
What can 5.17 billion regression fits tell us about artificial models of the human visual system?

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Rapid simultaneous advances in machine vision and cognitive neuroimaging
2 present an unparalleled opportunity to assess the current state of artificial models
3 of the human visual system. Here, we perform a large-scale benchmarking analysis
4 of 72 modern deep neural network models to characterize with robust statistical
5 power how differences in architecture and training task contribute to the prediction
6 of human fMRI activity across 16 distinct regions of the human visual system. We
7 find: one, that even stark architectural differences (e.g. the absence of convolution
8 in transformers and MLP-mixers) have very little consequence in emergent fits to
9 brain data; two, that differences in task have clear effects—with categorization and
10 self-supervised models showing relatively stronger brain predictivity across the
11 board; three, that feature reweighting leads to substantial improvements in brain
12 predictivity, without overfitting – yielding model-to-brain regression weights that
13 generalize at the same level of predictivity to brain responses over 1000s of new
14 images. Broadly, this work presents a lay-of-the-land for the emergent correspon-
15 dences between the feature spaces of modern deep neural network models and the
16 representational structure inherent to the human visual system.

17 1 Introduction

18 The pace of progress in computer vision poses a practical challenge for neuroscientists seeking to
19 assess state-of-the-art models in their ability to explain visual representation and behavior. New high-
20 performing models are released on a near-daily basis, and recent innovations (e.g. in self-supervised
21 learning [1, 2]) have created myriad new opportunities for productive synergy between the fields of
22 biological and machine vision. As such, methods for comparing the brain-predictivity of artificial
23 models using a predefined analysis pipeline (“neural benchmarking”) are critical in helping discern
24 the algorithmic innovations that may be meaningful with respect to the study of brain function.

25 Existing public neural benchmarking datasets have been limited to mouse and primate neurophysi-
26 ology (3, 4). Recent advances in the scale and quality of human neuroimaging datasets (5–7) now
27 present an opportunity to rigorously assess the state of deep neural network modeling as applied to
28 the human visual system.

29 Here we present a large-scale benchmark of dozens of state-of-the-art deep neural network models
30 in their prediction of human brain activity across the visual hierarchy. Aiming for coverage, our
31 survey attempts to document the current trends in how well different kinds of models, varying in both
32 task and architecture, learn features with brain-like response signatures. Our results complement
33 prior work examining different model predictivities (8–10), but at a significantly larger scale, and

34 incorporating a set of more modern models not yet fully accounted for in the benchmarking literature
35 (e.g. self-supervised models and vision transformers).

36 **2 Methods**

37 As the target of our neural benchmark, we use the Natural Scenes Dataset (NSD, [5]), a recent fMRI
38 dataset representing the most extensive sampling of visual responses in individual participants to
39 date (30,000 stimuli viewed per subject; 73,000 unique images total). Here we analyze only a small
40 fraction of this dataset, focusing on responses to 1,000 COCO stimuli that were shown to 4 subjects
41 at least 3 times, in a subset of ROIs along the visual hierarchy. We compare these responses with the
42 responses of 72 modern DNNs that vary in task and architecture (see Appendix for details).

43 We employ two methods for mapping the activations of model features within a layer to regions
44 of the brain – classical representational similarity analysis (RSA, [11]) and voxelwise-encoding
45 (re-weighted) RSA (12). Classical RSA considers all of the features from a given model layer equally
46 in computing the image-wise representational dissimilarity matrix (RDM), which is directly compared
47 with a given neural RDM. This method requires a fully-emergent match in population-level geometry
48 between a neural ROI and the full set of units in a model layer.

49 Voxelwise encoding RSA (veRSA), on the other hand, takes advantage of feature reweighting to
50 identify different model subspaces that correspond to the variance in different brain regions (13). To
51 implement voxel-wise encoding RSA, we use an efficient high-throughput model-fitting procedure,
52 first applying leave-one-out cross-validated ridge regression to map between a given model feature
53 space and the observed univariate activity pattern of each voxel; once we’ve collected a set of
54 predictions for the patterns of activity for each voxel in a given ROI, we compute an RDM from
55 these predictions and compare that RDM to the RDM in the brain. This re-weighted RSA procedure
56 requires massive parallelization, and entails performing a total of around 5.17 billion regression
57 fits (calculated by multiplying the total number of model layers we analyze by the total number of
58 voxels under consideration from the brain dataset). To assuage concerns of overfitting, we validate
59 the robustness of our fitted regressions by testing their generalizability to 1000 independent images
60 entirely removed from the training procedure.

61 **3 Results**

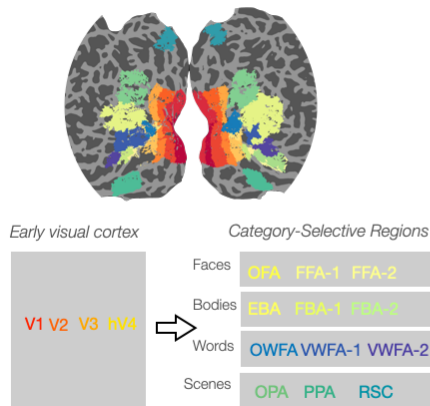
62 **3.1 Hierarchical Correspondence**

63 As a first step and sanity check, we ask: Does the seminal finding that the information processing
64 hierarchies in deep nets recapitulate the information processing hierarchy in the human visual system
65 (14–16) hold at scale and across a significantly diverse population of models? The answer is a
66 resounding affirmative (**Figure 1**): Using a purely data-driven aggregation procedure, we show that
67 the relative depth of the best-fitting model layer for each ROI seems to re-capitulate the human
68 visual hierarchy (e.g. early visual areas, followed by category-selective regions). This hierarchical
69 convergence holds even when breaking down the models by broad, divergent classes of architecture.

70 **3.2 Architecture Variation**

71 How do models with different architectures compare in their ability to predict the structure of human
72 brain responses across the visual system? Our particular survey of models, chosen deliberately
73 to reflect the diversity of modern object recognition (ImageNet-trained) architectures, allows for
74 numerous subdivisions, but perhaps the most prominent is between convolutional architectures (e.g.
75 VGG, ResNet, MobileNet, $n = 24$), vision transformers (e.g. Visformer, DeIT, $n = 13$) and MLP-
76 mixers (e.g. ResMLP, gMixer, $n = 5$). The latter two of these are more recent advents of computer
77 vision, and are defined by the lack of a convolutional inductive bias – once considered a cornerstone
78 of the link between biological and machine vision. Comparisons between these architectures (across
79 both classical RSA and voxel-encoding RSA) are shown in Figure 2.

A. Regions-of-Interest



B. Hierarchical Correspondence

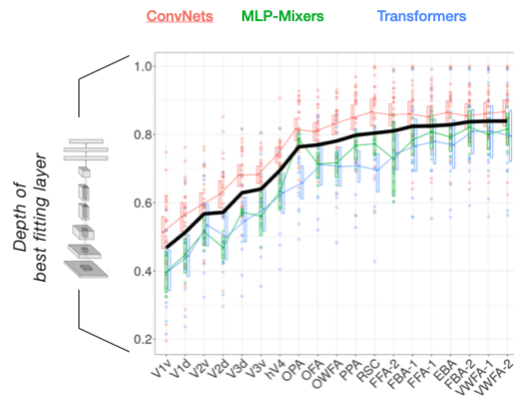


Figure 1: (A) Visualization of selected regions-of-interest on a flattened hemisphere. (B) Emergent hierarchical correspondence between the most predictive model layer and the hypothesized information processing hierarchy of the visual system. Regions along the x-axis are ordered by the average depth of the best predicting layer (across all models). Data are also broken down by the architectural distinctions of ConvNets, MLP-Mixers, and Transformers. Each point is the best performing layer from a given model, averaged over subjects.

80 To test for differences in predictivity, we use nonparametric ANOVAs. Without reweighting (classical RSA), there is a significant difference across ConvNets, MLP-Mixers, and Transformers (81 $\chi^2_{\text{Kruskal-Wallis}}(2) = 10.14; p_{\text{Holm}} = 0.02; \eta^2_{\text{ordinal}} = 0.27; CI_{95\%}[0.09; 0.49]$) in early visual areas, driven by a significant pairwise advantage of ConvNets over Transformers. With reweighting (82 veRSA), this difference disappears. Without reweighting, there is no significant difference (83 between architectures in higher-level cortical areas. With reweighting, there is a difference (84 $\chi^2_{\text{Kruskal-Wallis}}(2) = 10.59; p_{\text{Holm}} = 0.02; \eta^2_{\text{ordinal}} = 0.26; CI_{95\%}[0.09; 0.56]$), driven this time by (85 the pairwise superiority of both ConvNets and MLP-Mixers over Transformers. (86 (87

88 Behind these apparently significant effects is the numerical reality that the raw effect sizes in (89 both cases is effectively negligible – less than $r_{\text{Pearson}} = 0.01$ and $r_{\text{Pearson}} = 0.02$, respectively. (90 As such, the most striking effect here is not that of architecture, but of mapping method, which (91 substantially augments the predictive power of every model in our survey (with average gains of (92 $r_{\text{Pearson}} = 0.160; CI_{95\%}[0.152; 0.166]$) across model and ROI). In the most notable case, models (93 in EBA experience average gains of $r_{\text{Pearson}} = 0.265$. These improvements dwarf any difference (94 attributable to architecture, and underscore an important point: despite dramatic differences in the (95 design and algorithmic inductive biases of ConvNets, MLP-Mixers, and Transformers, there is little (96 consequence on the resulting brain predictivity (regardless of mapping method).

97 3.3 Task Variation

98 How does brain predictivity vary as a function of task? For a window into this question we consider (99 the 24 models from the Taskonomy project (17). These models share the same base architecture (100 ResNet-50) and visual diet, but are trained on 1 of 24 popular computer vision tasks. These tasks (101 are organized into 4 different categories (2D, 3D, Semantic and Geometric) according to what the (102 authors of the Taskonomy project call the models’ ‘transfer affinity’ – the degree to which a model (103 trained on one task supports transfer learning to another. The prediction levels of these models for (104 both classical and voxel-wise encoding RSA are shown in Figure 3.

105 Without reweighting, there is considerable variability across ROI in the tasks that are most predictive (106 of the brain, but the differences between the best task and the second-best task is minimal in most cases. (107 In V2, for example, a 2D task (edge detection) is the most predictive of the tasks at $r_{\text{Pearson}} = 0.226$, (108 but is closely followed by a 3D task (Keypoints) at $r_{\text{Pearson}} = 0.220$.

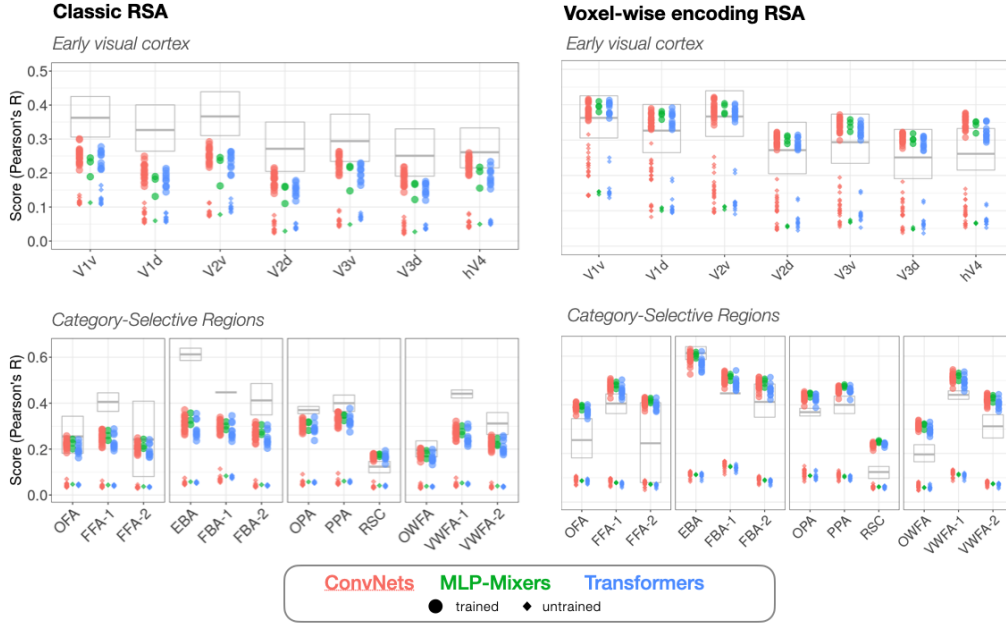


Figure 2: *Architecture variations. Model fits are shown along the y-axis, for early visual areas (top row) and category-selective areas (bottom row), for classical RSA (left) and voxel-wise encoding RSA (right). The gray boxes indicate an intersubject reference point (the average pairwise correlation of individual subject RDMs). Each dot is the best performing layer from a single model, with trained models in large circular points, and untrained counterparts in small diamonds.*

109 As in the case of architecture, feature reweighting (veRSA) leads to uniform improvement across
 110 models. Strikingly, however, object and scene classification gain disproportionately. The gains for
 111 object recognition are so substantial that it becomes the single most predictive task for all brain areas,
 112 often dominating by an impressively large margin, with a mean gain over the next best task (apart
 113 from scene classification) of $r_{\text{Pearson}} = 0.127; CI_{95\%}[0.122; 0.131]$ across all ROIs.

114 While these results point strongly to an advantage of category supervision in the formation of neurally
 115 predictive representation (at least in the case of veRSA), the self-supervised models (absent from
 116 Taskonomy) in our survey allow us to delve more deeply into whether the classification objective is
 117 the key driver of neural predictivity, or whether category-supervised models derive their advantage
 118 from the set of invariances that they learn in service of classification.

119 The predictive power of our self-supervised models strongly suggest the latter: regardless of mapping
 120 method, self-supervised models (especially recent contrastive ResNet-50 models such as SimCLR
 121 and BarlowTwins) tend to show a small but statistically significant advantage over a (recently
 122 revamped) category-supervised ResNet-50 [18]. For example, averaging across brain ROIs, SimCLR
 123 eeks out a mean gain of $r_{\text{Pearson}} = 0.013; CI_{95\%}[0.0106; 0.0192]$ in weighted RSA and a gain of
 124 $r_{\text{Pearson}} = 0.006; CI_{95\%}[0.002; 0.008]$ in classical RSA. While these results should not be interpreted
 125 as indicating superiority of self-supervision over category-supervision, they do indicate *parity* in
 126 prediction levels – a win for ethological plausibility (12, 19).

127 3.4 Generalization Tests

128 The sheer quantity of regression fits required to summarize the predictive performance of our model
 129 set, and the vast number of dimensions relative to data points, may raise concern: is this deep
 130 encoding pipeline massively overfitting, in spite of our cross-validation procedures? Or, are the
 131 estimates we derive truly a reasonable approximation (given the linking assumptions inherent in the
 132 analysis) of a given model’s brain predictivity?

A. Taskonomy Models

2-D		3-D		Geometric	Semantic	Random
Autoencoder	2D Segmentation	Occlusion Edges	2.5D Segmentation	Camera Pose	Object Classification	Random Weights
Denoising	Curvatures	3D Keypoints	Egomotion	Point Matching	Scene Classification	
Texture Edges	Euclidean Depth	Surface Normals	Camera Pose	Room Layout	Semantic Segmentation	
Inpainting	Z-Buffer Depth	Reshading	Jigsaw	Vanishing Point		
2D Keypoints						

B. Classic RSA

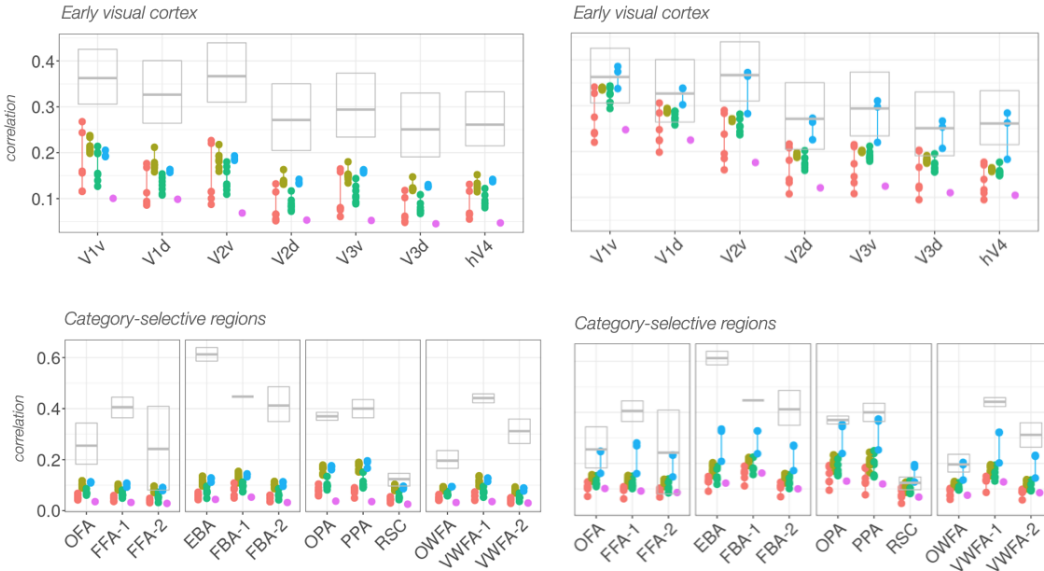


Figure 3: Effect of task on model-brain predictivity. (A) 24 Taskonomy models with a ResNet-50 architecture, grouped into 4 categories (2D, 3D, Geometric, or Semantic). An untrained (randomly-initialized) model, for comparison. The correlation between the model features and brain responses for early visual areas (top) and category-selective regions (bottom) is plotted, using classical RSA in (B) and reweighted RSA in (C). Gray box plots indicate the range of inter-subject RDM correlations using classical RSA.

133 To address this concern, we conducted a separate generalization test for each model, in which we
 134 selected the best performing layer (according to the original LOOCV score from our regression
 135 procedure) per subject, per ROI. For these layers, we then use reweighted RSA to compare brain and
 136 model feature spaces using a set of 1000 entirely held-out test images per subject. These images were
 137 never referenced or incorporated during training, and prediction scores on these images thus provide
 138 a measure of "pure" generalization.

139 Even with this more stringent test, we found little-to-no drop in accuracy in predicting brain represen-
 140 tation evoked by the 1000 unique test images per fMRI subject. When aggregating across subjects,
 141 models, and ROI, for example, the mean decrease in score on the unseen images was less than 1%
 142 ($r_{\text{Pearson}} = 0.0095$; $C I_{95\%} [0.00422; 0.0153]$). By adding a mere 103 million regression fits to our
 143 initial total of 5.17 billion, then, we can thus confirm definitively that our encoding models generalize
 144 to previously-unseen data. (A more detailed figure showing generalization across specific subjects
 145 and ROIs is shown in the Appendix.)

146 4 Discussion

147 So what can we learn about the human visual system from 5.17 billion regression fits? Broadly, it
 148 seems, there are two sets of answers, one more pessimistic, one more optimistic. On the side of
 149 pessimism, the lack of variation across architecture suggests that massive innovations in computer
 150 vision may often yield little to no change in our ability to predict the representational structure of
 151 biological vision, disrupting what was once prophesied to become a glorious feedback loop between

152 neuroscientific insight and engineering innovation. What’s more, the frequent variability in interpre-
153 tation across mapping method seems a potential pitfall if not accounted for with greater vigilance
154 and attention to theoretical commitment. On the side of optimism, it appears that more general,
155 algorithmic correspondences between DNNs and brains (especially in terms of the information
156 processing hierarchy) persist in spite of an increasingly rapid shift away from biological plausibility
157 in engineering. In opposite direction of this shift is a promising move *towards* ethological plausibility
158 – many cutting-edge models no longer rely on learning targets humans almost certainly do not share
159 (e.g. full category supervision). Not coincidentally, these models appear to be competitive predictors
160 of brain activity.

161 Current models are still far from capturing the kaleidoscopic complexity of biological visual systems.
162 Our goal in pursuing this large-scale benchmark is not to discern the "best" model of vision, but
163 rather to clarify what kinds of things are and are not important for building next-generation perceptual
164 models that will push our understanding of human vision further.

165 References

- 166 [1] Jure Zbontar, Li Jing, Ishan Misra, Yann LeCun, and Stéphane Deny. Barlow twins: Self-
167 supervised learning via redundancy reduction. *arXiv preprint arXiv:2103.03230*, 2021.
- 168 [2] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre H Richemond, Elena
169 Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Daniel Guo, Mohammad Gheshlaghi
170 Azar, et al. Bootstrap your own latent: A new approach to self-supervised learning. *arXiv*
171 *preprint arXiv:2006.07733*, 2020.
- 172 [3] Martin Schrimpf, Jonas Kubilius, Ha Hong, Najib J. Majaj, Rishi Rajalingham, Elias B. Issa,
173 Kohitij Kar, Pouya Bashivan, Jonathan Prescott-Roy, Franziska Geiger, Kailyn Schmidt, Daniel
174 L. K. Yamins, and James J. DiCarlo. Brain-score: Which artificial neural network for object
175 recognition is most brain-like? *bioRxiv preprint*, 2018.
- 176 [4] Colin Conwell, David Mayo, Michael A. Buice, Boris Katz, George A. Alvarez, and Andrei
177 Barbu. Neural regression, representational similarity, model zoology & neural taskonomy
178 at scale in rodent visual cortex. *bioRxiv*, 2021. doi: 10.1101/2021.06.18.448431. URL
179 <https://www.biorxiv.org/content/early/2021/06/18/2021.06.18.448431>.
- 180 [5] Emily Jean Allen, Ghislain St-Yves, Yihan Wu, Jesse L Breedlove, Logan T Dowdle, Bradley
181 Caron, Franco Pestilli, Ian Charest, J Benjamin Hutchinson, Thomas Naselaris, et al. A massive
182 7t fmri dataset to bridge cognitive and computational neuroscience. *bioRxiv*, 2021.
- 183 [6] Nadine Chang, John A Pyles, Austin Marcus, Abhinav Gupta, Michael J Tarr, and Elissa M
184 Aminoff. Bold5000, a public fmri dataset while viewing 5000 visual images. *Scientific data*, 6
185 (1):1–18, 2019.
- 186 [7] Oliver Contier, Martin N Hebart, Adam H Dickter, Lina Teichmann, Alexis Kidder, Anna
187 Coriveau, Charles Zheng, Maryam Vaziri-Pashkam, and Charles Baker. Things-fmri/meg: A
188 large-scale multimodal neuroimaging dataset of responses to natural object images. *Journal of*
189 *Vision*, 21(9):2633–2633, 2021.
- 190 [8] Seyed-Mahdi Khaligh-Razavi and Nikolaus Kriegeskorte. Deep supervised, but not unsuper-
191 vised, models may explain it cortical representation. *PLoS computational biology*, 10(11),
192 2014.
- 193 [9] Katherine R Storrs, Tim C Kietzmann, Alexander Walther, Johannes Mehrer, and Nikolaus
194 Kriegeskorte. Diverse deep neural networks all predict human inferior temporal cortex well,
195 after training and fitting. *Journal of Cognitive Neuroscience*, 33(10):2044–2064, 2021.
- 196 [10] Aria Wang, Michael Tarr, and Leila Wehbe. Neural taskonomy: Inferring the similarity of
197 task-derived representations from brain activity. In *Advances in Neural Information Processing*
198 *Systems*, pages 15475–15485, 2019.

- 199 [11] Nikolaus Kriegeskorte, Marieke Mur, and Peter A Bandettini. Representational similarity
200 analysis-connecting the branches of systems neuroscience. *Frontiers in systems neuroscience*
201 2:4, 2008.
- 202 [12] Talia Konkle and George A Alvarez. Beyond category-supervision: instance-level contrastive
203 learning models predict human visual system responses to objects. *bioRxiv*, 2021.
- 204 [13] Philipp Kaniuth and Martin N Hebart. Feature-reweighted rsa: A method for improving the fit
205 between computational models, brains, and behavior. *bioRxiv*, 2021.
- 206 [14] Daniel LK Yamins, Ha Hong, Charles F Cadieu, Ethan A Solomon, Darren Seibert, and James J
207 DiCarlo. Performance-optimized hierarchical models predict neural responses in higher visual
208 cortex. *Proceedings of the National Academy of Sciences* 111(23):8619–8624, 2014.
- 209 [15] Umut Güçlü and Marcel AJ van Gerven. Deep neural networks reveal a gradient in the
210 complexity of neural representations across the ventral stream. *Journal of Neuroscience* 35
211 (27):10005–10014, 2015.
- 212 [16] Umut Güçlü and Marcel AJ van Gerven. Increasingly complex representations of natural movies
213 across the dorsal stream are shared between subjects. *NeuroImage* 145:329–336, 2017.
- 214 [17] Amir R Zamir, Alexander Sax, William Shen, Leonidas J Guibas, Jitendra Malik, and Silvio
215 Savarese. Taskonomy: Disentangling task transfer learning. *Proceedings of the IEEE
216 Conference on Computer Vision and Pattern Recognition* pages 3712–3722, 2018.
- 217 [18] Ross Wightman, Hugo Touvron, and Hervé Jégou. Resnet strikes back: An improved training
218 procedure in timmarXiv preprint arXiv:2110.00476, 2021.
- 219 [19] Chengxu Zhuang, Siming Yan, Aran Nayebi, Martin Schrimpf, Michael C Frank, James J
220 DiCarlo, and Daniel LK Yamins. Unsupervised neural network models of the ventral visual
221 stream. *Proceedings of the National Academy of Sciences* 118(3), 2021.
- 222 [20] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr
223 Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. *European
224 conference on computer vision* pages 740–755. Springer, 2014.
- 225 [21] Jacob S Prince, John A Pyles, Michael J Tarr, and Kendrick N Kay. Glimsingle: a turnkey
226 solution for accurate single-trial fmri response estimation. *Journal of Vision* 21(9):2831–2831,
227 2021.
- 228 [22] Leyla Tarhan and Talia Konkle. Reliability-based voxel selection. *NeuroImage* 207:116350,
229 2020.
- 230 [23] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan,
231 Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas
232 Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy,
233 Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-
234 performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, A. Ché-
235 Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*
236 pages 8024–8035. Curran Associates, Inc., 2019. [http://papers.neurips.cc/paper/
237 9015-pytorch-an-imperative-style-high-performance-deep-learning-library.
238 pdf](http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf).
- 239 [24] Ross Wightman. Pytorch image models. [https://github.com/rwightman/
240 pytorch-image-models](https://github.com/rwightman/pytorch-image-models) , 2019.
- 241 [25] Priya Goyal, Quentin Duval, Jeremy Reizenstein, Matthew Leavitt, Min Xu, Benjamin
242 Lefaudeux, Mannat Singh, Vinicius Reis, Mathilde Caron, Piotr Bojanowski, Armand Joulin,
243 and Ishan Misra. Vissl <https://github.com/facebookresearch/vissl> , 2021.

- 244 [26] Alexander Sax, Bradley Emi, Amir R. Zamir, Leonidas J. Guibas, Silvio Savarese, and Jitendra
 245 Malik. Mid-level visual representations improve generalization and sample efficiency for
 246 learning visuomotor policies. 2018.
- 247 [27] Alexander Sax, Jeffrey O Zhang, Bradley Emi, Amir Zamir, Silvio Savarese, Leonidas Guibas,
 248 and Jitendra Malik. Learning to navigate using mid-level visual priors. arXiv preprint
 249 arXiv:1912.11121, 2019.
- 250 [28] Ping Li, Trevor J Hastie, and Kenneth W Church. Very sparse random projections. Proce-
 251 dings of the 12th ACM SIGKDD international conference on Knowledge discovery and data
 252 mining, pages 287–296, 2006.
- 253 [29] Nikolaus Kriegeskorte, Marieke Mur, Douglas A Ruff, Roozbeh Kiani, Jerzy Bodurka, Hossein
 254 Esteky, Keiji Tanaka, and Peter A Bandettini. Matching categorical object representations in
 255 inferior temporal cortex of man and monkey. *Neuron*, 60(6):1126–1141, 2008.
- 256 [30] Aran Nayebi, Nathan CL Kong, Chengxu Zhuang, Justin L Gardner, Anthony M Norcia, and
 257 Daniel LK Yamins. Unsupervised models of mouse visual cortex. arXiv, 2021.

258 Checklist

- 259 1. For all authors...
- 260 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
 261 contributions and scope? [\[Yes\]](#)
- 262 (b) Did you describe the limitations of your work? [\[Yes\]](#) See discussion section.
- 263 (c) Did you discuss any potential negative societal impacts of your work? [\[N/A\]](#)
- 264 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
 265 them? [\[Yes\]](#)
- 266 2. If you are including theoretical results...
- 267 (a) Did you state the full set of assumptions of all theoretical results? [\[N/A\]](#)
- 268 (b) Did you include complete proofs of all theoretical results? [\[N/A\]](#)
- 269 3. If you ran experiments...
- 270 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
 271 mental results (either in the supplemental material or as a URL)? [\[No\]](#) The code and
 272 the data will be released publicly upon publication.
- 273 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
 274 were chosen)? [\[Yes\]](#) See Methods section and the Appendix.
- 275 (c) Did you report error bars (e.g., with respect to the random seed after running exper-
 276 iments multiple times)? [\[Yes\]](#) See figures, figure captions, Methods section and the
 277 Appendix.
- 278 (d) Did you include the total amount of compute and the type of resources used (e.g.,
 279 type of GPUs, internal cluster, or cloud provider)? [\[Yes\]](#) See the Appendix section
 280 describing compute required
- 281 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 282 (a) If your work uses existing assets, did you cite the creators? [\[Yes\]](#)
- 283 (b) Did you mention the license of the assets? [\[N/A\]](#)
- 284 (c) Did you include any new assets either in the supplemental material or as a URL? [\[No\]](#)
- 285 (d) Did you discuss whether and how consent was obtained from people whose data you're
 286 using/curating? [\[N/A\]](#) All details relating to the consent processes for human subjects
 287 included in the Natural Scenes Dataset can be found in 5

- 288 (e) Did you discuss whether the data you are using/curating contains personally identifiable
289 information or offensive content? [N/A]
- 290 5. If you used crowdsourcing or conducted research with human subjects...
- 291 (a) Did you include the full text of instructions given to participants and screenshots, if
292 applicable? [N/A]
- 293 (b) Did you describe any potential participant risks, with links to Institutional Review
294 Board (IRB) approvals, if applicable? [N/A]
- 295 (c) Did you include the estimated hourly wage paid to participants and the total amount
296 spent on participant compensation? [N/A]

297 A Appendix

298 A.1 Human Brain Data

299 The Natural Scenes Dataset (NSD) is the largest effort to date to measure human brain responses with
300 functional magnetic resonance imaging (fMRI), reflecting measurements of 73,000 unique stimuli
301 from the Microsoft Common Objects in Context (COCO) dataset [20] at high resolution (7T field
302 strength, 1.33s TR, 1.8mm³ voxel size). In the present work, we analyze only a small fraction
303 of this dataset, focusing on responses to images that enable direct comparison between data from
304 different subjects. That is, we focus on the 1000 COCO stimuli that overlapped between subjects (the
305 "shared1000" images), and limit analyses to the 4 subjects (subjs 01, 02, 05, 07) for whom all 3 image
306 repetitions are available for the shared1000. For the generalization tests, we also select a random
307 unique set of 1000 images for each subject; these were not included in the "shared1000." All responses
308 were estimated using a custom GLM toolbox ("GLMsingle" [21]), which was applied during the
309 preprocessing of NSD time-series data, featuring optimized denoising and regularization procedures,
310 to accurately measure changes in neural activity in response to each experimental stimulus.

311 We focus our analyses to voxels within a set of predefined functional ROIs that span the visual
312 hierarchy (see [5] for details on the procedures used to define the ROIs). Further, to maximize
313 SNR of the target data, we implement a reliability-based voxel selection procedure that isolates
314 regions of the brain containing stable structure in their responses. To compute the split-half reli-
315 ability at a given voxel, we use 1,000 images from each subject (independent from the shared1000,
316 and from all images included in our main analyses and generalization tests), and take the aver-
317 age correlation in univariate response profiles over each pair of available image repetitions (e.g.
318 $\text{mean}(r(\text{rep1}; \text{rep2}); r(\text{rep2}; \text{rep3}); r(\text{rep1}; \text{rep3}))$). ROI voxels exceeding a reliability threshold of
319 $\text{Pearsonr} = 0.1$ were included in subsequent analyses. These procedures yield a matrix of dimension
320 (images, voxels, repetitions) for each subject's ROI, and we average over the 3 repetitions to yield the
321 final ROI data input into our neural benchmarking pipeline.

322 A.2 Candidate Deep Neural Network Models

323 In total, we survey a set of 72 distinct models (110 including the randomly-initialized versions of
324 certain of these models). These models are sourced from four different repositories: the Torchvision
325 (PyTorch) model zoo [23]; the pytorch-image-models (timm) library [24]; the VISSL (self-supervised)
326 model zoo [25]; and the Taskonomy (visualpriors) project [26, 27]. The first two of these
327 repositories offer pretrained versions of a large number of object recognition models with varying
328 architectures: including (classic and modern) convolutional networks, vision transformers, and MLP-
329 mixers. For each of these 'ImageNet' (object recognition) models, we include one trained and one
330 randomly initialized variant (using the initialization scheme the model authors recommended) so as to
331 assess the impact of ImageNet training on brain prediction, and as a sanity check. The self-supervised
332 models are mainly variants on a popular convolutional architecture (ResNet-50), though do include
333 some transformers (the 'DINO' models). The Taskonomy models consist of a core encoder-decoder
334 architecture trained on 24 different common computer vision tasks, ranging from autoencoding to
335 edge detection. These models are engineered in such a way that only the architecture of the decoder

336 varies across task, allowing us to assess (after detaching the encoder) what effect different kinds of
337 training has on predictive power, independent of model design.

338 A.3 Benchmarking Pipeline

339 Feature Extraction For each of our deep neural network models, we extract features in response to
340 each of our probe stimuli at each distinct layer of the network. Importantly, we define a layer here
341 as a distinct computational (sub)operation. This means, for example, that we treat convolution and
342 the rectified nonlinearity that follows it as two distinct feature maps. This is especially relevant in
343 the case of transformers, where the features inherent to the key - query - value computation of the
344 attention maps often differ substantially. At the end of our feature extraction procedure, we have for
345 each model and each model layer, a matrix of features of the dimension: number of images x number
346 of attended features from a given layer

347 Classical RSA (cRSA) As a first method of mapping deep neural network responses to voxel
348 responses, we use classical RSA, a nonparametric mapping method that quantifies the extent
349 similarity of the 'representational geometry' between two feature spaces, regardless of origin. To
350 compute this metric, we construct representational dissimilarity matrices (RDMs) using the pairwise
351 correlation distance ($1 - \text{Pearson}$) between the responses of a given neural ROI (image by voxel) or
352 a given model layer (image by unit) for all images being considered. We then compare these RDMs
353 by taking a second-order correlation (Pearson^2) between the attended upper-triangular portion
354 of each. This ultimately yields a matrix of correlation scores of dimension: (number of subjects x
355 number of ROIs x number of model layers x number of models). Classical RSA reflects the extent to
356 which the representational structure in each model layer naturally recapitulates the representational
357 structure in a visual cortical ROI, without alteration or feature reweighting.

358 Voxelwise Encoding RSA The following procedure yields the billions of regression fits we reference
359 in the title. The pipeline works as follows: first, we fit a regression for each voxel as a weighted
360 combination of model layer features. Given that the number of features in a layer sometimes number
361 in the millions, we employ sparse random projection [26] as a dimensionality-reduction procedure,
362 and then use ridge regression as a linear model to relate the model feature space to each voxel's tuning
363 function. Then, we use the voxel-encoding models to generate predicted activation profiles to the
364 complete set of held-out images, and correlate the subsequent predicted representational similarity
365 structure to that of the brain. For additional detail, see Section A.4.

366 We emphasize that this method contrasts with popular practices in primate and mouse benchmarking,
367 which treat predictivity of unit-level univariate response profiles as the key measure. However, fMRI
368 affords more systematic spatial sampling over the cortex. Thus, for the present analysis, rather
369 than taking the aggregate of single voxel fits as our key measure, we choose to treat the population
370 representational geometry over each ROI as our critical target for prediction. This multi-voxel
371 similarity structure provides different kinds of information about the format of population-level
372 coding than do individual units. (29).

373 Noise Ceilings and Reference Metrics

374 While powerful in the quantity and diversity of its images, the number of repetitions in image
375 presentation (3 per image) in the NSD dataset leaves little room to estimate a noise ceiling per voxel
376 with standard split-half reliability methods. Thus, as a reference metric for how well our models
377 are doing overall, we use inter-subject predictivity: a measure of how well the brain of one human
378 predicts the brain of another. Here, we took the average pairwise correlation of the individual subject
379 RDMs in a given ROI. For a more in-depth discussion of ceilings and reference points, as well as
380 experimental alternatives, see Section A.5.

381 A.4 Voxelwise Encoding RSA In-Depth

382 To predict the activity profile of each voxel, we first use Sparse Random Projection (SRP) [26]
383 project the model features generated in response to our 1000 probe images into a lower-dimensional

384 space. We use a dimensionality of 5960 projections—a number we choose using the Johnson-
385 Lindenstrauss lemma, which mathematically guarantees the preservation of pairwise distances in a
386 given space of operations (with a minimal distortion defined by a hyperparameter epsilon, which we
387 leave in all cases at the scikit-learn default of 0.1). We then perform a leave-one-out cross-validated
388 (LOOCV) ridge regression (cross-validating over images) to map these projections to the responses
389 of our voxels, obtaining a vector of predicted voxel responses that we then correlate with the true
390 voxel responses to obtain a score per voxel per model layer.

391 This leave-one-out cross-validation is performed in a single matrix operation often referred to as
392 generalized cross-validation, and is numerically equivalent to iterative leave-one-out, but is effectively
393 instantaneous. We iterate this regression procedure until we have a score for all voxels and all model
394 layers. No hyperparameter selection was performed over the course of the benchmarking, apart from
395 a minimal, exploratory grid search for a lambda parameter (of values $1e2, 1e3, 1e4, 1e5, 1e6, 1e7$)
396 on an AlexNet model (which we subsequently excluded from the main analysis). Thus, all feature
397 spaces were projected to 5960 sparse random projections, and all regressions were run with a lambda
398 penalty of $1e5$.

399 Rather than taking single-voxel scores as our key measure, we consider the geometry of the population
400 across the larger region of interest as a critical target for prediction. To do so, we use the predicted
401 responses from our voxel-wise encoding method to generate predicted representational dissimilarity
402 matrices. The logic behind this procedure is effectively to dispense with or otherwise transform
403 irrelevant features from the network via reweighting, such that new images are cast into a weighted
404 subspace of the original feature space. The representational geometry of this subspace serves as the
405 comparison to the brain. At the end of this procedure, we obtain a matrix of correlation scores of
406 dimension (number of subjects x number of ROIs x number of model layers x number of models).

407 A.5 Intersubject Predictivity and the Noise Ceiling

408 In general, the purpose of a noise ceiling is to estimate (at the level of an individual unit of prediction)
409 how reliable the response in that unit is across time. This metric allows us to then quantify how well
410 our response data at one point predicts our response data at another. One example such measure
411 relevant to fMRI is the Spearman-Brown-corrected split-half reliability of a voxel response over
412 sequential presentations of the same stimuli. However, this method tends to underestimate true voxel
413 reliability in regimes with few presentations.

414 The alternative we have provided here – the pairwise inter-subject representational similarity ref-
415 erence – is straightforward in its calculation, and computationally equivalent to the procedure for
416 benchmarking the models with classical RSA (which is to say, that subject RDMs and model RDMs
417 were computed in the exact same way, and compared using the same correlation metric).

418 As a reference point for weighted RSA, however, this threshold is perhaps a bit misleading – since
419 only the models benefit from the reweighting. One possible alternative, similar to work done recently
420 in the neural network modeling of mouse visual cortex²⁰, is to directly incorporate the neural
421 activity of human brains into a regression procedure wherein the regressand is the neural activity of a
422 target subject and the regressors are the neural activities of other subjects. This procedure has the
423 advantage of equating the set of computational (sub)operations that map model feature spaces to the
424 brain, and of providing similarly intuitive targets that undergird inferences over how much of the
425 variance in a target biological system we can capture with a system that is decidedly not biological.

426 As a preliminary test of weighted inter-subject predictivity, we consider another version of the
427 pairwise metric above, predicting single subjects using data from other single subjects. For each pair
428 of subjects, in which one is the target and the other is the contrast, we iterate over ROIs, gathering
429 all voxels from the contrast's ROI to serve as regressors in the prediction of activity in each of the

²⁰Note that after the SRP procedure there is no longer an interpretable mapping between individual model features and brain voxels; nonetheless, we have confirmed empirically that SRP procedure yields similar brain predictivity compared to a control analysis using AlexNet, a model whose feature map dimensionality is sufficiently low to run our encoding procedure without SRP.

Figure 4: Generalization scores across subject and ROI. Each point in red is the LOOCV score for a given model over the 1000 training images; each point in blue is the generalization to 1000 unseen images never incorporated into the training procedure.

430 target's ROI voxels. We repeat this procedure until we have predicted all voxels in a given target
431 subject with all possible contrast subjects. The mapping procedure in this case is exactly the same as
432 it was for the mapping of models, controlled even to the hyperparameter: we project the ROI voxel
433 activity from the contrast subject to 5960 sparse random projections, and regress these subjects to the
434 target voxel with a ridge regression set to a lambda penalty of

435 While this method equates each computational (sub)operation between model and human, the
436 intersubject predictivity threshold it establishes is even lower than the version without reweighting.
437 One reason for this may be that the brain activity from a single subject does not provide a sufficient
438 breadth of variance to benefit from the reweighting. As a first pass at rectifying this issue, we devised
439 a new measure, predicting each voxel from the concatenated activity of all voxels from all other
440 subjects in the target ROI, effectively creating a multi-human reference. While we are continuing
441 to assess, conceptually, whether such a reference point may be useful, the estimates it produces
442 for individual subjects are indeed far higher than the estimates of either the unweighted individual
443 subject-to-subject comparison or the corresponding weighted comparison, and is in most cases far
444 higher than the observed levels of model prediction. Figure 5 shows a comparison between the
445 different kinds of human reference points we compute.

446 A.6 Generalization Scores across Subject + ROI

447 Figure 4 shows the generalization scores across individual subjects and individual ROIs.

448 A.7 Compute Required

449 We used a single machine with 8 Nvidia RTX 3090 GPUs, 755gb of RAM, and 96 CPUs. GPUs were
450 used only for extracting model activations, and could (without major slowdown) be removed from
451 the analytic pipeline. Dimensionality reduction and regression computations were CPU and RAM
452 intensive. Replicating all of our results would take approximately two weeks on a similar machine.

