated separately for each layer and cross-validation fold (**Supplemental**
91 **Fig. 2**).

92 Finally, to determine whether the mapping between IT and DCNN activity
93 extrapolates to novel categories, IT activity from the test categories was held
94 out during training of the mapping. If decoding accuracies when recognizing
95 novel categories are high, this indicates that the IT to DCNN mapping is
96 zero-shot, capturing generic visual information in IT to generalize to novel
97 categories on which it was never trained.

98 To assess the overall amount of shared generic visual information between

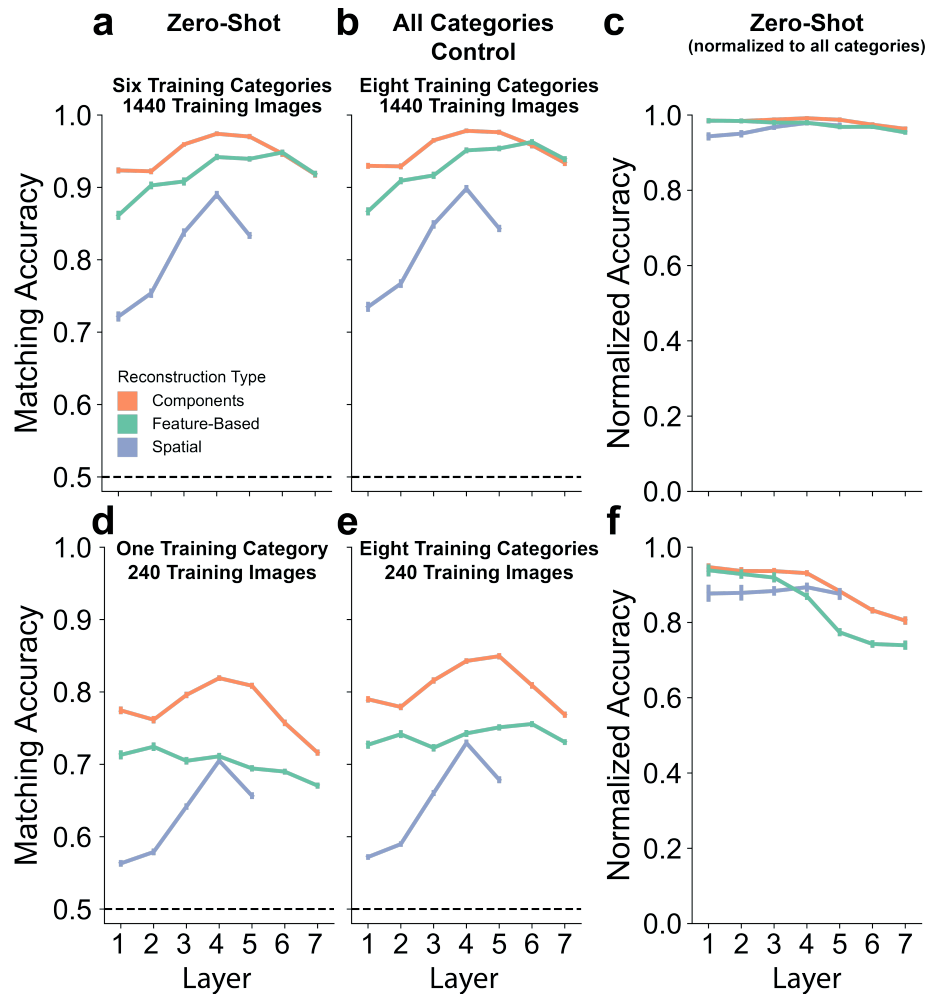


Figure 2. DCNN features reconstructed from IT activity match true DCNN features for the same images even when extrapolating across categories. **a.** High matching accuracies were achieved when neural activity from two test categories was held out during training (six training categories), indicating that IT to DCNN mappings indeed generalize across object category. **b.** As a control, a decoder was trained on neural responses from all eight categories, matching the overall number of training images to the number used for the zero-shot decoder. Strikingly similar results were obtained for the zero-shot and all categories control decoders. **c.** Zero-shot matching accuracies were normalized to calculate proportion of above-chance matching accuracy achieved by the zero-shot model relative to the all categories model. Normalized matching accuracies are close to ceiling for all layers and reconstruction types. **d, e, f.** Zero-shot, all categories control, and normalized results when only one category was used to train the zero-shot decoder. Again, the zero-shot decoder displayed highly similar results to the all categories control decoder. All accuracies for all decoders, DCNN feature-types, and layers are significant at $P < 0.001$ (permutation testing).

99 IT and DCNN activity, reconstructed DCNN features were matched to the true
100 DCNN features from the same images, relative to the true DCNN features from
101 every other image in the test set, for all possible pair-wise combinations. We
102 measured the matching accuracy, which ranges from 0.5 (chance) to 1 (perfect
103 reconstruction of features) (**Fig. 2**). We examined the extreme cases where
104 the maximum number of available training data (six categories) and minimum
105 number (one category) were used.

106 When the zero-shot decoder was trained on neural responses from six cate-
107 gories, matching accuracies for reconstructed DCNN components, feature-based
108 activity, and spatial activity were all well above chance (all $P_s < 0.001$, permu-
109 tation testing), with components and feature-based activity exhibiting higher
110 matching accuracies than spatial activity (**Fig. 2a**). As a control, we compared
111 the zero-shot decoder's performance to a decoder trained on responses from
112 all categories, matching the overall number of training images to the number
113 used for the zero-shot decoder. This all categories control decoder displayed
114 strikingly similar results to the zero-shot decoder (**Fig. 2b**, all $P_s < 0.001$,
115 permutation testing). A normalized matching accuracy, the proportion of
116 above-chance matching accuracies achieved by the zero-shot decoder relative to
117 the all categories control decoder, was calculated as $(\text{matching accuracy}_{\text{Zero-Shot}} - \text{chance}) / (\text{matching accuracy}_{\text{AllCategories}} - \text{chance})$. A normalized matching
118 accuracy of 0 indicates zero-shot performance was at chance, a value of 1
119 indicates the zero-shot matching accuracy was equal to the matching accuracy
120 for the model trained on all categories. Normalized matching accuracies were
121 all close to ceiling (**Fig. 2c**), indicating strikingly comparable accuracies
122 between the zero-shot and all categories control decoders.

124 Even when the zero-shot decoder was trained on neural responses from just
125 one category, matching accuracies remained well above chance (**Fig. 2d**, all P_s
126 < 0.001 , permutation testing). Again, after computing normalized accuracies
127 relative to an all categories control decoder matched for the number of training
128 images (**Fig. 2e**), we find normalized accuracies close to ceiling (**Fig. 2f**),
129 demonstrating strong generalization of the IT to DCNN mappings. Full
130 matching results for all possible numbers of training categories can be seen in
131 **Supplemental Fig. 3**.

132 Next, we assessed whether the information captured in the IT to DCNN
133 mappings is discriminative of object categories. For this purpose, we used the
134 same 8-choose-2 SVM classifiers, trained on DCNN activity (**Fig. 1d**), *without*
135 *retraining or fine-tuning on neural data*, to generate object category predictions
136 from IT-reconstructed DCNN features (**Fig. 3**). In other words, the model
137 presented with IT responses from novel categories and the task is decode which
138 category was presented to the monkey, even though the mapping to reconstruct
139 DCNN activations from IT was never exposed to neural responses from that
140 particular category.

141 When the zero-shot decoder was trained on six categories, we see significant

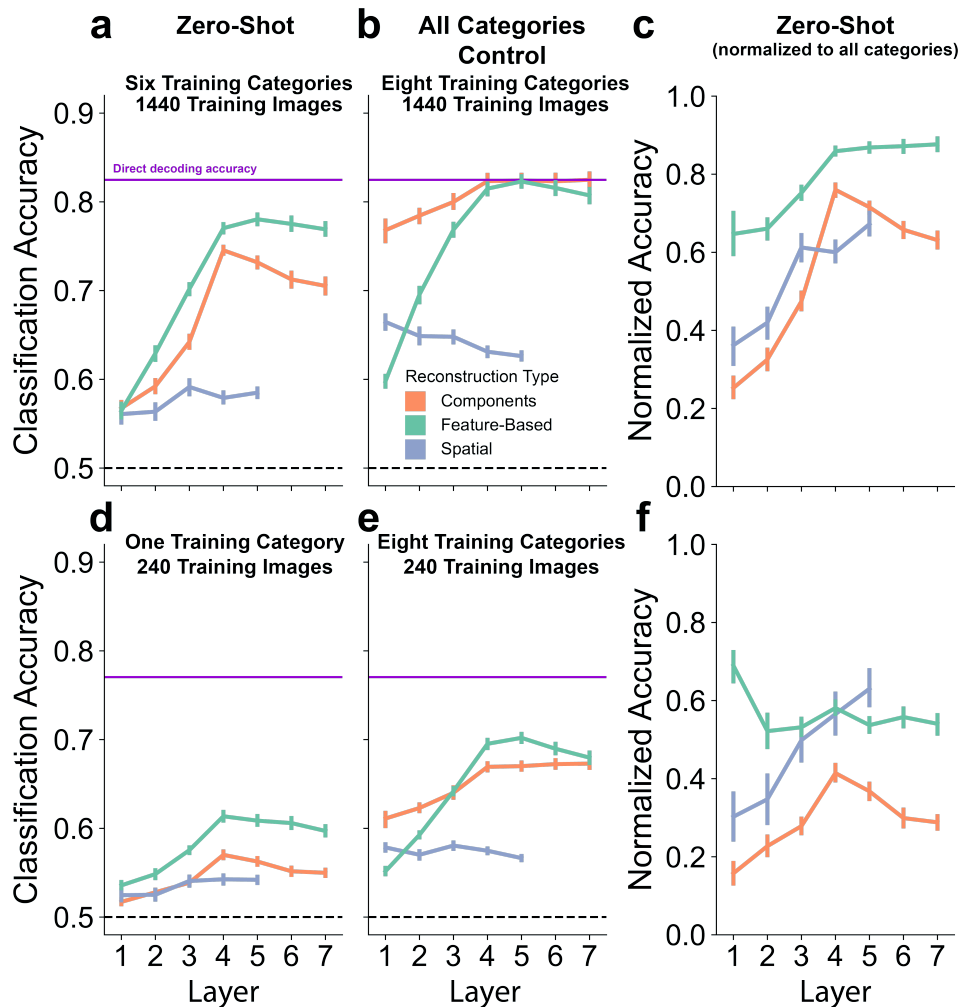


Figure 3. Object category can be decoded from IT-reconstructed DCNN features even when extrapolating to novel categories that were not used for training the IT to DCNN mapping. **a.** High classification accuracies were achieved when neural activity from two test categories was held out during training (six training categories), indicating that IT to DCNN mappings capture generic information about visual object category. Direct decoding accuracies (predicting category directly from IT responses using linear SVMs) are shown in purple. **b.** A control decoder, where the IT to DCNN mapping was learned using IT responses from all eight categories, showed a similar pattern of results, albeit with higher accuracies. **c.** Zero-shot classification accuracies were normalized to calculate proportion of above-chance classification accuracy achieved by the zero-shot model relative to the all categories model. Normalized classification accuracies were all greater than zero, and the best normalized accuracies for feature-based reconstructions achieved over 80% of the accuracies seen for the all categories control decoder. **d, e, f.** Zero-shot, all categories control, and normalized results when only one category was used to train the zero-shot decoder. Again, the zero-shot decoder displayed a similar pattern of results to the all categories control decoder but with lower accuracies. All accuracies for all decoders, DCNN feature-types, and layers are significant at $P < 0.001$ (permutation testing).

142 prediction for all reconstructed DCNN features and layers (**Fig. 3a**, all Ps
143 <0.001 , permutation testing). Feature-based reconstructions produced the
144 best zero-shot predictions and were well above chance. The feature-based
145 reconstructions show an increase in zero-shot prediction accuracy over the
146 first few layers before flat-lining, indicating an expected increase in shared
147 information about object category between IT and the DCNN across layers.
148 Again, we compare the zero-shot performance to a control decoder where the
149 IT to DCNN mapping was learned on responses from all categories (**Fig. 3b**,
150 all Ps <0.001 , permutation testing)) by calculating normalized classification
151 accuracies. These normalized classification accuracies were all well above
152 0, with the feature-based reconstructions from the later layers achieving a
153 proportion over 0.8 of the accuracies seen for the all categories decoder (**Fig.**
154 **3c**). These results demonstrate that information captured about the neural
155 code for object representation when learning the mapping from IT to DCNN
156 activity extrapolates to untrained novel object categories.

157 When a single category was used to train the zero-shot decoder, accuracies
158 were still all significantly above chance (**Fig. 3d**, all Ps <0.001 , permutation
159 testing). After normalizing these accuracies to an all categories control model
160 matched for training set size (**Fig. 3e**, all Ps <0.001 , permutation testing),
161 normalized accuracies were all greater than zero (**Fig. 3f**). Despite only ever
162 being exposed to neural responses from a single object category, this zero-shot
163 decoder was still able to make pair-wise category judgments for neural responses
164 from seven held-out novel categories, displaying extreme generalization not
165 previously reported for any neural decoding model from any imaging modality.
166 Full classification results for all possible number of training categories can be
167 seen in **Supplemental Fig. 4**.

168 Overall, mappings from IT to DCNN activity generalized across object
169 category and novel object categories could be predicted without the model
170 having prior exposure to responses from those categories, providing evidence
171 that DCNNs capture generic visual information in *rhesus macaque* IT, as
172 opposed to information that is restricted to the categories used for fitting.
173 Understanding the neural code for objects requires not only *interpolation* to
174 novel test items similar to those in the training set (as is standard practice), but
175 also *extrapolation* to completely novel shapes that are clearly distinct from those
176 in the training set. Such zero-shot generalization demonstrates that the model
177 has captured the inherent structure of information encoded in neural activity
178 and the relationship between the encoded features and object category. The
179 extreme case of successful generalization from just a single training category
180 (**Figs. 2d,f and 3d,f**) suggests a robust relationship between IT and DCNN
181 representations.

182 In studies linking DCNN features to brain activity, the DCNN units are
183 usually treated equally without regard to the native dimensions in DCNN
184 representations. Here, we separate feature-based and spatial activity, as well as

185 principal component scores summarizing the full spatio-featural activity space,
186 to provide greater clarity regarding which aspects of DCNN representations are
187 explaining the neural variance that carries information about object category.
188 We find that feature-based DCNN channel activity reconstructed from IT
189 activity carries the most generic visual information.

190 This work demonstrates the promise of zero-shot neural decoders from
191 electrophysiological recordings, which has broad applications. Theoretically,
192 zero-shot decoders are superior to standard decoders because they necessitate
193 a computational hypothesis for the neural code underlying the targeted process
194 and support maximal generalization. On the engineering side, they could
195 enable advances such as decoders from chronic neural recordings that can be
196 flexibly updated without new training data. Rather than directly mapping
197 neural responses to every desired output, neural activity would be mapped
198 to a convenient intermediate space that captures the relevant variance for
199 many different outputs. Predicting novel information from the decoder would
200 simply necessitate training a new computational model linking the intermediate
201 features to the new outputs, rather than collecting new neural recordings to
202 learn a direct mapping. As chronic neural recordings become commonplace,
203 building such flexible, generalizable neural decoding systems will become ever
204 more important.

205 Methods

206 Dataset

207 Details about the experimental setup, recording procedure, and pre-processing
208 can be found in [18]. Briefly, two awake *rhesus macaque* monkeys were passively
209 shown a rapid-serial-visual-presentation stream of 2560 grayscale images, each
210 presented 50 times (28 minimum repetitions) for 100 ms, depicting computer-
211 generated objects from eight categories (Airplanes, Animals, Boats, Cars,
212 Chairs, Faces, Fruits Tables) superimposed on arbitrary natural scene back-
213 grounds. Within each of the eight categories, there were eight unique objects
214 (40 images per object). The full stimulus set in [18] had three conditions of
215 objects: low-variation (same size, position, orientation across all background),
216 medium-variation (some variation in size, position, orientation across back-
217 grounds), and high-variation (high variation in size, position, orientation across
218 backgrounds). Here, only images and IT responses from the high-variation
219 condition were used (in which behavioral recognition for a monkey or machine
220 would be most difficult). Neural recordings were acquired from 168 visually-
221 selective IT units using multi-electrode arrays. Firing rates were calculated
222 from 70 to 170 ms post stimulus onset and averaged across repetitions.

223 **Cross-Validation Folds**

224 For all analyses, data were split into category-matched training and testing
225 folds. Data were split into folds according to object, so even within-category the
226 specific objects used in the various training and testing phases were independent.
227 DCNN activity from 75% of the images (six objects per category) was used to
228 map DCNN activity to category labels, define the PCA transform on DCNN
229 activity, and learn the mapping between IT and DCNN activity (Fig. 1).
230 Neural and DCNN activity from 25% of the images (two objects per category)
231 was used to test the models. Twenty unique train-test splits were used for all
232 analyses.

233 **DCNN Architecture**

234 We used VGG-16 trained for 1000-way object categorization on the ImageNet
235 dataset [19]. VGG-16 is a hierarchical DCNN with 21 convolutional, max-
236 pooling, and fully-connected layers (**Fig 1a**, middle column). Our analyses
237 focused on the five pooling layers (pool 1 = 802,816 units, pool2 = 401,408
238 units, pool3 = 200,704 units, pool4 = 100,352 units, pool5 = 25,088 units, fc6
239 = 4096 units, fc7 = 4096 units), which were selected to sample DCNN activity
240 from across the entire hierarchy. The full unit activity space for each layer
241 was reduced to 1000 principle component scores using PCA (75% train, 20%
242 test cross-validation splits). To isolate feature-based channel activity in the
243 pooling layers, we averaged units across the spatial dimensions (pool 1 = 64
244 channels, pool2 = 128 channels, pool3 = 256 channels, pool4 = 512 channels,
245 pool5 = 512 channels). To isolate spatial activity in the pooling layers, we
246 averaged unit activity across channels (pool1 = [112, 112], pool2 = [56, 56],
247 pool3 = [28, 28], pool4 = [14, 14], pool5 = [7, 7]). The fully-connected layers
248 are organized along a single channel dimension, so all 4096 units were included
249 as feature-based activity.

250 **DCNN Readout SVMs**

251 The relationship between DCNN activity and object category labels was learned
252 using linear support-vector-machines (SVMs). DCNN activity from 75% of the
253 images (six objects/category) were used to train the SVMs, and DCNN activity
254 from 25% of the images (two different objects/category) were used to test the
255 SVMs. Prior to training, each unit was normalized to have a mean of zero and
256 a standard deviation of one across images. The same scaling learned on the
257 training set was applied to the test set. Twenty-eight 8-choose-2 binary SVMs
258 were trained, one for every potential pair amongst the eight object categories
259 in the dataset. Binary classification was selected so the two test categories
260 could easily be held-out in the zero-shot condition. Hyper-parameters for
261 each binary classifier were optimized to maximize classification accuracy using

262 three-fold cross-validation within the training set. The same DCNN readout
263 SVMs were used to predict object category from DCNN activity for all analyses.
264 To emphasize, the DCNN readout SVMs were trained independent from any
265 electrophysiological recordings and were never exposed to IT activity until the
266 test phases.

267 **Decoding Methodology**

268 We used partial least squares regression (PLSR) with 25 components (as in [13])
269 to learn the mapping between IT activity and DCNN activity. IT and DCNN
270 activity from 75% of the images (six objects/category) were selected as the
271 training set to learn the IT to DCNN transformation. When IT and DCNN
272 activity from all eight categories were used to learn the transformation, the
273 same transformation was applied to test IT activity from all categories. For
274 zero-shot decoding, one to six categories were used to learn the transformation,
275 and the transformation was applied to the held-out test categories. To assess
276 how well the IT to DCNN transformation extrapolates to novel categories in
277 the most extreme conditions, we used all possible numbers of categories (one to
278 six, step-size one) to learn the IT to DCNN transformation. In all of the above
279 versions, the IT activity transformed into DCNN activity was passed into the
280 DCNN readout SVMs to generate the final object category predictions.

281 **Matching Analysis**

282 To assess the accuracy of the reconstructions and obtain a measure of the
283 overall amount of shared generic information between IT and DCNN activity,
284 we matched DCNN features reconstructed from IT to true DCNN features from
285 the same image. In a pairwise fashion, the reconstructed DCNN features were
286 correlated (Pearson) with the true DCNN features for the same image and the
287 true DCNN features from another image in the test set. If the within-image
288 correlation is greater than the between-images correlation, that comparison
289 was scored as a hit. For a given target image, this comparison was made for
290 every other image in the test set, and the same procedure was applied using
291 each test image as the target image. We averaged across all comparisons to
292 get a matching accuracy (50% chance). Significance was determined using
293 permutation testing. Reconstructed DCNN features were permuted 1000 times
294 relative to their image labels and the full analysis was run to derive a null
295 distribution of matching accuracies. P was defined as the proportion of matching
296 accuracies from this null distribution that are greater than the true matching
297 accuracy. The matching analysis was run for each DCNN feature-type and
298 layer.

299 **Classification Analysis**

300 To assess whether DCNN features reconstructed from IT activity contain
301 information about object category, the reconstructions were fed into the 8-
302 choose-2 linear SVMs trained on true DCNN activity. The SVMs were not
303 modified or fine-tuned on any neural activity. Significance was derived using
304 permutation testing, again permuting the reconstructed DCNN features relative
305 to image labels to derive empirical null distributions of classification accuracies.
306 The classification analysis was run for each DCNN feature-type and layer.

307 **Normalizing Zero-Shot Accuracies**

308 To better assess zero-shot accuracies and compare across conditions with
309 different numbers of training categories, we calculated normalized zero-shot
310 accuracies. The scores account for the proportion of above-chance accuracy
311 present in a zero-shot decoder relative to a control decoder trained on neural
312 responses for all categories. In all comparisons, the number of training images
313 for the all categories control decoder was matched to the number of training
314 images for the zero-shot decoder.

315 **Acknowledgments**

316 This work was supported by the Center for Brains, Minds and Machines, funded
317 by NSF Science and Technology Centers Award CCF-1231216, and NIH Grant
318 R01EY026025. M.M.C is funded by NIMH 108591. We thank Jim DiCarlo
319 for sharing the IT recording data used in this experiment. We thank Tyler
320 Bonnen, Kasper Vincken, Yaoda Xu, and Daeyeol Lee for helpful comments on
321 this work.

322 **Author Contributions**

323 T.P.O’C. conceived of the study. T.P.O’C. and G.K. designed the study.
324 T.P.O’C. performed the analyses. T.P.O’C. and G.K. wrote the manuscript
325 with contributions from M.M.C.

References

1. N. Kriegeskorte and P. K. Douglas, “Interpreting encoding and decoding models,” *Current opinion in neurobiology*, vol. 55, pp. 167–179, 2019.
2. M. Palatucci, D. Pomerleau, G. E. Hinton, and T. M. Mitchell, “Zero-shot learning with semantic output codes,” in *Advances in Neural Information Processing Systems 22* (Y. Bengio, D. Schuurmans, J. D. Lafferty, C. K. I.

- Williams, and A. Culotta, eds.), pp. 1410–1418, Curran Associates, Inc., 2009.
3. T. Horikawa and Y. Kamitani, “Generic decoding of seen and imagined objects using hierarchical visual features,” *Nature Communications*, vol. 8, p. 15037, 2017.
 4. F. Pereira, B. Lou, B. Pritchett, S. Ritter, S. J. Gershman, N. Kanwisher, M. Botvinick, and E. Fedorenko, “Toward a universal decoder of linguistic meaning from brain activation,” *Nature communications*, vol. 9, no. 1, p. 963, 2018.
 5. T. P. O’Connell and M. M. Chun, “Predicting eye movement patterns from fmri responses to natural scenes,” *Nature communications*, vol. 9, no. 1, p. 5159, 2018.
 6. G. K. Anumanchipalli, J. Chartier, and E. F. Chang, “Speech synthesis from neural decoding of spoken sentences,” *Nature*, vol. 568, no. 7753, p. 493, 2019.
 7. N. Kriegeskorte, “Deep neural networks: a new framework for modeling biological vision and brain information processing,” *Annual review of vision science*, vol. 1, pp. 417–446, 2015.
 8. D. L. Yamins and J. J. DiCarlo, “Using goal-driven deep learning models to understand sensory cortex,” *Nature neuroscience*, vol. 19, no. 3, p. 356, 2016.
 9. M. Riesenhuber and T. Poggio, “Hierarchical models of object recognition in cortex,” *Nature Neuroscience*, vol. 2, no. 11, pp. 1019–1025, 1999.
 10. D. Mely and T. Serre, “Towards a theory of computation in the visual cortex,” in *Computational and Cognitive Neuroscience of Vision* (Q. Zhao, ed.), pp. 59–84, Springer, 2016.
 11. D. A. Pospisil, A. Pasupathy, and W. Bair, “‘artiphysiology’ reveals v4-like shape tuning in a deep network trained for image classification,” *eLife*, vol. 7, 2018.
 12. Y. LeCun, Y. Bengio, and G. Hinton, “Deep learning,” *nature*, vol. 521, no. 7553, p. 436, 2015.
 13. D. L. Yamins, H. Hong, C. F. Cadieu, E. A. Solomon, D. Seibert, and J. J. DiCarlo, “Performance-optimized hierarchical models predict neural responses in higher visual cortex,” *Proceedings of the National Academy of Sciences*, vol. 111, no. 23, pp. 8619–8624, 2014.

14. M. Schrimpf, J. Kubilius, H. Hong, N. J. Majaj, R. Rajalingham, E. B. Issa, K. Kar, P. Bashivan, J. Prescott-Roy, K. Schmidt, *et al.*, “Brain-score: Which artificial neural network for object recognition is most brain-like?,” *bioRxiv*, p. 407007, 2018.
15. L. Arend, Y. Han, M. Schrimpf, P. Bashivan, K. Kar, T. Poggio, J. J. DiCarlo, and X. Boix, “Single units in a deep neural network functionally correspond with neurons in the brain: preliminary results,” tech. rep., Center for Brains, Minds and Machines (CBMM), 2018.
16. N. Logothetis and D. Sheinberg, “Visual object recognition,” *Annual Review of Neuroscience*, vol. 19, pp. 577–621, 1996.
17. C. Connor, S. Brincat, and P. A., “Transformation of shape information in the ventral pathway,” *Current Opinion in Neurobiology*, vol. 17, pp. 140–147, 2007.
18. N. J. Majaj, H. Hong, E. A. Solomon, and J. J. DiCarlo, “Simple learned weighted sums of inferior temporal neuronal firing rates accurately predict human core object recognition performance,” *Journal of Neuroscience*, vol. 35, no. 39, pp. 13402–13418, 2015.
19. K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” *arXiv preprint arXiv:1409.1556*, 2014.

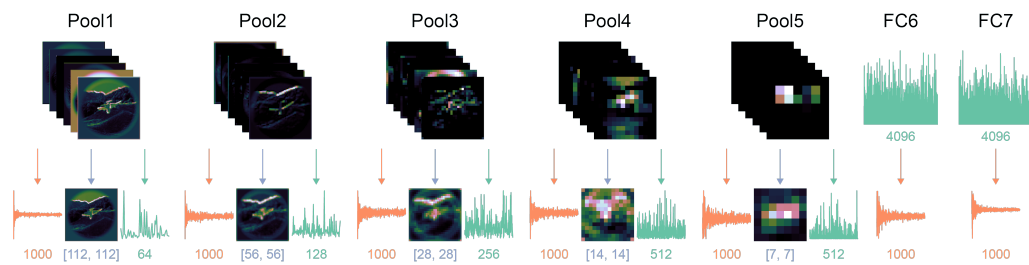


Figure S1. Features extracted from each layer of VGG-16. Components (orange) were defined using principle components analysis and the number of components (1000) was matched across layers. Spatial activity (blue) was defined by averaging across the channel dimension in the native unit-activity space. Feature-based activity (green) was defined by averaging unit activity across the spatial dimensions. The final dimensionality of each feature-type is shown.

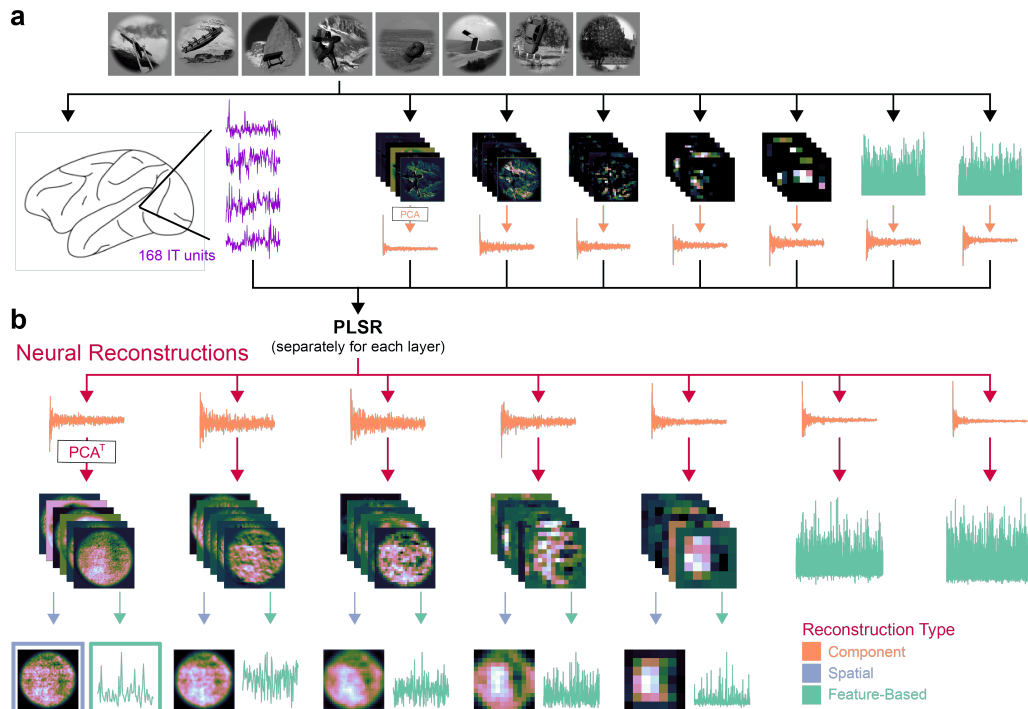


Figure S2. **a.** Two rhesus macaque monkeys viewed images from 8 object categories while IT responses were recorded using multi-electrode arrays. DCNN activity for each image was computed using VGG-16, and the full unit activity for each layer was encoded into 1000 components using PCA. **b.** Using partial least squares regression (PLSR), linear mappings were learned from IT response patterns to DCNN components. These mappings were applied to held-out data (twenty 75% train, 25% test splits) to decode DCNN components from IT activity. The decoded components were multiplied by the transpose of the PCA transformation to reconstruct the full space of DCNN activity for each layer. In the five pooling layers, full reconstructions were averaged across channels (blue arrows) to calculate spatial reconstructions and across spatial dimensions (green arrows) to calculate spatially invariant feature-based reconstructions.

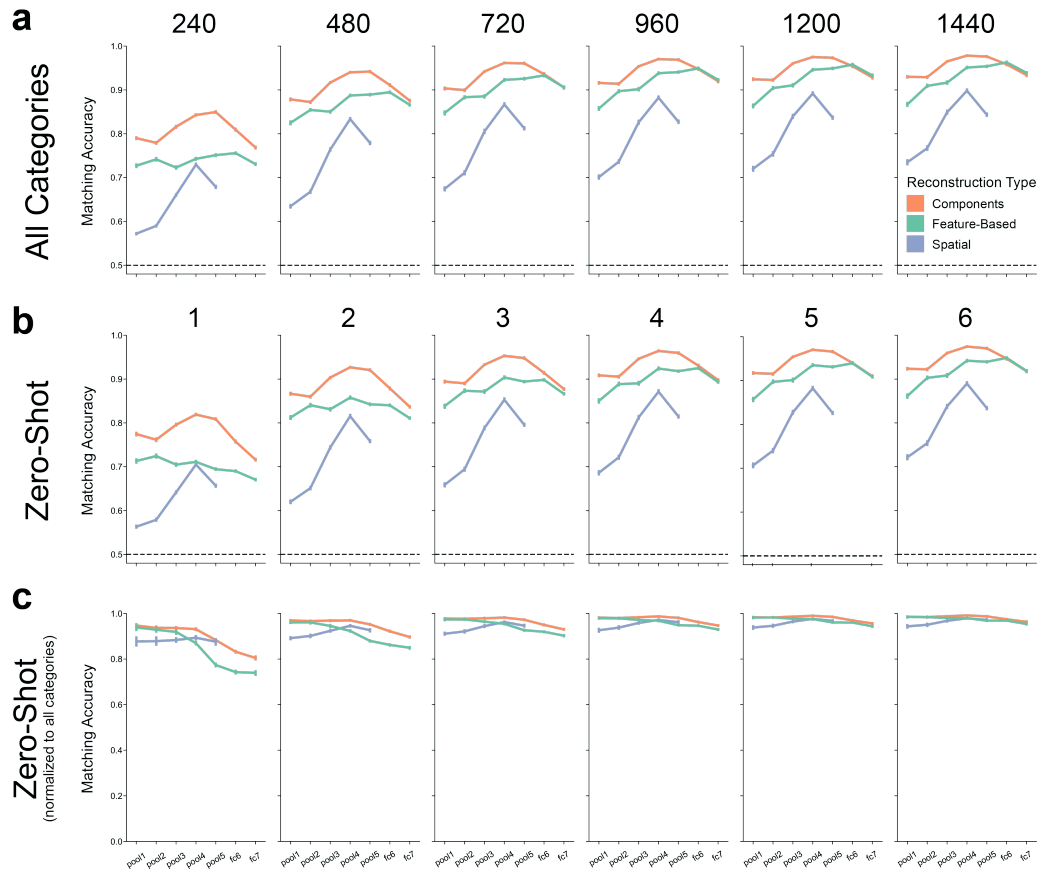


Figure S3. (expanding on **Fig. 2**) **Matching results for all possible number of training categories.** **a.** Model-based matching accuracy when training on neural activity from all categories (similar to **Fig. 2**). The number of training images, shown along the top of each column, are matched to the number of training images for each zero-shot condition below. The subplot for 240 training images corresponds to **Fig. 2e** and the subplot for 1440 training images corresponds to **Fig. 2b**. **b.** Zero-shot matching accuracy for all possible number of training categories. The number of training categories are shown at the top of each column. The subplot for one training category corresponds to **Fig. 2d** and the subplot for six training categories corresponds to **Fig. 2a**. **c.** Zero-shot matching accuracy normalized by the all categories matching accuracy. The subplot for one training category corresponds to **Fig. 2f** and the subplot for six training categories corresponds to **Fig. 2c**.

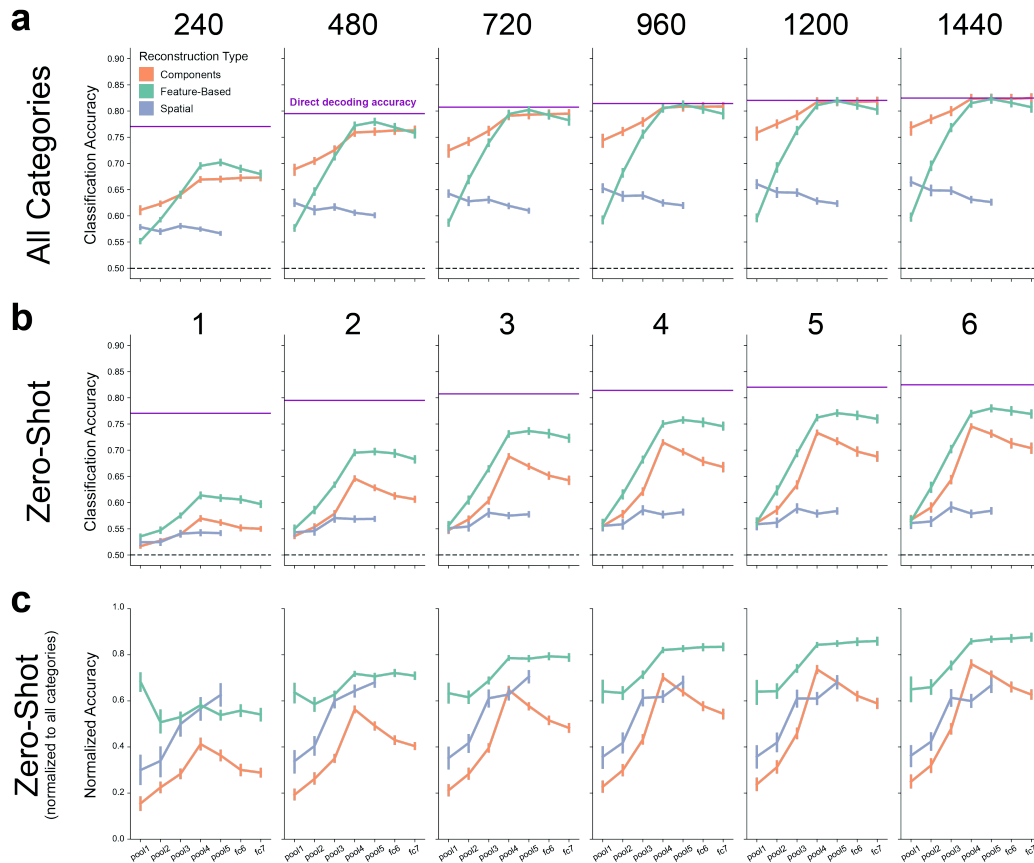


Figure S4. (expanding on **Fig. 3**) **Classification results for all possible number of training categories.** **a.** Model-based classification accuracy training on neural activity from all categories (similar to **Fig. 3**). The number of training images, shown along the top of each column, were matched to the number of training images for each zero-shot condition below. The subplot for 240 training images corresponds to **Fig. 3e** and the subplot for 1440 training images corresponds to **Fig. 3b**. **b.** Zero-shot classification accuracy for all possible number of training categories. The number of training categories are shown at the top of each column. The subplot for one training category corresponds to **Fig. 3d** and the subplot for six training categories corresponds to **Fig. 3a**. **c.** Zero-shot classification accuracy normalized by the all categories classification accuracy. The subplot for one training category corresponds to **Fig. 3f** and the subplot for six training categories corresponds to **Fig. 3c**.