

Model manifold analysis suggests the human visual brain is less like an optimal classifier and more like a feature bank

Colin Conwell

MIT CSAIL

CONWELL@MIT.EDU

Michael F. Bonner

Johns Hopkins University

MF BONNER@JHU.EDU

Editors: List of editors' names

Abstract

What do deep neural network (DNN) models actually tell us about the computational principles of visual information-processing in the biological brain? A now common finding in visual neuroscience is that many different kinds of DNNs – with different architectures, tasks, and training diets – are all comparably performant predictors of image-evoked brain activity in the ventral visual cortex. This relative parity of highly diverse models may at first seem to undermine the common intuition that we can use these models to infer the computational principles that govern the visual brain. In this work, we show to the contrary that comparable brain-predictivity does not preclude the differentiation of these same models in terms of the underlying manifold geometries that define them. To do this, we assess 12 manifold geometry metrics computed across a diverse set of 117 DNN models, curated to include multiple tasks, architectures, and input diets. We then use these metrics to predict how well each model aligns with occipitotemporal cortex (OTC) activity from the human fMRI Natural Scenes Dataset. We find that *manifold signal-to-noise ratio* (a metric previously associated with few-shot learning) is a robust predictor of downstream brain-alignment and supersedes both other manifold geometry metrics (i.e. *manifold capacity*) and downstream task-performance (e.g. top-k recognition accuracy) across multiple different image sets (e.g. ImageNet21K versus Places365) and controlled model comparisons (e.g. assessments across ImageNet-1K trained architectural variants only). These results add to a growing body of evidence that the ventral visual stream serves as a basis set (or feature vocabulary) for object recognition rather than as the actual locus of recognition *per se*.

1. Background + Methods

High-level ventral visual cortex is widely considered the primary substrate of object recognition in the human brain [17; 18; 27]. This notion has been further reinforced in recent years by the now-seminal finding that deep neural network (DNN) models supporting image classification are the most predictive models of ventral visual brain activity to date [52; 43; 13]. One lingering issue with this formulation, however, is the under-specification of what we mean by the word “recognition” as implemented by the biological brain. What *precise* neural-computational mechanisms instantiate this process? And how might we use deep neural network models to elucidate them? In models, the computational mechanisms of ‘recognition’ are explicitly (i.e. algorithmically) defined. Canonical, ‘end-to-end’ image classification models (e.g. AlexNet [30]) ‘recognize’ images by the direct, nonlinear mapping

of a tensorized image onto one-hot-encodings (point indices) of human-defined category labels. Self-supervised, contrastive-learning models (e.g. SimCLR [9]) and masked decoding models (e.g. DINO, [8]) also support recognition, but do so first by learning (without labels) the invariances and selectivities that define the axes of an input image space more generally. ‘Zero-shot’- multimodal classifiers (e.g. CLIP [40]) ‘recognize’ images by way of a top-k similarity score between the language embeddings of category labels and the vision embeddings of candidate images.

Identifying which of these (or many other) motifs is most brain-like has been a process somewhat confounded by the fact that many (if not most) of these models often achieve roughly comparable (often high) measures of ‘representational alignment’ [50] to human visual cortex activity recorded via neuroimaging and neurophysiology alike [13; 49; 24; 38; 14]. One response to this relative parity (inspired by the neuro-connectionist research programme [19]) has been the use of controlled model comparisons designed to better parse the differences between models by *grouping them* in ways empirically isolate the representational influence of broader computational design principles such as sensory diet (training data), architecture, and task [15]. In this work, we attempt to better parse the differences between models more directly by *sorting them* along scalar axes of representational structure defined by manifold geometry analysis [12; 11; 46]. In so doing, we attempt as well to more directly interrogate how this structure relates to the *function* of object recognition, and what our neural networks may or may be telling us about the nature of this structure in the human ventral visual system.

Related work The use of manifold geometrics to understand representational structure (brain-like or otherwise) is a technique simultaneously steeped in a rich history of theory, and an actively developing area of interest in modern machine learning and cognitive neuroscience [35; 48; 7; 21; 31; 47; 45]. Extended discussion of these related works is available in Appendix A1.

Model & metric curation Our general approach was to first curate N=12 metrics of manifold geometry from two distinct lines of work: N=4 metrics from Chung et al. [12] (derived from replica mean field theory) and N=8 metrics from Sorscher et al. [46] (derived from principal components analysis). We then curated a set of N=117 candidate deep neural network models spanning different visual diets, architectures, and tasks, based on the work of Conwell et al. [15]. (See Appendix A2 for a list of model sources.)

Measuring brain-alignment Model-to-brain predictivity scores (henceforth, brain-alignment) was computed for each model first by computing the average voxelwise encoding score over a 50% train split (N=500) of the ‘Shared NSD1000’ image set [1] for each individual layer with field-standard regularized ridge regression. Models were then compared on the basis of the average encoding scores over the test split (N=500) images for the most brain-like layer (as determined via nested cross-validation on the train split).

Manifold geometry analysis Manifold geometry metrics (from both Chung et al. [12] and Sorscher et al. [46]) were computed using N=50 ‘concept manifolds’ (the representational matrices of 50 images from each of 50 category labels from ImageNet-1K/21K [16; 42], or Places365 [54]). To ensure proper evaluation of these metrics on our candidate representations, we used the exact GitHub source code associated with the defining works of each: [schung039/neural-manifolds-replicaMFT](#) [12; 48] and [bsorsch/geometry-fewshot-learning](#) [46].

2. Results + Analysis

Primary analysis: explaining brain-alignment with manifold geometry How well do differences in underlying manifold geometries explain downstream brain-alignment? To answer this, we compute the rank-order correlation ($\rho_{Spearman}$) between the models’ brain-predictivity (mean encoding) scores and each of our 12 manifold metrics (e.g. manifold capacity, effective dimensionality) computed from the N=50 concept manifolds of ImageNet-21K used in [48]; results from this analysis are shown in Appendix Figure 1C and Table 1.

6 / 12 metrics showed significant rank-order correlation with brain-predictivity at $p < 0.01$. 3 of these metrics (manifold signal-to-noise ratio, capacity, and signal) showed significant, positive correlations. Three (manifold radius, dimensionality, and within-concept radius) showed significant, negative correlations. The overall strongest predictor of downstream brain-predictivity was the manifold signal-to-noise ratio from Sorscher et al. [46], with a strikingly high rank-order correlation of $\rho_{Spearman} [\pm 95\%BCI] = 0.798 [0.731, 0.895]$, $p = 4.70e-27$. Not far behind this, however, was the negatively correlating manifold radius metric from Chung et al. [12], at $\rho = -0.724 [-0.857, -0.623]$, and the positively correlating manifold capacity metric from Chung et al. [12], at $\rho = 0.616 [0.493, 0.779]$. In sum, the primary takeaway from this first analysis is that manifold metrics can indeed be the strong predictors of brain-alignment that theorists and empiricist alike have proposed they might be [46; 51; 7; 21; 22].

Querying manifold metric robustness & interpretability across experimental subconditions The strong, significant correlation of *multiple* manifold metrics with downstream brain-alignment does, however, raise new questions. One of these is: When metrics describe otherwise divergent properties of a manifold’s geometry (but both explain downstream brain-alignment to similar degrees), which geometry should we take to be more brain-like? In our survey, two metrics in particular – manifold capacity and signal-to-noise ratio (henceforth SNR) – instantiate a rather palpable case of this algorithmic ambiguity. In the extreme (as in the final output layer of end-to-end trained category-supervised DNNs), high manifold capacity signals the complete collapse of all category information into fully separate, single point-estimates of category identity [12; 48; 11] – in other words, an optimal classifier. In a similar extreme, manifold signal-to-noise ratio also collapses to perfect point-estimates of category identity (all signal, no noise). As described by Sorscher et al. [46], however, mid-to-high range values of signal-to-noise ratio better describe the conditions of better few-shot learning algorithms (something end-to-end category-supervised neural classifiers tend notably *not* to be). Which of these metrics, then, better describes the object-recognition-supporting representational motifs instantiated in our DNN models of the ventral visual stream? To better resolve this ambiguity, we assessed the correlations of manifold capacity and SNR to downstream brain-predictivity in a series of experimental sub-conditions designed both to test the robustness of these correlations and to disambiguate the somewhat competing representational hypotheses they entail. (Details from Experiments 1, 2 and 3 are shown in Appendix Tables 2 and 3.)

Experiment 1: Measuring robustness across model (sub)sets In our first experiment, we assessed the correlations of manifold capacity and SNR in increasingly smaller,

targeted subsets of our otherwise diversely sampled model set. The first subset we assessed was a subset of models we call the "high-performing" set: effectively all models above a notable visual elbow in brain-predictivity first described in [15], but seen also in our sample (see Figure 1C-Left). Quantifying this elbow with a segmented regression analysis yielding a breakpoint of $\psi = 84.1$ [72.4, 94.9], we defined this 'high-performing' set as the N=84 models with average encoding scores of $r_{\text{Pearson}} = 0.336$ or higher. Here, already, manifold capacity begins to diverge from manifold SNR in its rank-order correlation with brain-predictivity, with capacity diminishing to the point of non-significance at $\rho = 0.14$, $p = 0.222$ and manifold SNR remaining substantial and significant at $\rho = 0.50$, $p < .001$. In an even smaller subset of models varying only in architecture (the N=53 category-supervised ImageNet-1K-trained models from the Torchvision model zoo [36]), the trend is similar, with ρ for manifold capacity diminishing to 0.16, $p = 0.262$ and manifold SNR remaining high at 0.471, $p < 0.001$. In short, manifold SNR persists as a predictor of downstream brain-alignment even in very restricted ranges; manifold capacity does not.

Experiment 2: Layer-wise analysis of category-recognition models To better understand the difference between manifold capacity and manifold SNR we were observing in this smaller subset, we next probed variation in the correlation of manifold metric and brain-predictivity in layers beyond the most brain-predictive layer selected by our initial cross-validation, and in particular, the 'last hidden layer' feeding into the one-hot, category-encoding output. In effect, in this smaller subset of category-supervised models, this layer instantiates the representation most directly responsible for the 'recognition' behavior the model will output for any given input, and by association, is the layer we might presumably observe the highest covariance between manifold capacity and manifold SNR. And indeed, what we observe here is that the correlations of both manifold capacity and manifold SNR with downstream brain-predictivity change dramatically. Here, in this final hidden vector of models trained to collapse category information onto the single points of the output layer, manifold capacity and manifold SNR are shown to be strong, significant, *negative* predictors of downstream brain-alignment, with $\rho = -0.606$, $p < .001$ and -0.663 , respectively. This sign-reversal corresponds to substantial increases in the scalar values of both metrics relative to the most brain-predictive layers (with shifts in the max values of manifold capacity increasing 182.17% from .129 to .235 and signal-to-noise ratio increasing 209.12% from 4.131 to 8.639.) What is happening, in effect, is that the more the models are successful in collapsing category information to single point-estimates at this final hidden layer, the less predictive of downstream brain-alignment they will be. Notably, the trend is similar if, as with our manifold metrics, we correlate the ImageNet-1K classification accuracy of these models with downstream brain-predictivity, a trend we find (in line with recent work, e.g. [33]) to be strongly negative across our 53 models ($\rho = -0.63$, $p < .001$).

Experiment 3: Differential, IID/OOD concept manifold sampling In a final experiment, we recomputed each of our manifold metrics with two new sets of N=50 concept manifolds (N=50 test set images each): one from the object categories of ImageNet1K (versus the ImageNet21K sample we use in our main analysis, following the protocol and codebase of Stephenson et al. [48]), and another from the *scene* categories of Places365. The logic here is that these instantiate two different levels of 'generalization' for our category-supervised ImageNet-1K-trained models, one nearer (in-distribution, IID), one farther (*pos-*

sibly out-of-distribution, OOD). Manifold capacity, in this regime, should be higher for those concepts that are IID (i.e. the ImageNet1K sample) than OOD (the Places365), again instantiating the progressive tightening of category information towards single point-estimates. Manifold SNR, on the other hand, will also decrease. But supporting few-shot-learning, as it nominally does, manifold SNR will also maintain information that bridges the gaps between separable categories. Is this the kind of information that boosts manifold SNR’s explanatory power for downstream brain-alignment? The results of this experiment suggest it might be: In the most brain-predictive layers, for example, we observe that manifold SNR remains a significant, positive predictor of downstream brain-alignment in both the new ImageNet1K concept manifold sample *and* in the Places365 sample, with $\rho = 0.457$, $p < .001$ and 0.319 , $p = .002$, respectively. Manifold capacity (as in the ImageNet21K sample) is *not* significantly predictive of downstream brain-alignment in either of these cases.

3. Discussion

What factors make for a ‘good’ neural network model of the visual brain? Since the adoption in visual neuroscience of the task-optimized deep convolutional neural network model more than a decade ago [53], the dominant – and in some ways most empirically defensible answer – has largely been ‘prediction’: Better models of the visual brain are those models whose internal representations most accurately *predict* the activity patterns of the biological brain. For those seeking downstream control or causal perturbation of biological systems [4; 44], this answer may be sufficient. For those seeking ‘understanding’, the search remains for other forms of explanatory variables that supplement raw prediction with the parsimony of theories articulable in formal or natural language [41; 19; 22].

In this work, we attempt to instrumentalize the emergent framework of neural manifold geometry [11] to better understand the underlying *structural* factors that make certain models of ventral visual cortex ‘better’ (i.e. more predictive of brain activity) than others – in effect, by using the metric scalars of manifold geometry to more directly link *representation* to *function*. Through this lens, we return to the seminal question of how the representation in ventral visual cortex relates to the function of object recognition (i.e. the ‘readout’ of a category label). First validating the second-order predictive power of manifold geometry metrics (i.e. the strength of their rank-order correlation with downstream brain-alignment), we find multiple candidate metrics that nevertheless instantiate divergent hypotheses about the object-recognition-supporting representations of the ventral stream. Testing these hypotheses in a series of experimental subconditions, we find that a metric (manifold capacity) whose value scales with representational convergence toward separable, single point-estimates of category identity is less robust in predicting downstream brain-alignment than a metric (manifold SNR) that accounts for more graded forms of invariance and separability that still subserve recognition, but equally subserve the few-shot learning of new categories. Taken together, these results add to a growing body of evidence from across multiple experimental modalities [34; 26; 5; 6; 20; 10; 37] that suggest the ventral visual stream may be less like an optimal classifier (i.e. the locus of recognition itself) and more like a feature bank (i.e. the vocabulary of whatever compositional process of recognition involves downstream.) Further elaboration of concurrent work as well as the limitations of the current work are available in Appendix A1.

4. Code + Data

Source code and data for this work will be made available in our project GitHub:

github.com/ColinConwell/BMM-Geometrics

References

- [1] Emily Jean Allen, Ghislain St-Yves, Yihan Wu, Jesse L Breedlove, Logan T Dowdle, Bradley Caron, Franco Pestilli, Ian Charest, J Benjamin Hutchinson, Thomas Naselaris, et al. A massive 7t fmri dataset to bridge cognitive and computational neuroscience. *bioRxiv*, 2021.
- [2] Adrien Bardes, Jean Ponce, and Yann LeCun. Vicreg: Variance-invariance-covariance regularization for self-supervised learning. *arXiv preprint arXiv:2105.04906*, 2021.
- [3] Adrien Bardes, Jean Ponce, and Yann LeCun. Vicregl: Self-supervised learning of local visual features. *Advances in Neural Information Processing Systems*, 35:8799–8810, 2022.
- [4] P. Bashivan, K. Kar, and J.J. DiCarlo. Neural population control via deep image synthesis. *Science*, 364, 2019. doi: 10.1126/science.aav9436. URL <https://doi.org/10.1126/science.aav9436>.
- [5] Tyler Bonnen, Daniel LK Yamins, and Anthony D Wagner. When the ventral visual stream is not enough: A deep learning account of medial temporal lobe involvement in perception. *Neuron*, 109(17):2755–2766, 2021.
- [6] Tyler Bonnen, Anthony D Wagner, and Daniel LK Yamins. Medial temporal cortex supports object perception by integrating over visuospatial sequences. *Cognition*, 262: 106135, 2025.
- [7] Abdulkadir Canatar, Jenelle Feather, Albert Wakhloo, and SueYeon Chung. A spectral theory of neural prediction and alignment. *Advances in Neural Information Processing Systems*, 36:47052–47080, 2023.
- [8] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging Properties in Self-Supervised Vision Transformers. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2021. doi: 10.1109/ICCV48922.2021.00951.
- [9] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International Conference on Machine Learning*, pages 1597–1607. PMLR, 2020. arXiv preprint arXiv:2002.05709.
- [10] Zirui Chen and Michael F Bonner. Universal dimensions of visual representation. *Science Advances*, 11(27):eadw7697, 2025.

- [11] SueYeon Chung and LF Abbott. Neural population geometry: An approach for understanding biological and artificial neural networks. *Current opinion in neurobiology*, 70:137–144, 2021.
- [12] SueYeon Chung, Daniel D Lee, and Haim Sompolinsky. Classification and geometry of general perceptual manifolds. *Physical Review X*, 8(3):031003, 2018.
- [13] Radoslaw Martin Cichy, Gemma Roig, Alex Andonian, Kshitij Dwivedi, Benjamin Lahner, Alex Lascelles, Yalda Mohsenzadeh, Kandan Ramakrishnan, and Aude Oliva. The algonauts project: A platform for communication between the sciences of biological and artificial intelligence. *arXiv preprint arXiv:1905.05675*, 2019.
- [14] David Colaço, Russell A Poldrack, Aliya Rumana, Thierault Jordan, Daniel C Burnston, Sofie Valk, Philipp Haueis, Daniel Marguiles, and Carl F Craver. Do dnns explain the visual system? guidelines for a better debate about explanation. 2025.
- [15] Colin Conwell, Jacob S Prince, Kendrick N Kay, George A Alvarez, and Talia Konkle. A large-scale examination of inductive biases shaping high-level visual representation in brains and machines. *Nature communications*, 15(1):9383, 2024.
- [16] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. doi: 10.1109/CVPR.2009.5206848.
- [17] James J DiCarlo and David D Cox. Untangling invariant object recognition. *Trends in cognitive sciences*, 11(8):333–341, 2007. doi: 10.1016/j.tics.2007.06.010. Publisher: Elsevier.
- [18] J.J. DiCarlo, D. Zoccolan, and N.C. Rust. How Does the Brain Solve Visual Object Recognition? *Neuron*, 73(3):415–434, 2012. doi: 10.1016/j.neuron.2012.01.010. URL <https://doi.org/10.1016/j.neuron.2012.01.010>.
- [19] Adrien Doerig, Rowan P Sommers, Katja Seeliger, Blake Richards, Jenann Ismael, Grace W Lindsay, Konrad P Kording, Talia Konkle, Marcel AJ Van Gerven, Nikolaus Kriegeskorte, et al. The neuroconnectionist research programme. *Nature Reviews Neuroscience*, pages 1–20, 2023.
- [20] Fenil R Doshi and Talia Konkle. Cortical topographic motifs emerge in a self-organized map of object space. *Science Advances*, 9(25):eade8187, 2023.
- [21] Eric Elmoznino and Michael F Bonner. High-performing neural network models of visual cortex benefit from high latent dimensionality. *PLoS computational biology*, 20(1):e1011792, 2024.
- [22] Jenelle Feather, Meenakshi Khosla, N Murty, and Aran Nayebi. Brain-model evaluations need the neuroai turing test. *arXiv preprint arXiv:2502.16238*, 2025.
- [23] Priya Goyal, Quentin Duval, Jeremy Reizenstein, Matthew Leavitt, Min Xu, Benjamin Lefaudeaux, Mannat Singh, Vinicius Reis, Mathilde Caron, Piotr Bojanowski,

- Armand Joulin, and Ishan Misra. VISSL, 2021. URL <https://github.com/facebookresearch/vissl>.
- [24] Yena Han, Tomaso A Poggio, and Brian Cheung. System identification of neural systems: If we got it right, would we know? In *International conference on machine learning*, pages 12430–12444. PMLR, 2023.
- [25] Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan Taori, Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hannaneh Hajishirzi, Ali Farhadi, and Ludwig Schmidt. OpenCLIP, 2021. URL <https://doi.org/10.5281/zenodo.5143773>.
- [26] Akshay V Jagadeesh and Justin L Gardner. Texture-like representation of objects in human visual cortex. *Proceedings of the National Academy of Sciences*, 119(17):e2115302119, 2022.
- [27] Nancy Kanwisher. Functional specificity in the human brain: a window into the functional architecture of the mind. *Proceedings of the National Academy of Sciences*, 107(25):11163–11170, 2010. doi: 10.1073/pnas.1005062107. Publisher: National Acad Sciences.
- [28] Nancy Kanwisher, Meenakshi Khosla, and Katharina Dobs. Using artificial neural networks to ask ‘why’ questions of minds and brains. *Trends in Neurosciences*, 2023.
- [29] Talia Konkle and George A Alvarez. A self-supervised domain-general learning framework for human ventral stream representation. *Nature Communications*, 13(1):1–12, 2022. doi: 10.1038/s41467-022-28091-4. Publisher: Nature Publishing Group.
- [30] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. ImageNet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems*, pages 1097–1105, 2012.
- [31] Michael Kuoch, Chi-Ning Chou, Nikhil Parthasarathy, Joel Dapello, James J. DiCarlo, Haim Sompolinsky, and SueYeon Chung. Probing biological and artificial neural networks with task-dependent neural manifolds. In Yuejie Chi, Gintare Karolina Dziugaite, Qing Qu, Atlas Wang Wang, and Zhihui Zhu, editors, *Conference on Parsimony and Learning*, volume 234 of *Proceedings of Machine Learning Research*, pages 395–418. PMLR, 03–06 Jan 2024. URL <https://proceedings.mlr.press/v234/kuoch24a.html>.
- [32] Dongxu Li, Junnan Li, Hung Le, Guangsen Wang, Silvio Savarese, and Steven C.H. Hoi. LAVIS: A one-stop library for language-vision intelligence. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, pages 31–41, Toronto, Canada, July 2023. Association for Computational Linguistics. URL <https://aclanthology.org/2023.acl-demo.3>.
- [33] Drew Linsley, Ivan F Rodriguez, Thomas Fel, Michael Arcaro, Saloni Sharma, Margaret Livingstone, and Thomas Serre. Performance-optimized deep neural networks

- are evolving into worse models of inferotemporal visual cortex. *arXiv preprint arXiv:2306.03779*, 2023.
- [34] Bria Long, Chen-Ping Yu, and Talia Konkle. Mid-level visual features underlie the high-level categorical organization of the ventral stream. *Proceedings of the National Academy of Sciences*, 115(38):E9015–E9024, 2018. doi: 10.1073/pnas.1719616115. Publisher: National Acad Sciences.
 - [35] Bruno A Olshausen, David J Field, and others. Sparse coding of natural images produces localized, oriented, bandpass receptive fields. *Submitted to Nature. Available electronically as ftp://redwood.psych.cornell.edu/pub/papers/sparse-coding.ps*, 1995. Publisher: Citeseer.
 - [36] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019. URL <http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf>.
 - [37] Jacob S Prince, George A Alvarez, and Talia Konkle. Contrastive learning explains the emergence and function of visual category-selective regions. *Science Advances*, 10(39):ead1776, 2024.
 - [38] Jacob S Prince, Colin Conwell, George A Alvarez, and Talia Konkle. A case for sparse positive alignment of neural systems. In *ICLR 2024 Workshop on Representational Alignment*, 2024.
 - [39] Jacob S Prince, Bin Xu Wang, Akshay V Jagadeesh, Thomas Fel, Emily Lo, George A Alvarez, Margaret S Livingstone, and Talia Konkle. Feature accentuation along the encoding axes of it neurons uncovers hidden differences in model-brain alignment. *Journal of Vision*, 25(9):2215–2215, 2025.
 - [40] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, and others. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763, 2021. tex.organization: PMLR arXiv preprint arXiv:2103.00020.
 - [41] Blake A Richards, Timothy P Lillicrap, Philippe Beaudoin, Yoshua Bengio, Rafal Bogacz, Amelia Christensen, Claudia Clopath, Rui Ponte Costa, Archy de Berker, Surya Ganguli, and others. A deep learning framework for neuroscience. *Nature Neuroscience*, 22(11):1761–1770, 2019. doi: 10.1038/s41593-019-0520-2. Publisher: Nature Publishing Group US New York.

- [42] Tal Ridnik, Emanuel Ben-Baruch, Asaf Noy, and Lihi Zelnik-Manor. Imagenet-21k pretraining for the masses. *arXiv preprint arXiv:2104.10972*, 2021.
- [43] Martin Schrimpf, Jonas Kubilius, Ha Hong, Najib J. Majaj, Rishi Rajalingham, Elias B. Issa, Kohitij Kar, Pouya Bashivan, Jonathan Prescott-Roy, Franziska Geiger, Kailyn Schmidt, Daniel L. K. Yamins, and James J. DiCarlo. Brain-Score: Which Artificial Neural Network for Object Recognition is most Brain-Like? *bioRxiv preprint*, 2018. doi: 10.1101/407007.
- [44] Thomas Serre. Deep learning: the good, the bad, and the ugly. *Annual Review of Vision Science*, 5:399–426, 2019. doi: 10.1146/annurev-vision-091718-014951. Publisher: Annual Reviews.
- [45] Zhenan Shao, Yiqing Zhou, and Diane M Beck. Human visual robustness emerges from manifold disentanglement in the ventral visual stream. *Journal of Vision*, 25(9):1967–1967, 2025.
- [46] Ben Sorscher, Surya Ganguli, and Haim Sompolinsky. Neural representational geometry underlies few-shot concept learning. *Proceedings of the National Academy of Sciences*, 119(43):e2200800119, 2022.
- [47] Ghislain St-Yves, Kendrick Kay, and Thomas Naselaris. Variation in the geometry of concept manifolds across human visual cortex. *bioRxiv*, pages 2024–11, 2024.
- [48] Cory Stephenson, Jenelle Feather, Suchismita Padhy, Oguz Elibol, Hanlin Tang, Josh McDermott, and SueYeon Chung. Untangling in invariant speech recognition. In *Advances in Neural Information Processing Systems*, pages 14368–14378, 2019.
- [49] Vighnesh Subramaniam, Colin Conwell, Christopher Wang, Gabriel Kreiman, Boris Katz, Ignacio Cases, and Andrei Barbu. Revealing vision-language integration in the brain with multimodal networks. *ArXiv*, pages arXiv–2406, 2024.
- [50] Ilia Sucholutsky, Lukas Muttenthaler, Adrian Weller, Andi Peng, Andreea Bobu, Been Kim, Bradley C Love, Erin Grant, Jascha Achterberg, Joshua B Tenenbaum, et al. Getting aligned on representational alignment. *arXiv preprint arXiv:2310.13018*, 2023.
- [51] Greta Tuckute, Jenelle Feather, Dana Boebinger, and Josh H McDermott. Many but not all deep neural network audio models capture brain responses and exhibit correspondence between model stages and brain regions. *Plos Biology*, 21(12):e3002366, 2023.
- [52] Daniel LK Yamins and James J DiCarlo. Using goal-driven deep learning models to understand sensory cortex. *Nature neuroscience*, 19(3):356, 2016. doi: 10.1038/nn.4244. Publisher: Nature Publishing Group.
- [53] Daniel LK Yamins, Ha Hong, Charles F Cadieu, Ethan A Solomon, Darren Seibert, and James J DiCarlo. Performance-optimized hierarchical models predict neural responses in higher visual cortex. *Proceedings of the National Academy of Sciences*, 111(23):8619–8624, 2014. doi: 10.1073/pnas.1403112111. Publisher: National Acad Sciences.

- [54] Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10 million image database for scene recognition. *IEEE transactions on pattern analysis and machine intelligence*, 40(6):1452–1464, 2017.

Appendix + Supplementary Information

A1. Precedent, Parallel, and Concurrent Work; Limitations

The use of manifold geometrics in cognitive neuroscience and machine learning has recently resurged, but more generally has been a key methodology in both of these disciplines for many decades (consider e.g. [35] for a paradigmatic example). More recent work in this domain has seen manifold geometry analysis applied to *speech* recognition in (biological and artificial) auditory neural networks [48]; in theoretical work on the limits or underpinnings of representational alignment more generally [7; 21]; to the characterizing of task-optimized models predicting activity in macaque visual cortex [31]; and (most directly relevant to the current analysis) to the characterizing of brain-like representational structure in human ventral visual system activity sourced from the Natural Scenes Dataset [47; 45]. St-Yves et al. [47] in particular is notable for the use of the same manifold signal-to-noise and dimensionality metrics [46] we leverage in this analysis, deployed in service of characterizing how geometry varies across different (sub)regions of the ventral and dorsal stream alike.

Our work builds and extends on these approaches by explicitly yoking the more distal, ecological insights of the neuro-connectionist, many-model / controlled comparisons approach with the more proximate, structural insights of manifold geometry - and in so doing, ideally, to get just a little closer to unifying the goals of ‘prediction’ and ‘understanding’ that so often seem in tension in the application of DNN models to neuroscience data [28].

In finding that a metric more in tune with few-shot learning than optimal classification seems to be a better explanatory variable of a representation’s downstream alignment to the ventral stream, our work seems to resonate well with recent findings from multiple other research programmes. Whether it be the finding that the biological ventral stream in macaques may be more ‘texture-like’ than previously assumed on the basis of shape-based behavioral biases [26]; the greater impairment to object recognition from damage to medial temporal lobe than ventral stream in double dissociation neuropsychology experiments [5; 6]; the demonstration of emergent category-like topographic structure from self-organized maps learned over self-supervised natural image models [20]; or even more simply the now robustly reproduced finding (inherent to these results as well) that the later layers of category-recognizing deep neural network models are less brain-predictive than more intermediate layers [15], the ventral stream seems in many cases to be preserving information in ways that diverge from the strict motif of category invariance and separability that one might assume were one to assume the function of this region is the kind of classification subserved by the models that best predict its activity.

A major limitation of the current work (and similar efforts), of course, is that the inferences we make about the manifold geometry of the brain are inferences we make by proxy of the brain-predictive models that are the actual targets of our analysis. Increasingly, however, we have seen that however much ‘proxy’ the inferences we make from models may be, these inferences are often directly convertible into the kinds of casual / perturbational empiricism that are the gold standard of scientific understanding for a target biological system. Concurrent work is already applying manifold geometry analysis and similar towards the goal of robustifying brain-predictive models through explicit, manifold metric-guided representational alignment and neural control [45; 39]. Future work should continue to leverage the structural grip of manifold analysis with the high-throughput de-

velopment of ever-more competent task-optimized models to accelerate the incisive style of this empiricism even further.

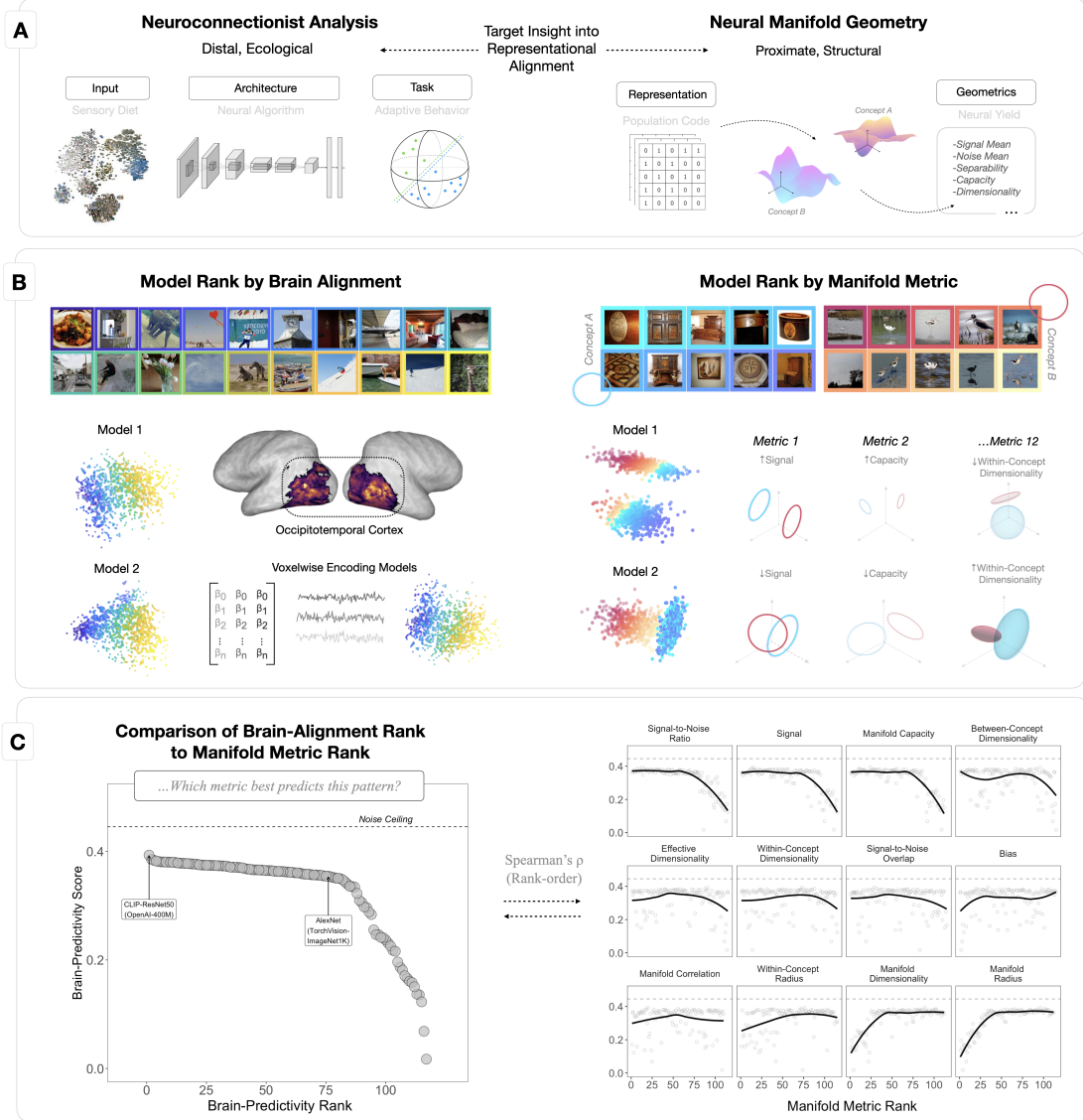


Figure 1: An overview of our motivation, methodology, and primary analysis. An extended caption is available in the [paragraph below](#).

Figure 1 extended caption In **A**, we schematize the key factors that define and contrast the two animating frameworks of our model-to-brain comparisons. With its focus on the influence of input, architecture, and task, neuroconnectionist-style analysis [19; 15] seeks

to use models as proxies of the design constraints (i.e. ‘pressures’) that could in principle have shaped the emergence of the representational structure we observe in the biological brain. Neural manifold geometry [11; 46] seeks to use models as a direct empirical substrate for probing how the differences in the structure of representation (both within and across models) contribute to differences in downstream behavior (i.e. classification, or in this case, brain-alignment / predictivity). In **B**, we schematize our primary method for using manifold geometry to interpret brain-alignment in high-level ventral visual stream (occipitotemporal cortex, or OTC). We first rank models in terms of their brain-alignment scores by computing average voxelwise encoding scores with field-standard cross-validated ridge-regressions for the ‘Shared-NSD1000’ [1]. We then take these same models and rank them according to each of our curated manifold geometry metrics computed over ‘concept manifolds’ (representational matrices) composed from ImageNet categories. In **C**, we show the primary outcome of this analysis: A rank-order comparison (Spearman’s ρ correlation) between the rank of each model according to its brain-predictivity, and the rank of each model according to its associated manifold geometry metrics. As shown, the trend in the brain-predictivity plot on the left (with brain-predictivity score in units of Pearson’s r on the y axis, and brain-predictivity rank on the x axis) is better and worse captured by the various metrics in the subplots on the right (with manifold metric rank in place of the the brain-predictivity on the x axis), which are sorted from top to bottom by Spearman’s ρ .

A2. Model Selection

Following the method of Conwell et al. [15], we curated a set of 117 deep neural network (DNN) models spanning different visual input diets (training data), architectures, and tasks. These models were sourced from the following repositories:

- the Torchvision model zoo [36];
- the VISSL model zoo [23];
- the DINO collection [25];
- the OpenAI CLIP collection [40];
- the OpenCLIP model zoo [25];
- the VicReg(-L) collections; [2; 3];
- the Salesforce-LAVIS model zoo [32];
- and Open-IPCL collection [29]

Additional information in the Project GitHub: [ColinConwell/BMM-Geometrics](#)

Manifold Metric	$r_{\text{Spearman}} \pm 95\% \text{ CI}$	p
<i>Chung et al.</i>		
Capacity	0.643 [0.527, 0.796]	.001
Correlation	-0.077 [-0.266, 0.107]	0.42
Dimensionality	-0.594 [-0.747, -0.471]	.001
Radius	-0.721 [-0.856, -0.62]	.001
<i>Sorscher et al.</i>		
Signal-to-Noise Ratio	0.774 [0.706, 0.873]	.001
Signal	0.569 [0.445, 0.724]	.001
Between-Concept Dimensionality	0.269 [0.103, 0.45]	0.003
Effective Dimensionality	0.154 [-0.018, 0.336]	0.097
Signal-to-Noise Overlap	0.133 [-0.04, 0.317]	0.153
Within-Concept Dimensionality	0.073 [-0.104, 0.26]	0.433
Within-Concept Radius	-0.239 [-0.429, -0.068]	0.009
Bias	-0.242 [-0.409, -0.088]	0.009

Table 1: Rank-order correlations between brain-predictivity and manifold geometry metrics, with bootstrapped 95% confidence intervals shown in brackets; results from the Primary Analysis.

	Manifold Capacity	Signal-to-Noise Ratio
	$\rho_{\text{Spearman}}, p$	$\rho_{\text{Spearman}}, p$
All Surveyed Models	0.62, $p < 0.001$	0.80, $p < 0.001$
High-Performing	0.14, $p = 0.222$	0.52, $p < 0.001$
ImageNet1K-Supervised	0.16, $p = 0.262$	0.47, $p < 0.001$

Table 2: A comparison of the rank-order correlations between the manifold capacity and manifold signal-to-noise ratio (SNR) metrics across progressively smaller subsets of models; results from Experiment 1.

	Manifold Capacity		Signal-to-Noise Ratio	
	Range	r_{Pearson}, p	Range	r_{Pearson}, p
ImageNet1K				
Best Layer	0.049 - 0.129	-0.070, $p = 0.637$	2.038 - 4.131	0.457, $p < 0.001$
Last Layer	0.098 - 0.235	-0.606, $p < 0.001$	1.622 - 8.639	-0.663, $p < 0.001$
ImageNet21K				
Best Layer	0.052 - 0.125	-0.132, $p = 0.373$	2.092 - 4.062	0.468, $p < 0.001$
Last Layer	0.095 - 0.160	-0.515, $p < 0.001$	1.833 - 4.467	-0.294, $p = 0.032$
Places365 (Scenes)				
Best Layer	0.047 - 0.095	-0.145, $p = 0.327$	1.677 - 3.446	0.319, $p = 0.02$
Last Layer	0.070 - 0.108	-0.515, $p < 0.001$	1.305 - 2.956	-0.158, $p = 0.259$

Table 3: Comparisons of manifold capacity and signal-to-noise ratio between the peak (i.e. most brain-predictive) layer and the last layer in 3 different probe datasets. r_{Pearson} is the correlation between the manifold metric value and brain-predictivity (mean encoding score) of each layer; results from Experiments 2 and 3.

