
Generalized Shape Metrics on Neural Representations

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Understanding the operation of biological and artificial networks remains a difficult
2 and important challenge. To identify general principles of neural computation,
3 researchers are increasingly interested in surveying large collections of networks
4 that are trained on, or biologically adapted to, similar tasks. A standardized set of
5 analysis tools is now needed to identify how network-level covariates—such as
6 architecture, learning algorithm, anatomical brain region, model organism, etc.—
7 systematically impact neural representations (hidden layer activations). Here, we
8 provide a rigorous foundation for these analyses by defining a broad family of
9 metric spaces that quantify representational dissimilarity between networks. Us-
10 ing this framework, we modify existing measures based on canonical correlation
11 analysis and centered kernel alignment to satisfy the triangle inequality, formu-
12 late a novel metric that respects the inductive biases in convolutional layers, and
13 identify approximate Euclidean embeddings that enable network representations to
14 be incorporated into essentially any off-the-shelf machine learning method. We
15 demonstrate these methods on large-scale datasets from biology (Allen Institute
16 Brain Observatory) and deep learning (NAS-Bench-101), and identify relation-
17 ships between neural representations that are interpretable in terms of anatomical
18 hierarchies and model performance.

19 1 Introduction

20 The extent to which different deep networks or neurobiological systems use equivalent represen-
21 tations in support of similar task demands is a topic of persistent interest in machine learning
22 and neuroscience [1]. Several methods including linear regression [2, 3], canonical correlation
23 analysis (CCA; [4, 5]), representational similarity analysis (RSA; [6]), and centered kernel align-
24 ment (CKA; [7, 8]) have been used to quantify the similarity of hidden layer activation patterns.
25 These similarity measures are often interpreted on an ordinal scale and are employed to compare a
26 limited number of networks—e.g., they can indicate whether networks A and B are more or less
27 similar than networks A and C . While these comparisons have yielded important insights [2–7, 9–11],
28 the underlying methodologies have not been extended to quantitative analyses spanning thousands of
29 networks.

30 To unify existing approaches and enable more sophisticated analyses, we draw on ideas from statistical
31 shape analysis [12–14] to develop dissimilarity measures that are proper metrics—i.e., are symmetric
32 and respect the triangle inequality—on neural representations. This enables several off-the-shelf
33 methods with theoretical guarantees for classification (e.g. k -nearest neighbors, [15]) and clustering
34 (e.g. hierarchical clustering [16]) on network representations. Existing similarity measures can
35 violate the triangle inequality, which complicates these downstream analyses [17–19]. However, by
36 introducing simple modifications to these methods, we can satisfy the triangle inequality and view
37 existing approaches as special cases of the framework we outline. This framework also facilitates the
38 discovery of new metrics, including metrics specialized to convolutional layer representations.

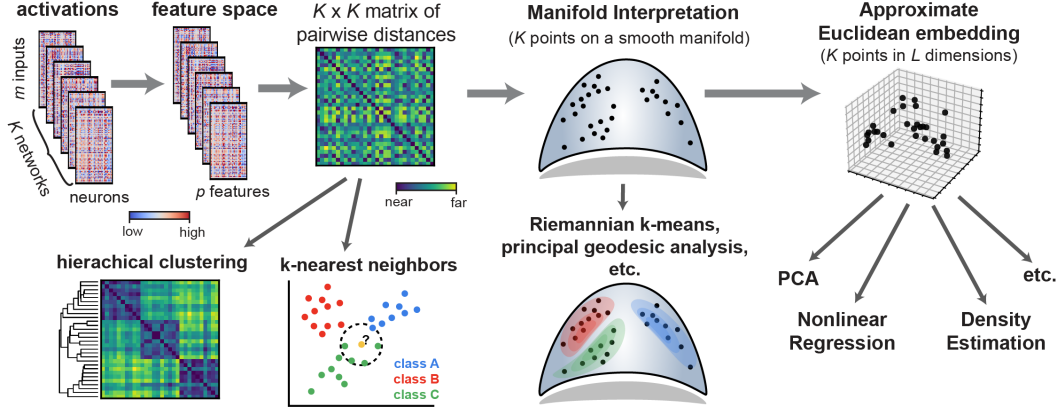


Figure 1: Machine learning workflows enabled by generalized shape metrics.

39 Moreover, we show empirically that neural representations can be embedded with low distortion
 40 into Euclidean spaces with a moderate number of dimensions. This enables an even broader variety
 41 of modeling approaches for regression, classification, and unsupervised analysis. Applying this
 42 approach to representations of visual inputs in mice (Allen Brain Observatory; [20]), we predict an
 43 anatomically derived feature of brain regions (“hierarchy score”) based on features derived from
 44 representational similarity. We also apply this approach to analyze hidden layer representations in a
 45 database of 432K deep networks (NAS-Bench-101; [21]), and find a surprising degree of correlation
 46 between early layer and deep layer representations. One consequence of this is that we can predict a
 47 network’s test set accuracy using only embedded representations of early layers.

48 Overall, we provide a theoretical grounding which explains why existing similarity measures are
 49 useful (they are often close to metric spaces, and can be modified to fulfill metric space axioms
 50 precisely), draw new connections between these approaches and other established research areas [14,
 51 22], utilize this framework to propose novel metrics, and demonstrate a practical and general-purpose
 52 machine learning workflow that scales to datasets containing thousands of networks.

53 2 Methods

54 This section outlines a workflow (Fig. 1) to analyze representations across large collections of
 55 networks. After briefly summarizing prior approaches (sec. 2.1), we cover background material
 56 on metric spaces and discuss their theoretical advantages over existing dissimilarity measures (sec.
 57 2.2). We then present a class of metrics that capture these advantages (sec. 2.3) and cover a special
 58 case that is suited to convolutional layers (sec. 2.4). We conclude by developing even stronger
 59 geometric insights (sec. 2.6). First, we note that certain metrics correspond to geodesic distances on
 60 shape manifolds [23], enabling manifold optimization and learning algorithms to be applied to neural
 61 representations. Moreover, we can use multidimensional scaling methods [24] to find Euclidean
 62 embeddings that approximate this metric space, which we demonstrate empirically in Section 3.

63 2.1 Prior work and problem setup

64 Neural network representations are often summarized over a set of m reference inputs (e.g. test set
 65 images). Let $\mathbf{X}_i \in \mathbb{R}^{m \times n_i}$ and $\mathbf{X}_j \in \mathbb{R}^{m \times n_j}$ denote the responses of two networks (with n_i and n_j
 66 neurons, respectively) to a collection of these inputs. Quantifying the similarity between \mathbf{X}_i and \mathbf{X}_j
 67 is complicated by the fact that, while the m inputs are the same, there is no direct correspondence
 68 between the neurons. Even if $n_i = n_j$, the typical Frobenius inner product, $\langle \mathbf{X}_i, \mathbf{X}_j \rangle = \text{Tr}[\mathbf{X}_i^\top \mathbf{X}_j]$,
 69 and metric, $\|\mathbf{X}_i - \mathbf{X}_j\| = \langle \mathbf{X}_i - \mathbf{X}_j, \mathbf{X}_i - \mathbf{X}_j \rangle^{1/2}$, fail to capture the desired notion of dissimilarity.
 70 For instance, let $\mathbf{\Pi}$ denote some $n \times n$ permutation matrix and let $\mathbf{X}_i = \mathbf{X}_j \mathbf{\Pi}$. Intuitively, we
 71 should consider \mathbf{X}_i and \mathbf{X}_j to be identical in this case since the ordering of neurons is arbitrary. Yet,
 72 clearly $\|\mathbf{X}_i - \mathbf{X}_j\| \neq 0$, except in very special cases.

73 Representational similarity analysis (RSA; [6]) addresses this problem by computing the dissimilarity
 74 between network responses for each pair of inputs, resulting in an $m \times m$ representational dissimilarity

75 matrix (RDM) for each network. Then, RSA uses the Spearman rank correlation between the
 76 RDMs from two networks to assess their similarity. A simpler approach is to use linear regression
 77 to predict the responses of one network given the responses of another, using the coefficient of
 78 determination (R^2) as a measure of similarity [2, 3]. In machine learning, canonical correlation
 79 analysis (CCA; [4, 5]) and centered kernel alignment (CKA; [7, 8]) have been used to quantify the
 80 similarity of hidden layer activation patterns. These procedures allow for various linear and nonlinear
 81 transformations of representations, but without a theoretical framework it is unclear how to choose
 82 among them, use their outputs for downstream tasks, or generalize them to new domains.

83 2.2 Feature space mapping, metrics, and equivalence relations

84 Our first contribution will be to establish formal notions of distance (metrics) between neural
 85 representations. To accommodate the common scenario when the number of neurons varies across
 86 networks (i.e. when $n_i \neq n_j$), we must first map the representations into a common feature space so
 87 that a single distance function may be defined between them. For each set of representations, \mathbf{X}_i ,
 88 we suppose there is a mapping into a p -dimensional feature space, $\mathbf{X}_i \mapsto \mathbf{X}_i^\phi$, where $\mathbf{X}_i^\phi \in \mathbb{R}^{m \times p}$.
 89 In the special case where all networks have equal size, $n_1 = n_2 = \dots = n$, we can express
 90 the feature mapping as a single function $\phi : \mathbb{R}^{m \times n} \mapsto \mathbb{R}^{m \times p}$, so that $\mathbf{X}_i^\phi = \phi(\mathbf{X}_i)$. When
 91 networks have dissimilar sizes, we can reduce the representations down to some common dimension
 92 $p \leq \min(n_1, \dots, n_K)$ using, for example, PCA [4]. Alternatively, we can zero pad representations up
 93 to $p = \max(n_1, \dots, n_K)$ dimensions, which preserves the geometry of all network representations.

94 Next, we seek to establish *metrics* within the feature space, which are distance functions that satisfy:

$$\text{Equivalence: } d(\mathbf{X}_i^\phi, \mathbf{X}_j^\phi) = 0 \iff \mathbf{X}_i^\phi \sim \mathbf{X}_j^\phi \quad (1)$$

$$\text{Symmetry: } d(\mathbf{X}_i^\phi, \mathbf{X}_j^\phi) = d(\mathbf{X}_j^\phi, \mathbf{X}_i^\phi) \quad (2)$$

$$\text{Triangle Inequality: } d(\mathbf{X}_i^\phi, \mathbf{X}_j^\phi) \leq d(\mathbf{X}_i^\phi, \mathbf{X}_k^\phi) + d(\mathbf{X}_k^\phi, \mathbf{X}_j^\phi) \quad (3)$$

95 for all \mathbf{X}_i^ϕ , \mathbf{X}_j^ϕ , and \mathbf{X}_k^ϕ in the feature space. The symbol ‘ \sim ’ denotes an *equivalence relation*
 96 between two elements. That is, the expression $\mathbf{X}_i^\phi \sim \mathbf{X}_j^\phi$ means that “ \mathbf{X}_i^ϕ is equivalent to \mathbf{X}_j^ϕ .”
 97 Formally, distance functions satisfying Eqs. (1) to (3) define a metric over a quotient space defined by
 98 the equivalence relation (see Supplement A). Intuitively, by specifying different equivalence relations
 99 we can account for symmetries in network representations, such as permutations over arbitrarily
 100 labeled neurons (other options are discussed below in sec. 2.3).

101 Metrics quantify dissimilarity in a way that agrees with our intuitive notion of distance. For exam-
 102 ple, Eq. (2) ensures that the distance from \mathbf{X}_i^ϕ to \mathbf{X}_j^ϕ is the same as the distance from \mathbf{X}_j^ϕ to \mathbf{X}_i^ϕ .
 103 Linear regression is an approach that violates this condition—the similarity measured by R^2 depends
 104 on which network is treated as the dependent variable. Further, Eq. (3) ensures that distances are
 105 self-consistent in the sense that if two elements (\mathbf{X}_i^ϕ and \mathbf{X}_j^ϕ) are close to a third (\mathbf{X}_k^ϕ), then they
 106 are necessarily close to each other. Existing measures based on CCA, RSA, and CKA, are symmetric,
 107 but do not satisfy the triangle inequality. By modifying these approaches, we avoid potential pitfalls
 108 and can leverage theoretical guarantees on learning in proper metric spaces [15–19].

109 2.3 Generalized shape metrics and group invariance

110 In this section, we outline a new framework to quantify representational dissimilarity, which leverages
 111 a well-developed mathematical literature on *shape spaces* [12–14]. The key idea is to treat $\mathbf{X}_i^\phi \sim \mathbf{X}_j^\phi$
 112 if and only if there exists a linear transformation \mathbf{T} within a set of allowable transformations \mathcal{G} , such
 113 that $\mathbf{X}_i^\phi = \mathbf{X}_j^\phi \mathbf{T}$. Much of shape analysis literature focuses on the special case where \mathcal{G} is the
 114 special orthogonal group $\mathcal{SO}(p) = \{\mathbf{R} \in \mathbb{R}^{p \times p} \mid \mathbf{R}^\top \mathbf{R} = \mathbf{I}, \det(\mathbf{R}) = 1\}$, meaning that \mathbf{X}_i^ϕ and
 115 \mathbf{X}_j^ϕ are equivalent if there is a p -dimensional rotation (without reflection) that relates them. Standard
 116 shape analysis further considers each \mathbf{X}_i^ϕ to be a mean-centered ($(\mathbf{X}_i^\phi)^\top \mathbf{1} = \mathbf{0}$) and normalized
 117 ($\|\mathbf{X}_i^\phi\|_F = 1$) version of the raw landmark locations held in $\mathbf{X}_i \in \mathbb{R}^{m \times p}$ (an assumption that we
 118 will relax). Defining $\mathbb{S}^{m \times p}$ to be the hypersphere of $m \times p$ matrices with unit Frobenius norm, we
 119 see that $\mathbf{X}_i^\phi \in \mathbb{S}^{m \times p}$. In this context, \mathbf{X}_i^ϕ is called a “pre-shape.” By removing rotations from a
 120 pre-shape, $[\mathbf{X}_i^\phi] = \{\mathbf{S} \in \mathbb{S}^{m \times p} \mid \mathbf{S} \sim \mathbf{X}_i^\phi\}$ for pre-shape \mathbf{X}_i^ϕ , we recover its “shape.”

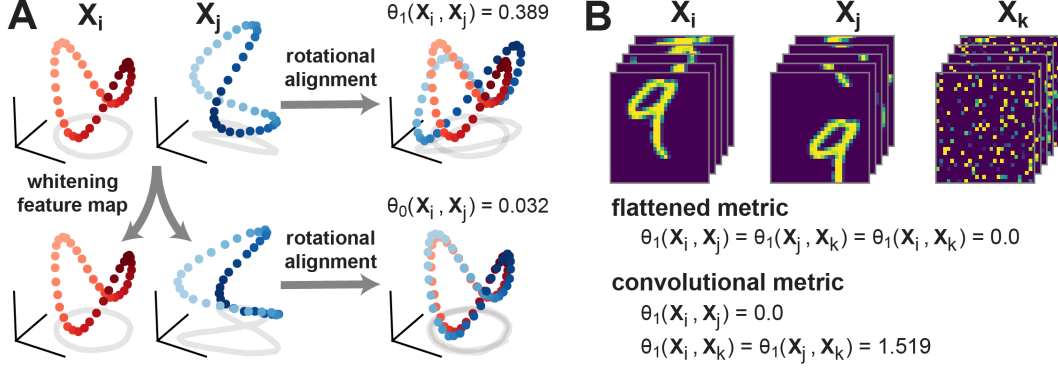


Figure 2: (A) Schematic illustration of metrics with rotational invariance (top), and linear invariance (bottom). Red and blue dots represent a pair of network representations X_i and X_j , which correspond to m points in n -dimensional space. (B) Demonstration of convolutional metric on toy data. Flattened metrics (e.g. [4, 7]) that ignore convolutional layer structure treat permuted images (X_k , right) as equivalent to images with coherent spatial structure (X_i and X_j , left and middle). A convolutional metric, Eq. (11), distinguishes between these cases while still treating X_i and X_j as equivalent (obeying translation invariance).

121 To quantify dissimilarity in network representations, we generalize this notion of shape to include
 122 other feature mappings (beyond mean-centering and normalization) and alignment operations (beyond
 123 rotations without reflection). The minimal distance within the feature space, after optimizing over
 124 alignments, defines a metric under suitable conditions (Fig. 2A). This results in a broad variety of
 125 *generalized shape metrics* (see also, ch. 18 of [14]), which fall into two categories as formalized by
 126 the pair of propositions below.

127 **Proposition 1.** Let $X_i^\phi \in \mathbb{R}^{m \times p}$, and let \mathcal{G} be a group of linear isometries on $\mathbb{R}^{m \times p}$. Then,

$$d(X_i^\phi, X_j^\phi) = \min_{T \in \mathcal{G}} \|X_i^\phi - X_j^\phi T\| \quad (4)$$

128 defines a metric, where $X_i^\phi \sim X_j^\phi$ if and only if there is a $T \in \mathcal{G}$ such that $X_i^\phi = X_j^\phi T$.

129 **Proposition 2.** Let $X_i^\phi \in \mathbb{S}^{m \times p}$, and let \mathcal{G} be a group of linear isometries on $\mathbb{S}^{m \times p}$. Then,

$$\theta(X_i^\phi, X_j^\phi) = \min_{T \in \mathcal{G}} \arccos \langle X_i^\phi, X_j^\phi T \rangle \quad (5)$$

130 defines a metric, where $X_i^\phi \sim X_j^\phi$ if and only if there is a $T \in \mathcal{G}$ such that $X_i^\phi = X_j^\phi T$.

131 Proofs are provided in Supplement B. The essential difference between d and θ is that the former uses
 132 the Euclidean metric to measure distance in the feature space, while the latter uses angular distance
 133 (geodesic distance on a sphere). The former typically enables easier computations, but the latter has
 134 appealing theoretical properties, including a close connection to CCA that we outline below.

135 Two key conditions appear in these propositions. First, \mathcal{G} must be a *group* of functions. This means \mathcal{G}
 136 is as a set that contains the identity function, is closed under composition ($T_1 T_2 \in \mathcal{G}$ for any $T_1 \in \mathcal{G}$
 137 and $T_2 \in \mathcal{G}$), and whose elements are invertible by other members of the set (if $T \in \mathcal{G}$ then $T^{-1} \in \mathcal{G}$).

138 Second, every $T \in \mathcal{G}$ must be an *isometry*, meaning that $\|X_i^\phi - X_j^\phi\| = \|X_i^\phi T - X_j^\phi T\|$ for all
 139 $T \in \mathcal{G}$ and all elements of the feature space. On $\mathbb{R}^{m \times p}$ and $\mathbb{S}^{m \times p}$, all linear isometries are orthogonal
 140 transformations. Further, the set of orthogonal transformations, $\mathcal{O}(p) = \{Q \in \mathbb{R}^{p \times p} : Q^\top Q = I\}$,
 141 defines a well-known group. Thus, the condition that \mathcal{G} is a group of isometries is equivalent to \mathcal{G}
 142 being a subgroup of $\mathcal{O}(p)$ —i.e., a subset of $\mathcal{O}(p)$ satisfying the group axioms.

143 Intuitively, the condition that \mathcal{G} is a group ensures that the alignment procedure is symmetric—i.e. it
 144 is equivalent to consider transforming X_i^ϕ to match X_j^ϕ , or vice versa. Further, the condition that
 145 each $T \in \mathcal{G}$ is an isometry ensures that the properties of the underlying metric space (Euclidean or
 146 angular distance) are preserved.

147 These propositions define a broad class of metrics as we enumerate below. For simplicity, we assume
 148 that $n_i = n_j = n$ in the examples below, with the understanding that a PCA or zero-padding
 149 preprocessing step has been performed in the case of dissimilar network sizes.

150 **Permutation invariance** The most stringent notion of representational similarity is to demand that
 151 neurons are one-to-one matched across networks. If we set $\mathbf{X}_i^\phi = \mathbf{X}_i$ for all i , then:

$$d_{\mathcal{P}}(\mathbf{X}_i, \mathbf{X}_j) = \min_{\mathbf{\Pi} \in \mathcal{P}} \|\mathbf{X}_i - \mathbf{X}_j \mathbf{\Pi}\| \quad (6)$$

152 defines a metric by Proposition 1, since the set of permutation matrices, $\mathcal{P}(n)$ is a subgroup of $\mathcal{O}(n)$.
 153 To evaluate this metric we must optimize over the set of neuron permutations to align the two
 154 networks. This can be reformulated (see Supplement C) as a linear assignment problem, which is a
 155 fundamental problem in combinatorial optimization (see [25] for review). Exploiting an algorithm
 156 due to Jonker and Volgenant [26, 27] we can solve this problem in $O(n^3)$ time. The overall runtime
 157 for evaluating Eq. (6) is $O(mn^2 + n^3)$, since we must evaluate $\mathbf{X}_i^\top \mathbf{X}_j$ to formulate the assignment
 158 problem.

159 **Rotation scaling invariance** Let $\mathbf{C} = \mathbf{I}_m - (1/m)\mathbf{1}\mathbf{1}^\top$ denote an $m \times m$ centering matrix, and
 160 consider the feature mapping which mean-centers the columns, $\mathbf{X}_i^\phi = \mathbf{C}\mathbf{X}_i$. Then,

$$d_1(\mathbf{X}_i^\phi, \mathbf{X}_j^\phi) = \min_{\mathbf{Q} \in \mathcal{O}} \|\mathbf{X}_i^{\phi_1} - \mathbf{X}_j^{\phi_1} \mathbf{Q}\| \quad (7)$$

161 defines a metric by Proposition 1, and is equivalent to the *Procrustes size-and-shape distance* with
 162 reflections [14]. Further,

$$\theta_1(\mathbf{X}_i^\phi, \mathbf{X}_j^\phi) = \min_{\mathbf{Q} \in \mathcal{O}} \arccos \frac{\langle \mathbf{X}_i^{\phi_1}, \mathbf{X}_j^{\phi_1} \mathbf{Q} \rangle}{\|\mathbf{X}_i^{\phi_1}\| \|\mathbf{X}_j^{\phi_1}\|} \quad (8)$$

163 defines a metric by Proposition 2, and is equivalent to the Riemannian distance on Kendall’s shape
 164 space [14]. To evaluate Eqs. (7) and (8), we must optimize over the set of orthogonal matrices to
 165 find the best alignment. This also maps onto a fundamental optimization problem known as the
 166 *orthogonal Procrustes problem* [28, 29], which can be solved in closed form in $O(n^3)$ time. As in
 167 the permutation-invariant metric described above, the overall runtime is $O(mn^2 + n^3)$.

168 **Linear invariance** Consider a partial whitening transformation, parameterized by $0 \leq \alpha \leq 1$:

$$\mathbf{X}^{\phi_\alpha} = \mathbf{C}\mathbf{X}(\alpha\mathbf{I} + (1 - \alpha)(\mathbf{X}^\top \mathbf{C}\mathbf{X})^{-1/2}) \quad (9)$$

169 Note that $\mathbf{X}^\top \mathbf{C}\mathbf{X}$ is the empirical covariance matrix of \mathbf{X} . Thus, when $\alpha = 0$, Eq. (9) corre-
 170 sponds to ZCA whitening [30], which intuitively removes invertible linear transformations from the
 171 representations. When $\alpha = 1$, Eq. (9) reduces to the mean-centering feature map used above.

172 This feature space leads to a metric that is related to CCA. First, let $1 \geq \rho_1 \geq \dots \geq \rho_n \geq 0$ denote
 173 the singular values of $(\mathbf{X}_i^{\phi_\alpha})^\top (\mathbf{X}_j^{\phi_\alpha}) / \|\mathbf{X}_i^{\phi_\alpha}\| \|\mathbf{X}_j^{\phi_\alpha}\|$. Using orthogonal alignments, we see:

$$\theta_\alpha(\mathbf{X}_i^\phi, \mathbf{X}_j^\phi) = \min_{\mathbf{Q} \in \mathcal{O}} \arccos \frac{\langle \mathbf{X}_i^{\phi_\alpha}, \mathbf{X}_j^{\phi_\alpha} \mathbf{Q} \rangle}{\|\mathbf{X}_i^{\phi_\alpha}\| \|\mathbf{X}_j^{\phi_\alpha}\|} = \arccos\left(\frac{1}{n} \sum_\ell \rho_\ell\right) \quad (10)$$

174 defines a metric by Proposition 2. The values ρ_1, \dots, ρ_n correspond to canonical correlation co-
 175 efficients when $\alpha = 0$ [31] and ridge regularized canonical correlations when $\alpha > 0$ [32] (see
 176 Supplement C for further notes). Past works have suggested to use the average canonical correlation
 177 as a similarity measure [4, 5]. This heuristic corresponds to $\cos \theta_0(\mathbf{X}_i, \mathbf{X}_j)$, but applying \arccos is
 178 needed to satisfy the triangle inequality. Since the covariance is often ill-conditioned or singular in
 179 practice, setting $\alpha > 0$ to regularize the calculation is typically necessary.

180 **Nonlinear invariances** We discuss feature maps that enable nonlinear notions of equivalence, and
 181 which relate to kernel CCA [33], in Supplement C.

182 2.4 Metrics for convolutional layers

183 In deep networks for image processing, each convolutional layer produces a $h \times w \times c$ array
 184 of activations, whose axes respectively correspond to image height, image width, and channels
 185 (number of convolutional filters). If stride-1 circular convolutions are used, then applying a circular

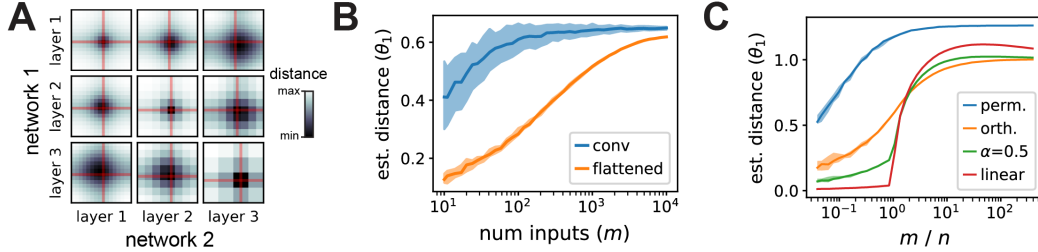


Figure 3: (A) Each heatmap shows a brute-force search over the shift parameters along the width and height dimensions of a pair of convolutional layers compared across two networks. The optimal shifts are typically close to zero (red lines). (B) Impact of sample size, m , on flattened and convolutional metrics with orthogonal invariance. The convolutional metric approaches its final value faster than the flattened metric, which is still increasing even at the full size of the CIFAR-10 test set ($m = 10^4$). (C) Impact of sample density, m/n , on metrics invariant to permutation, orthogonal, regularized linear ($\alpha = 0.5$), and linear transformations. Shaded regions mark the 10th and 90th percentiles across shuffled repeats. Further details for all simulations are provided in Supplement E.

186 shift along either spatial dimension produces the same shift in the layer’s output. It is natural to
 187 reflect this property, known as translation equivariance [22], in the equivalence relation on layer
 188 representations. Supposing that the feature map preserves the shape of the activation tensor, we have
 189 $\mathbf{X}_k^\phi \in \mathbb{R}^{m \times h \times w \times c}$ for $k \in 1, \dots, K$. Letting $\mathcal{S}(n)$ denote the group of n -dimensional circular shifts
 190 (a subgroup of the permutation group) and ‘ \otimes ’ denote the Kronecker product, we propose:

$$\mathbf{X}_i^\phi \sim \mathbf{X}_j^\phi \iff \text{vec}(\mathbf{X}_i^\phi) = (\mathbf{I} \otimes \mathbf{S}_1 \otimes \mathbf{S}_2 \otimes \mathbf{Q}) \text{vec}(\mathbf{X}_j^\phi) \quad (11)$$

191 for some $\mathbf{S}_1 \in \mathcal{S}(h)$, $\mathbf{S}_2 \in \mathcal{S}(w)$, $\mathbf{Q} \in \mathcal{O}(c)$, as the desired equivalence relation. This relation allows
 192 for orthogonal invariance across the channel dimension but only shift invariance across the spatial
 193 dimensions. The mixed product property of Kronecker products, $(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) = \mathbf{AB} \otimes \mathbf{CD}$,
 194 ensures that the overall transformation maintains the group structure and remains an isometry.
 195 Figure 2B uses a toy dataset (stacked MNIST digits) to show that this metric is sensitive to differences
 196 in spatial activation patterns, but insensitive to coherent spatial translations across channels. In
 197 contrast, metrics that ignore the convolutional structure (as in past work [4, 7]) treat very different
 198 spatial patterns as identical representations.

199 Evaluating Eq. (11) requires optimizing over spatial shifts in conjunction with solving a Procrustes
 200 alignment. If we fit the shifts by an exhaustive brute-force search, the overall runtime is $O(mh^2w^2c^2 +$
 201 $hwc^3)$, which is costly if this calculation is repeated across a large collection of networks. In practice,
 202 we observe that the optimal shift parameters are typically close to zero (Fig. 3A). This motivates the
 203 more stringent equivalence relation:

$$\mathbf{X}_i^\phi \sim \mathbf{X}_j^\phi \iff \text{vec}(\mathbf{X}_i^\phi) = (\mathbf{I} \otimes \mathbf{I} \otimes \mathbf{I} \otimes \mathbf{Q}) \text{vec}(\mathbf{X}_j^\phi) \quad \text{for some } \mathbf{Q} \in \mathcal{Q}, \quad (12)$$

204 which has a more manageable runtime of $O(mhwc^2 + c^3)$. To evaluate the metrics implied by
 205 Eq. (12), we can simply reshape each \mathbf{X}_k^ϕ from a $(m \times h \times w \times c)$ tensor into a $(mhw \times c)$ matrix
 206 and apply the Procrustes alignment procedure as done above for previous metrics. In contrast, the
 207 “flattened metric” in Fig. 2B reshapes the features into a $(m \times hwc)$ matrix, resulting in a more
 208 computationally expensive alignment that runs in $O(mh^2w^2c^2 + h^3w^3c^3)$ time.

209 2.5 How large of a sample size is needed?

210 An important issue, particularly in neurobiological applications, is to determine the number of
 211 network inputs, m , and neurons, n , that one needs to measure to accurately infer the distance between
 212 two network representations [11]. Reasoning about these questions rigorously requires a probabilistic
 213 perspective of shape distances, which we outline in Supplement D. Intuitively, looser equivalence
 214 relations involve a more flexible class of alignment operations, and thus require more sampled inputs
 215 to prevent overfitting. Figure 3B-C show that this intuition holds in practice for data from deep
 216 convolutional networks. Metrics with looser equivalence relations—the “flattened” metric in panel B,
 217 or e.g. the linear metric in panel C—converge slower to a stable estimate as m is increased.

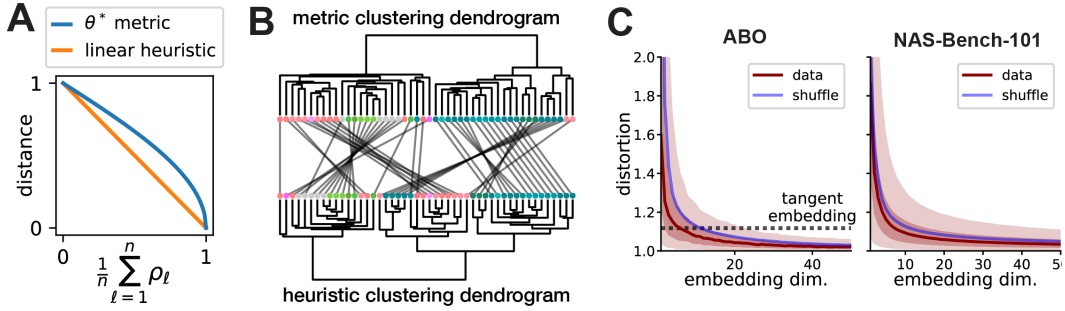


Figure 4: (A) Comparison of metric and linear heuristic. (B) Metric and linear heuristic produce discordant hierarchical clusterings of ABO dataset. Leaves represent brain areas that are clustered by representational similarity (see Fig. 1C), colored by Allen reference atlas, and ordered to maximize dendrogram similarities of adjacent leaves. In the middle, grey lines connect leaves corresponding to the same brain region across the two dendrograms. (C) ABO and NAS-Bench-101 datasets can be accurately embedded into Euclidean spaces. Dark red line shows median distortion. Light red shaded region corresponds to 5th to 95th percentiles of distortion, dark red shaded corresponds to interquartile range. The mean distortion of a null distribution over representations (blue line) was generated by shuffling the m inputs independently in each network.

218 2.6 Modeling approaches and conceptual insights

219 Generalized shape metrics facilitate several new approaches and perspectives. In many cases, the
 220 set of K neural representations can be visualized as K points on a smooth manifold (see Fig. 1).
 221 This holds rigorously due to the *quotient manifold theorem* [34] so long as \mathcal{G} is not a finite set (e.g.
 222 corresponding to permutation) and all matrices are full rank in the feature space. This geometric
 223 intuition can be made even stronger when \mathcal{G} corresponds to a connected manifold, such as $\mathcal{SO}(p)$. In
 224 this case, it can be shown that the geodesic distance between two neural representations coincides
 225 with the metrics we defined in Propositions 1 and 2 (see Supplement C, and [14]). This result extends
 226 to the well-documented manifold structure of *Kendall’s shape space* [23].

227 Viewing neural representations as points on a manifold is not a purely theoretical exercise—several
 228 statistical models can be adapted to manifold-valued data (e.g. principal geodesic analysis [35]
 229 provides a generalization of PCA), and additional adaptations are an area of active research [36].
 230 However, there is generally no simple connection between these curved geometries and the flat
 231 geometries of Euclidean or Hilbert spaces [37].¹ Unfortunately, the majority of off-the-shelf machine
 232 learning tools are incompatible with the former and require the latter. Thus, we can resort to a heuristic
 233 approach: the set of K representations can be embedded into a Euclidean space that approximately
 234 preserves the pairwise shape distances. One possibility, employed widely in shape analysis, is to
 235 embed points in the tangent space of the manifold at a reference point [40, 41]. Another approach,
 236 which we demonstrate below with favorable results, is to optimize the embedding directly [24, 42].

237 3 Applications and Results

238 We analyzed two large-scale public datasets spanning neuroscience (Allen Brain Observatory, ABO;
 239 Neuropixels - visual coding experiment; [20]) and deep learning (NAS-Bench-101; [21]). We
 240 constructed the ABO dataset by pooling recorded neurons from $K = 48$ anatomically defined brain
 241 regions across all sessions; each $\mathbf{X}_k \in \mathbb{R}^{m \times n}$ was a dimensionally reduced matrix holding the
 242 neural responses (summarized by $n = 100$ principal components) to $m = 1600$ movie frames (120
 243 second clip, “natural movie three”). The full NAS-Bench-101 dataset contains 423,624 architectures;
 244 however, we analyze a subset of $K = 2000$ networks for simplicity. In this application each $\mathbf{X}_k \in$
 245 $\mathbb{R}^{m \times n}$ is a representation from a specific network layer, with $(m, n) \in \{(32^2 \times 10^5, 128), (16^2 \times$
 246 $10^5, 256), (8^2 \times 10^5, 512), (10^5, 512)\}$. Here, n corresponds to the number of channels and m is
 247 the product of the number of test set images (10^5) and the height and width dimensions of the
 248 convolutional layer—i.e., we use equivalence relation in Eq. (12) to evaluate dissimilarity.

¹However, see [38] for a conjectured relationship and [39] for a result in the special case of 2D shapes.

249 **Triangle inequality violations can occur in practice when using existing methods.** As mentioned
 250 above, a dissimilarity measure based on the mean canonical correlation, $1 - \sum_{\ell} \rho_{\ell}/n$, has been used
 251 in past work [5, 9]. We refer to this as the “linear heuristic.” A slight reformulation of this calculation,
 252 $\arccos(\sum_{\ell} \rho_{\ell}/n)$, produces a metric that satisfies the triangle inequality (see Eq. (10)). Figure 4A
 253 compares these calculations as a function of the average (regularized) canonical correlation: one can
 254 see that $\arccos(\cdot)$ is approximately linear when the mean correlation is near zero, but highly nonlinear
 255 when the mean correlation is near one. Thus, we reasoned that triangle inequality violations are more
 256 likely to occur when K is large and when many network representations are close to each other. Both
 257 ABO and NAS-Bench-101 datasets satisfy these conditions, and in both cases we observed triangle
 258 inequality violations by the linear heuristic with full regularization ($\alpha = 1$): 17/1128 network pairs in
 259 the ABO dataset had at least one triangle inequality violation, while 10128/100000 randomly sampled
 260 network pairs contained violations in the NAS-Bench-101 Stem layer dataset. We also examined a
 261 standard version of RSA that quantifies similarity via Spearman’s rank correlation coefficient [6].
 262 Similar to the results above, we observed violations in 14/1128 pairs of networks in the ABO dataset.

263 Overall, these results suggest that generalized shape metrics correct for triangle inequality violations
 264 that do occur in practice. Depending on the dataset, these violations may be rare ($\sim 1\%$ occurrence
 265 in ABO) or relatively common ($\sim 10\%$ in the Stem layer of NAS-Bench-101). While we do not
 266 expect these violations to invalidate the qualitative conclusions of past work, these differences can
 267 produce quantitative discrepancies in downstream analyses. For example, the dendrograms produced
 268 by hierarchical clustering differ depending on whether one uses the linear heuristic or the shape
 269 distance ($\sim 85.1\%$ dendrogram similarity as quantified by the method in [43]; see Fig. 4B).

270 **Neural representation metric spaces can be approximated by Euclidean spaces.** Having estab-
 271 lished that neural representations can be viewed as elements in a metric space, it is natural to ask if this
 272 metric space is, loosely speaking, “close to” a Euclidean space. We used standard multidimensional
 273 scaling methods (SMACOF, [24]; implementation in [44]) to obtain a set of embedded vectors,
 274 $\mathbf{y}_i \in \mathbb{R}^L$, for which $\theta_1(\mathbf{X}_i^{\phi}, \mathbf{X}_j^{\phi}) \approx \|\mathbf{y}_i - \mathbf{y}_j\|$ for $i, j \in 1, \dots, K$. The embedding dimension L is
 275 a freely chosen hyperparameter. This embedding problem admits multiple formulations and optimiza-
 276 tion strategies [42], which could be systematically explored in future work. Our simple approach
 277 already yields promising results: in both datasets, we find that moderate embedding dimensions
 278 ($L \approx 20$) is sufficient to produce high-quality embeddings. Following past work [45], we quantify
 279 the embedding distortions multiplicatively as:

$$\max(\theta_1(\mathbf{X}_i^{\phi}, \mathbf{X}_j^{\phi})/\|\mathbf{y}_i - \mathbf{y}_j\|; \|\mathbf{y}_i - \mathbf{y}_j\|/\theta_1(\mathbf{X}_i^{\phi}, \mathbf{X}_j^{\phi})) \quad (13)$$

280 for each pair of networks $i, j \in 1, \dots, K$. Plotting the distortions as a function of L (Fig. 4C), we see
 281 that they rapidly decrease, such that 95% of pairwise distances are distorted by, at most, $\sim 5\%$ (ABO
 282 data) or 10% (NAS-Bench-101) for sufficiently large L . Past work [9] has used multidimensional
 283 scaling heuristically to visualize collections of network representations in $L = 2$ dimensions. Our
 284 results here suggest that such a small value of L , while being amenable to visualization, results in a
 285 highly distorted embedding. It is noteworthy that the situation improves dramatically when L is even
 286 modestly increased. While we cannot visualize these higher-dimensional vector embeddings, we can
 287 use them as features for downstream modeling tasks. This is well-motivated as an approximation to
 288 performing model inference in the true metric space that characterizes neural representations [45].

289 **Anatomical structure and hierarchy is reflected in ABO representations.** We can now collect
 290 the L -dimensional vector embeddings of K network representations into a matrix $\mathbf{Z} \in \mathbb{R}^{K \times L}$. We
 291 expect the distance between any two rows, $\|\mathbf{z}_i - \mathbf{z}_j\|$, to closely reflect the distance between network
 292 representations i and j in shape space. We applied PCA to \mathbf{Z} to visualize the $K = 48$ brain regions
 293 and found that anatomically related brain regions indeed were closer together in the embedded space
 294 (Fig. 5A): cortical and sub-cortical regions are separated along PC 1, and different layers of the same
 295 region (e.g. layers 2/3, 4, 5, and 6a of VISp) are clustered together. As expected from the distortions
 296 shown in Fig. 4C, performing multidimensional scaling directly ($L = 2$, as done in [9]) results in
 297 a qualitatively different outcome (see Supplement E). Additionally, we used \mathbf{Z} to fit an ensembled
 298 kernel regressor (see Supplement E) to predict an anatomical hierarchy score (defined in [46]) from
 299 the embedded vectors. Overall, these results demonstrate that the geometry of the learned embedding
 300 is scientifically interpretable and can be exploited for novel analyses, such as nonlinear regression.
 301 To our best knowledge, the fine scale anatomical parcellation used here is novel in the context of
 302 representational similarity studies. The heightened scale and complexity of these data motivate the
 303 use of methods grounded in metric spaces that are amenable to detailed geometric interpretations.

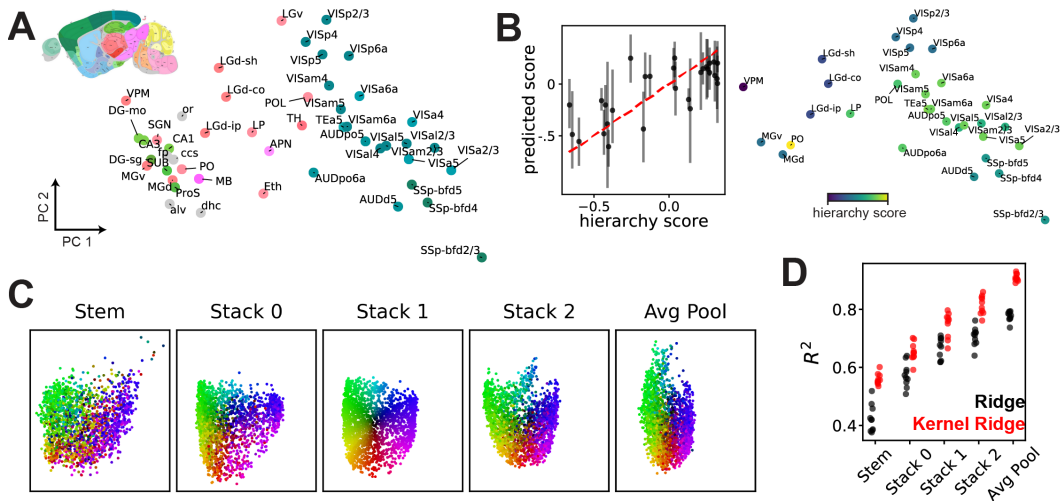


Figure 5: (A) PCA visualization of representations across 48 brain regions in the ABO dataset. Areas are colored by the reference atlas (see inset), illustrating a functional clustering of regions that maps onto anatomy. (B) *Left*, kernel regression predicts anatomical hierarchy [46] from embedded representations (see Supplement E). *Right*, PCA visualization of 31 areas labeled with hierarchy scores. (C) PCA visualization of 2000 network representations (a subset of NAS-Bench-101) across five layers, showing global structure is preserved across layers. Each network is colored by its position in the “Stack 1” layer (the middle of the architecture). (D) Embeddings of NAS-Bench-101 representations are predictive of test set accuracy, *even in very early layers*.

304 **NAS-Bench-101 representations show persistent structure across layers.** Since we collected
 305 representations across five layers in each deep network, the embedded representation vectors form
 306 a set of five $K \times L$ matrices, $\{Z_1, Z_2, Z_3, Z_4, Z_5\}$. We aligned these embeddings by rotations in
 307 \mathbb{R}^L , and then performed PCA to visualize the $K = 2000$ network representations from each layer
 308 in a common low-dimensional space. We observe that many features of the global structure are
 309 remarkably well-preserved—two networks that are close together in the Stack1 layer are assigned
 310 similar colors in Fig. 5C, and are likely to be close together in the other four layers. This preservation
 311 of representational similarity across layers suggests that even early layers contain signatures of
 312 network performance, which we expect to be present in the AvgPool layer. Indeed, when we fit ridge
 313 and RBF kernel ridge regressors to predict test set accuracy from representation embeddings, we see
 314 that even early layers support moderately good predictions (Fig. 5D). This is particularly surprising
 315 for the Stem layer. This is the first layer in each network, and its architecture is identical for all
 316 networks. Thus, the differences that are detected in the Stem layer result only from differences in
 317 backpropagated gradients. Again, these results demonstrate the ability of generalized shape metrics
 318 to incorporate neural representations into analyses with greater scale (K corresponding to thousands
 319 of networks) and complexity (nonlinear kernel regression) than has been previously executed.

320 4 Conclusion and Limitations

321 We demonstrated how to scale investigations into neural representations to large collections of
 322 networks by grounding analyses in proper metric spaces. An important limitation of our work, as
 323 well as the past works we build upon, is the possibility that representational geometry may only be
 324 loosely tied to the underlying mechanisms of network function [9]. On the other hand, analyses of
 325 representational geometry provide intermediate-level insights into network function [47]. Further,
 326 as we showed, these analyses can be systematically and quantitatively studied across thousands of
 327 networks—a much larger scale than is typically considered. Our work exploited connections to previ-
 328 ously unrecognized shape analysis literature [12–14], while suggesting only modest modifications to
 329 existing similarity measures like CCA (see Fig. 4A). Further, several of the metrics we utilize can
 330 be viewed as geodesic distances on Riemannian manifolds [23]. Future work would ideally exploit
 331 methods that are rigorously adapted to such manifolds, which are being actively developed [36].
 332 Nonetheless, we found that optimized Euclidean embeddings, while only approximate, provide a
 333 practical off-the-shelf solution for large-scale surveys of neural representations.

334 **References**

- 335 [1] David GT Barrett, Ari S Morcos, and Jakob H Macke. “Analyzing biological and artificial
336 neural networks: challenges with opportunities for synergy?” *Current Opinion in Neurobiology*
337 55 (2019). Machine Learning, Big Data, and Neuroscience, pp. 55–64.
- 338 [2] Daniel L. K. Yamins, Ha Hong, Charles F. Cadieu, Ethan A. Solomon, Darren Seibert, and
339 James J. DiCarlo. “Performance-optimized hierarchical models predict neural responses
340 in higher visual cortex”. *Proceedings of the National Academy of Sciences* 111.23 (2014),
341 pp. 8619–8624.
- 342 [3] Santiago A. Cadena, Fabian H. Sinz, Taliah Muhammad, Emmanouil Froudarakis, Erick Cobos,
343 Edgar Y. Walker, Jake Reimer, Matthias Bethge, Andreas Tolias, and Alexander S. Ecker.
344 “How well do deep neural networks trained on object recognition characterize the mouse visual
345 system?” *NeurIPS Workshop Neuro AI* (2019).
- 346 [4] Maithra Raghu, Justin Gilmer, Jason Yosinski, and Jascha Sohl-Dickstein. “SVCCA: Singular
347 Vector Canonical Correlation Analysis for Deep Learning Dynamics and Interpretability”.
348 *Advances in Neural Information Processing Systems 30*. Ed. by I. Guyon, U. V. Luxburg,
349 S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett. Curran Associates, Inc.,
350 2017, pp. 6076–6085.
- 351 [5] Ari Morcos, Maithra Raghu, and Samy Bengio. “Insights on representational similarity in
352 neural networks with canonical correlation”. *Advances in Neural Information Processing*
353 *Systems 31*. Ed. by S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and
354 R. Garnett. Curran Associates, Inc., 2018, pp. 5727–5736.
- 355 [6] Nikolaus Kriegeskorte, Marieke Mur, and Peter Bandettini. “Representational similarity analy-
356 sis - connecting the branches of systems neuroscience”. *Frontiers in Systems Neuroscience* 2
357 (2008), p. 4.
- 358 [7] Simon Kornblith, Mohammad Norouzi, Honglak Lee, and Geoffrey Hinton. “Similarity of
359 Neural Network Representations Revisited”. Ed. by Kamalika Chaudhuri and Ruslan Salakhut-
360 dinov. Vol. 97. Proceedings of Machine Learning Research. Long Beach, California, USA:
361 PMLR, 2019, pp. 3519–3529.
- 362 [8] Corinna Cortes, Mehryar Mohri, and Afshin Rostamizadeh. “Algorithms for learning kernels
363 based on centered alignment”. *The Journal of Machine Learning Research* 13.1 (2012),
364 pp. 795–828.
- 365 [9] Niru Maheswaranathan, Alex Williams, Matthew Golub, Surya Ganguli, and David Sus-
366 sillo. “Universality and individuality in neural dynamics across large populations of recurrent
367 networks”. *Advances in Neural Information Processing Systems 32*. Ed. by H. Wallach, H.
368 Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett. Curran Associates, Inc.,
369 2019, pp. 15629–15641.
- 370 [10] Thao Nguyen, Maithra Raghu, and Simon Kornblith. *Do Wide and Deep Networks Learn the*
371 *Same Things? Uncovering How Neural Network Representations Vary with Width and Depth*.
372 2020.
- 373 [11] Jianghong Shi, Eric Shea-Brown, and Michael Buice. “Comparison Against Task Driven
374 Artificial Neural Networks Reveals Functional Properties in Mouse Visual Cortex”. *Advances in*
375 *Neural Information Processing Systems 32*. Ed. by H. Wallach, H. Larochelle, A. Beygelzimer,
376 F. d’Alché-Buc, E. Fox, and R. Garnett. Curran Associates, Inc., 2019, pp. 5764–5774.
- 377 [12] Christopher G. Small. *The statistical theory of shape*. Springer series in statistics. New York:
378 Springer, 1996.
- 379 [13] David George Kendall, Dennis Barden, Thomas K Carne, and Huiling Le. *Shape and shape*
380 *theory*. New York: Wiley, 1999.
- 381 [14] Ian L. Dryden and Kantilal Mardia. *Statistical shape analysis with applications in R*. Chichester,
382 UK Hoboken, NJ: John Wiley & Sons, 2016.
- 383 [15] Peter N Yianilos. “Data structures and algorithms for nearest neighbor search in general metric
384 spaces”. *Proceedings of the fourth annual ACM-SIAM Symposium on Discrete algorithms*.
385 1993, pp. 311–321.
- 386 [16] Sanjoy Dasgupta and Philip M Long. “Performance guarantees for hierarchical clustering”.
387 *Journal of Computer and System Sciences* 70.4 (2005), pp. 555–569.

- 388 [17] Saaid Baraty, Dan A. Simovici, and Catalin Zara. “The Impact of Triangular Inequality
389 Violations on Medoid-Based Clustering”. *Foundations of Intelligent Systems*. Ed. by Marzena
390 Kryszkiewicz, Henryk Rybinski, Andrzej Skowron, and Zbigniew W. Raś. Berlin, Heidelberg:
391 Springer Berlin Heidelberg, 2011, pp. 280–289.
- 392 [18] Fei Wang and Jimeng Sun. “Survey on distance metric learning and dimensionality reduction
393 in data mining”. *Data Mining and Knowledge Discovery* 29.2 (2015), pp. 534–564.
- 394 [19] C. Chang, W. Liao, Y. Chen, and L. Liou. “A Mathematical Theory for Clustering in Metric
395 Spaces”. *IEEE Transactions on Network Science and Engineering* 3.1 (2016), pp. 2–16.
- 396 [20] Joshua H. Siegle et al. “Survey of spiking in the mouse visual system reveals functional
397 hierarchy”. *Nature* 592.7852 (2021), pp. 86–92.
- 398 [21] Chris Ying, Aaron Klein, Eric Christiansen, Esteban Real, Kevin Murphy, and Frank Hutter.
399 “NAS-Bench-101: Towards Reproducible Neural Architecture Search”. *Proceedings of the*
400 *36th International Conference on Machine Learning*. Ed. by Kamalika Chaudhuri and Ruslan
401 Salakhutdinov. Vol. 97. Proceedings of Machine Learning Research. PMLR, 2019, pp. 7105–
402 7114.
- 403 [22] Taco Cohen and Max Welling. “Group Equivariant Convolutional Networks”. *Proceedings of*
404 *The 33rd International Conference on Machine Learning*. Ed. by Maria Florina Balcan and
405 Kilian Q. Weinberger. Vol. 48. Proceedings of Machine Learning Research. New York, New
406 York, USA: PMLR, 2016, pp. 2990–2999.
- 407 [23] David G. Kendall. “Shape Manifolds, Procrustean Metrics, and Complex Projective Spaces”.
408 *Bulletin of the London Mathematical Society* 16.2 (1984), pp. 81–121.
- 409 [24] Ingwer Borg and Patrick JF Groenen. *Modern multidimensional scaling: Theory and applica-*
410 *tions*. Springer Science & Business Media, 2005.
- 411 [25] Rainer Burkard, Mauro Dell’Amico, and Silvano Martello. *Assignment Problems*. Society for
412 Industrial and Applied Mathematics, 2012.
- 413 [26] R Jonker and A Volgenant. “A shortest augmenting path algorithm for dense and sparse linear
414 assignment problems”. *Computing* 38.4 (1987), pp. 325–340.
- 415 [27] David F. Crouse. “On implementing 2D rectangular assignment algorithms”. *IEEE Transac-*
416 *tions on Aerospace and Electronic Systems* 52.4 (2016), pp. 1679–1696.
- 417 [28] Peter H. Schönemann. “A generalized solution of the orthogonal procrustes problem”. *Psy-*
418 *chometrika* 31.1 (1966), pp. 1–10.
- 419 [29] J. C. Gower and Garnt B. Dijkstra. *Procrustes problems*. Oxford New York: Oxford
420 University Press, 2004.
- 421 [30] Agnan Kessy, Alex Lewin, and Korbinian Strimmer. “Optimal whitening and decorrelation”.
422 *The American Statistician* 72.4 (2018), pp. 309–314.
- 423 [31] Harold Hotelling. “Relations Between Two Sets of Variates”. *Biometrika* 28.3/4 (1936),
424 pp. 321–377.
- 425 [32] Hrishikesh D Vinod. “Canonical ridge and econometrics of joint production”. *Journal of*
426 *econometrics* 4.2 (1976), pp. 147–166.
- 427 [33] P. L. Lai and C. Fyfe. “Kernel and Nonlinear Canonical Correlation Analysis”. *International*
428 *Journal of Neural Systems* 10.05 (2000). PMID: 11195936, pp. 365–377.
- 429 [34] John M. Lee. *Introduction to smooth manifolds*. 2nd ed. Graduate texts in mathematics 218.
430 New York ; London: Springer, 2013.
- 431 [35] P. Thomas Fletcher and Sarang Joshi. “Riemannian geometry for the statistical analysis of
432 diffusion tensor data”. *Signal Processing* 87.2 (2007). Tensor Signal Processing, pp. 250–262.
- 433 [36] Nina Miolane, Nicolas Guigui, Alice Le Brigant, Johan Mathe, Benjamin Hou, Yann Thanwer-
434 das, Stefan Heyder, Olivier Peltre, Niklas Koep, Hadi Zaatiti, Hatem Hajri, Yann Cabanes,
435 Thomas Gerald, Paul Chauchat, Christian Shewmake, Daniel Brooks, Bernhard Kainz, Claire
436 Donnat, Susan Holmes, and Xavier Pennec. “Geomstats: A Python Package for Riemannian
437 Geometry in Machine Learning”. *Journal of Machine Learning Research* 21.223 (2020),
438 pp. 1–9.
- 439 [37] Aasa Feragen, Francois Lauze, and Soren Hauberg. “Geodesic Exponential Kernels: When
440 Curvature and Linearity Conflict”. *Proceedings of the IEEE Conference on Computer Vision*
441 *and Pattern Recognition (CVPR)*. 2015.

- 442 [38] Aasa Feragen and Søren Hauberg. “Open Problem: Kernel methods on manifolds and metric
443 spaces. What is the probability of a positive definite geodesic exponential kernel?” *29th Annual
444 Conference on Learning Theory*. Ed. by Vitaly Feldman, Alexander Rakhlin, and Ohad Shamir.
445 Vol. 49. Proceedings of Machine Learning Research. Columbia University, New York, New
446 York, USA: PMLR, 2016, pp. 1647–1650.
- 447 [39] Sadeep Jayasumana, Mathieu Salzmann, Hongdong Li, and Mehrtaash Harandi. “A Framework
448 for Shape Analysis via Hilbert Space Embedding”. *Proceedings of the IEEE International
449 Conference on Computer Vision (ICCV)*. 2013.
- 450 [40] Ian L Dryden and Kanti V Mardia. “Multivariate shape analysis”. *Sankhyā: The Indian Journal
451 of Statistics, Series A* (1993), pp. 460–480.
- 452 [41] F. James Rohlf. “Shape Statistics: Procrustes Superimpositions and Tangent Spaces”. *Journal
453 of Classification* 16.2 (1999), pp. 197–223.
- 454 [42] Akshay Agrawal, Alnur Ali, and Stephen Boyd. “Minimum-Distortion Embedding”. *arXiv
455* (2021).
- 456 [43] Alexander J. Gates, Ian B. Wood, William P. Hetrick, and Yong-Yeol Ahn. “Element-centric
457 clustering comparison unifies overlaps and hierarchy”. *Scientific Reports* 9.1 (2019), p. 8574.
- 458 [44] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P.
459 Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher,
460 M. Perrot, and E. Duchesnay. “Scikit-learn: Machine Learning in Python”. *Journal of Machine
461 Learning Research* 12 (2011), pp. 2825–2830.
- 462 [45] Leena Chennuru Vankadara and Ulrike von Luxburg. “Measures of distortion for machine
463 learning”. *Advances in Neural Information Processing Systems*. Ed. by S. Bengio, H. Wallach,
464 H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett. Vol. 31. Curran Associates, Inc.,
465 2018.
- 466 [46] Julie A. Harris et al. “Hierarchical organization of cortical and thalamic connectivity”. *Nature
467* 575.7781 (2019), pp. 195–202.
- 468 [47] Jess B. Hamrick and Shakir Mohamed. “Levels of Analysis for Machine Learning”. *Proceed-
469 ings of the ICLR 2020 Workshop on Bridging AI and Cognitive Science*. 2020.

470 Checklist

- 471 1. For all authors...
- 472 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
473 contributions and scope? [Yes]
- 474 (b) Did you describe the limitations of your work? [Yes] See final section on conclusions
475 and limitations
- 476 (c) Did you discuss any potential negative societal impacts of your work? [No] We cannot
477 foresee any negative impacts.
- 478 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
479 them? [Yes]
- 480 2. If you are including theoretical results...
- 481 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 482 (b) Did you include complete proofs of all theoretical results? [Yes] They are included in
483 the supplement.
- 484 3. If you ran experiments...
- 485 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
486 mental results (either in the supplemental material or as a URL)? [Yes]
- 487 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
488 were chosen)? [Yes] They are provided in the supplement. The vast majority of models
489 were pre-trained, with details provided in [21].
- 490 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
491 ments multiple times)? [Yes] See, e.g., Fig. 5B
- 492 (d) Did you include the total amount of compute and the type of resources used (e.g., type
493 of GPUs, internal cluster, or cloud provider)? [Yes] See Supplement E.

- 494 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 495 (a) If your work uses existing assets, did you cite the creators? [Yes] See, e.g., our citations
- 496 to [20, 21, 44]
- 497 (b) Did you mention the license of the assets? [N/A] We do not believe this is needed for
- 498 our use cases.
- 499 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
- 500 We have included a small code package in the supplement
- 501 (d) Did you discuss whether and how consent was obtained from people whose data you're
- 502 using/curating? [N/A] No assets involve data collected from human subjects and the
- 503 assets are all freely shared
- 504 (e) Did you discuss whether the data you are using/curating contains personally identifiable
- 505 information or offensive content? [Yes]
- 506 5. If you used crowdsourcing or conducted research with human subjects...
- 507 (a) Did you include the full text of instructions given to participants and screenshots, if
- 508 applicable? [N/A]
- 509 (b) Did you describe any potential participant risks, with links to Institutional Review
- 510 Board (IRB) approvals, if applicable? [N/A]
- 511 (c) Did you include the estimated hourly wage paid to participants and the total amount
- 512 spent on participant compensation? [N/A]