

---

# Estimating shape distances on neural representations with limited samples

---

Dean A. Pospisil<sup>1</sup> Brett W. Larsen<sup>2,3,4</sup> Sarah E. Harvey<sup>2,3</sup> Alex H. Williams<sup>2,3</sup>

<sup>1</sup>Princeton University, Princeton, NJ, 08544; dp4846@princeton.edu

<sup>2</sup>New York University, Center for Neural Science, New York, NY, 10003

<sup>3</sup>Flatiron Institute, Center for Computational Neuroscience, New York, NY, 10010

<sup>4</sup>Flatiron Institute, Center for Computational Mathematics, New York, NY, 10010  
{brettlarsen, sharvey, awilliams}@flatironinstitute.org

## Abstract

Measuring geometric similarity between high-dimensional network representations is a topic of longstanding interest to neuroscience and deep learning. Although many methods have been proposed, only a few works have rigorously analyzed their statistical efficiency or quantified estimator uncertainty in data-limited regimes. Here, we derive upper and lower bounds on the worst-case convergence of standard estimators of *shape distance*—a measure of representational dissimilarity proposed by Williams et al. [30]. These bounds reveal the challenging nature of the problem in high-dimensional feature spaces. To overcome these challenges, we introduce a new method-of-moments estimator with a tunable bias-variance tradeoff. We show that this estimator achieves superior performance to standard estimators, particularly in high-dimensional settings. Thus, we lay the foundation for a rigorous statistical theory for high-dimensional shape analysis, and we contribute a new estimation method well-suited to practical scientific settings.

## 1 Introduction

Many approaches have been proposed to quantify similarity in neural network representations. Some popular methods include canonical correlations analysis [21], centered kernel alignment [CKA; 13], representational similarity analysis [RSA; 14], and shape metrics [30]. Each of these approaches takes in a set of high-dimensional measurements from two networks—e.g., hidden layer activations or measured biological responses—and outputs a (dis)similarity score. Shape distances additionally satisfy the triangle inequality, thus enabling downstream algorithms for clustering and nearest-neighbor regression that leverage metric space structure [30]. These measures have numerous applications including comparisons of artificial and biological systems [10, 24], comparisons of neural activity across different animal species [15], quantifying how hidden layer activity differs across deep network architectures [18, 19], and many more [see 11, for review]

In many practical settings, these measures must be estimated over a finite set of sampled networks inputs. However, with the noteworthy exception of research on RSA [4, 22, 28], there is little work on quantifying uncertainty (e.g. through confidence intervals) on estimators of representational similarity. This poses a serious obstacle to adoption of these methods, particularly in experimental neuroscience where there is a hard limit on the number of conditions that can be feasibly sampled [23, 29].

We address these concerns in the context of measuring shape distances between neural representations [Fig. 1; 30]. First we obtain analytic upper and lower bounds on the accuracy of typical “plug-in estimates” of shape distance as a function of the number of samples,  $M$ , and the dimension of the representation,  $N$ . We then propose a new method-of-moments estimator with an explicit and tunable tradeoff between estimator bias and variance to overcome the limitations of the plug-in estimator.

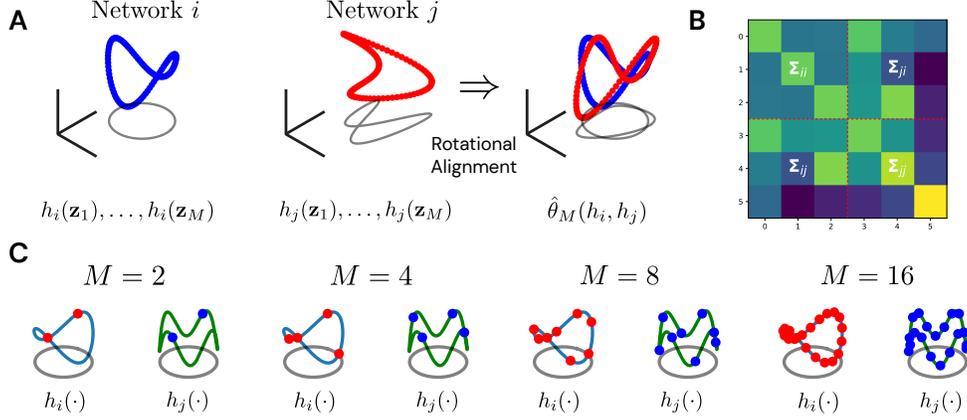


Figure 1: **(A)** Classical shape distances [9] can be used to provide a rotation-invariant distance between neural representations [30]. Given two labelled points clouds in  $N$ -dimensional space (*left* and *middle*), the distance is computed after an optimal orthogonal transformation is chosen to align the point clouds (*right*). In this visual example the point clouds trace out a low-dimensional manifold. **(B)** Heatmap shows the covariances ( $\Sigma_{ii}$ ,  $\Sigma_{jj}$ ) and cross-covariance ( $\Sigma_{ij}$ ) of the 3D representations in panel A. Shape distances can be re-expressed in terms of these quantities (see eq. 2.5, 2.6). **(C)** Our ability to estimate the shape distance is related to  $M$ , the number of stimuli. As  $M$  increases (*left* to *right*) the number of sampled points along the underlying manifold increases, and we are better able to resolve shape differences between the representations.

## 2 Results

We begin by reviewing generalized shape distances and the standard plug-in estimator (extended background can be found in App. A) Based on our theoretical characterization of the plug-in estimator in App. B, we find that plug-in estimates rapidly converge onto their expected value, but the expected error decays moderately slowly (i.e. the estimators have low variance and high bias). We thus introduce a method-of-moments estimator with tunable bias (Sec. 2.2) to overcome these shortcomings. We characterize the behavior of both estimators on synthetic (Sec. 2.3) and neural data (App. D.2). Finally, we discuss the implications of our results in Sec. 3.

### 2.1 Problem Setting

We consider initially a simple setting where each neural network is a deterministic map (for the stochastic setting, see appendix A.2). A collection of  $K$  neural systems can then be viewed as a set of functions, each denoted  $h_i : \mathcal{Z} \mapsto \mathbb{R}^N$  for  $i \in \{1, \dots, K\}$ . Here,  $\mathcal{Z}$  is a feature space and  $N$  can be interpreted as the number of neurons in each system (e.g. the size of a hidden layer in an artificial network, or the number of recorded neurons in a biological experiment).

Motivated by the shape theory literature [9, 30], we consider estimating the **Procrustes size-and-shape distance**,  $\rho$ , and **Riemannian shape distance**,  $\theta$ , between neural representations. Let  $h_i$  and  $h_j$  denote neural systems that are mean-centered and bounded:

$$\mathbb{E}[h_i(\mathbf{z})] = \mathbb{E}[h_j(\mathbf{z})] = \mathbf{0} \quad \text{and} \quad \|h_i(\mathbf{z})\|_2, \|h_j(\mathbf{z})\|_2 < B\sqrt{N} \quad \text{almost surely.} \quad (2.1)$$

for some constant  $B > 0$ . Here, the expectations are taken over  $\mathbf{z} \sim P$ , for some distribution  $P$  over network inputs. The Procrustes and Riemannian shape distances can be defined [App. D in 30]:

$$\rho(h_i, h_j) = \min_{\mathbf{Q} \in \mathcal{O}(N)} \sqrt{\mathbb{E} \|h_i(\mathbf{z}) - \mathbf{Q}h_j(\mathbf{z})\|_2^2} \quad (2.2)$$

$$\theta(h_i, h_j) = \min_{\mathbf{Q} \in \mathcal{O}(N)} \cos^{-1} \left( \frac{\mathbb{E}[h_i(\mathbf{z})^\top \mathbf{Q}h_j(\mathbf{z})]}{\sqrt{\mathbb{E}[h_i(\mathbf{z})^\top h_i(\mathbf{z})] \mathbb{E}[h_j(\mathbf{z})^\top h_j(\mathbf{z})]}} \right) \quad (2.3)$$

where  $\mathcal{O}(N)$  denotes the set of  $N \times N$  orthogonal matrices.

It is well-known that the optimal orthogonal alignment in eqs. (2.2) and (2.3) can be identified in closed form (App. A). Leveraging this, we can use the covariance and cross-covariance matrices,

$$\mathbf{\Sigma}_{ii} = \mathbb{E}[h_i(\mathbf{z})h_i(\mathbf{z})^\top], \quad \mathbf{\Sigma}_{jj} = \mathbb{E}[h_j(\mathbf{z})h_j(\mathbf{z})^\top], \quad \mathbf{\Sigma}_{ij} = \mathbb{E}[h_i(\mathbf{z})h_j(\mathbf{z})^\top], \quad (2.4)$$

to reformulate the squared Procrustes distance and cosine shape similarity:

$$\rho^2(h_i, h_j) = \text{Tr}[\mathbf{\Sigma}_{ii}] + \text{Tr}[\mathbf{\Sigma}_{jj}] - 2\|\mathbf{\Sigma}_{ij}\|_* \quad (2.5)$$

$$\cos \theta(h_i, h_j) = \frac{\|\mathbf{\Sigma}_{ij}\|_*}{\sqrt{\text{Tr}[\mathbf{\Sigma}_{ii}] \text{Tr}[\mathbf{\Sigma}_{jj}]}} \quad (2.6)$$

where  $\|\mathbf{\Sigma}_{ij}\|_*$  denotes the nuclear norm (or Shatten 1-norm) of the cross-covariance matrix.

Suppose we are given  $M$  independent and identically distributed network inputs  $\mathbf{z}_1, \dots, \mathbf{z}_M \sim P$ . We can estimate the generalized shape distances by substituting the empirical covariances:

$$\hat{\mathbf{\Sigma}}_{ii} = \frac{1}{M} \sum_{m=1}^M h_i(\mathbf{z}_m)h_i(\mathbf{z}_m)^\top, \quad \hat{\mathbf{\Sigma}}_{jj} = \frac{1}{M} \sum_{m=1}^M h_j(\mathbf{z}_m)h_j(\mathbf{z}_m)^\top, \quad \hat{\mathbf{\Sigma}}_{ij} = \frac{1}{M} \sum_{m=1}^M h_i(\mathbf{z}_m)h_j(\mathbf{z}_m)^\top \quad (2.7)$$

to approximate the true covariances appearing in eqs. (2.5) and (2.6). Thus,

$$\hat{\rho}^2(h_i, h_j) = \text{Tr}[\hat{\mathbf{\Sigma}}_{ii}] + \text{Tr}[\hat{\mathbf{\Sigma}}_{jj}] - 2\|\hat{\mathbf{\Sigma}}_{ij}\|_* \quad (2.8)$$

$$\cos \hat{\theta}(h_i, h_j) = \frac{\|\hat{\mathbf{\Sigma}}_{ij}\|_*}{\sqrt{\text{Tr}[\hat{\mathbf{\Sigma}}_{ii}] \text{Tr}[\hat{\mathbf{\Sigma}}_{jj}]}} \quad (2.9)$$

define plug-in estimators for the squared Procrustes and cosine Riemannian shape distances.

## 2.2 A new estimator with controllable bias

The plug-in estimator of  $\|\mathbf{\Sigma}_{ij}\|_*$  has low variance but large and slowly decaying bias (see theorems B.2 and B.1). Here we develop an alternative estimator that is nearly unbiased.

First, note that the eigenvalues of  $\mathbf{\Sigma}_{ij}\mathbf{\Sigma}_{ij}^\top$  correspond to the squared singular values of  $\mathbf{\Sigma}_{ij}$ . Thus,  $\text{Tr}[(\mathbf{\Sigma}_{ij}\mathbf{\Sigma}_{ij}^\top)^{1/2}] = \|\mathbf{\Sigma}_{ij}\|_*$ , and so we can reduce our problem to estimating the trace of  $(\mathbf{\Sigma}_{ij}\mathbf{\Sigma}_{ij}^\top)^{1/2}$ , which is symmetric. Leveraging ideas from a well-developed literature [1], we proceed to define the  $p^{\text{th}}$  moment of this matrix as:

$$W_p = \text{Tr}[(\mathbf{\Sigma}_{ij}\mathbf{\Sigma}_{ij}^\top)^p] = \sum_{n=1}^N \lambda_n^p \quad (2.10)$$

where  $\lambda_1, \dots, \lambda_N$  denote the eigenvalues of  $\mathbf{\Sigma}_{ij}\mathbf{\Sigma}_{ij}^\top$ . Now, for any function  $f: \mathbb{R} \mapsto \mathbb{R}$  and symmetric matrix  $\mathbf{S}$  with eigenvalues  $\lambda_1, \dots, \lambda_N$ , we define<sup>1</sup>  $\text{Tr}[f(\mathbf{S})] = \sum_i f(\lambda_i)$ . So long as  $f$  is reasonably well-behaved, we can approximate it using a truncated power series with  $P$  terms. Thus, with  $\mathbf{S} = \mathbf{\Sigma}_{ij}\mathbf{\Sigma}_{ij}^\top$  and  $f(x) = \sqrt{x}$ :

$$\|\mathbf{\Sigma}_{ij}\|_* = \text{Tr}[(\mathbf{\Sigma}_{ij}\mathbf{\Sigma}_{ij}^\top)^{1/2}] \approx \sum_{n=1}^N \sum_{p=0}^P \gamma_p \lambda_n^p = \sum_{p=0}^P \gamma_p \sum_{n=1}^N \lambda_n^p = \sum_{p=0}^P \gamma_p W_p \quad (2.11)$$

where  $\gamma_0, \dots, \gamma_P$  are scalar coefficients.

In summary, we can estimate  $\|\mathbf{\Sigma}_{ij}\|_*$  by (a) specifying an estimator of the top eigenmoments,  $W_1, \dots, W_P$ , and (b) specifying a desired set of scalar coefficients  $\gamma_0, \dots, \gamma_P$ . To estimate the eigenmoments, we adapt procedures described by Kong and Valiant [12] to obtain unbiased estimates for each moment,  $\hat{W}_1, \dots, \hat{W}_P$  (see App. C). To select the scalar coefficients, we propose an optimization procedure that trades off between bias and variance in the estimate of  $\|\mathbf{\Sigma}_{ij}\|_*$ . Our starting point is the usual bias-variance decomposition:

$$\mathbb{E} \left[ \left( \|\mathbf{\Sigma}_{ij}\|_* - \sum_p \gamma_p \hat{W}_p \right)^2 \right] = \left( \mathbb{E} \left[ \|\mathbf{\Sigma}_{ij}\|_* - \sum_p \gamma_p \hat{W}_p \right] \right)^2 + \text{Var} \left[ \sum_p \gamma_p \hat{W}_p \right]. \quad (2.12)$$

<sup>1</sup>This is a common convention to extend scalar functions [see e.g. 20, sec. 1.2.6].

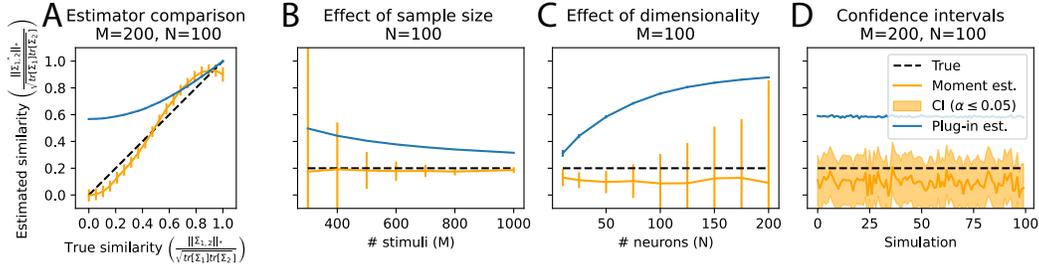


Figure 2: Validation of estimator on synthetic data. **(A)** The moment based estimator (orange) compared to plug-in estimator (blue) in simulation with standard deviation bars calculated across simulations. Estimators are evaluated at 20 linearly spaced ground truth similarity score values. **(B)** Effect of increasing sample size when moment estimator is constrained to have a bias less than 5%. **(C)** Effect of increasing dimensionality. **(D)** Demonstration of conservative confidence intervals that account for variance and maximal bias of moment estimator. We do not include CIs for the plug in estimator (implied by theorem B.1) because for small sample sizes, the theoretical bounds on estimator bias always contain far more than the entire allowable interval  $([0, 1])$ .

Since  $\mathbb{E}[\hat{W}_p] = W_p = \sum_n \lambda_n^p$ , the first term above (i.e. the “bias”) simplifies and is upper-bounded:

$$\left( \mathbb{E} \left[ \|\Sigma_{ij}\|_* - \sum_p \gamma_p \hat{W}_p \right] \right)^2 = \left( \sum_n \left( \lambda_n^{1/2} - \sum_p \gamma_p \lambda_n^p \right) \right)^2 \leq \max_{0 \leq x \leq 1} \left( N \left( x^{1/2} - \sum_p \gamma_p x^p \right) \right)^2$$

The inequality follows from replacing each term in the sum over  $n$  with the worst case approximation error of the polynomial expansion (given here as the maximization over  $x$ ). Thus, we seek to:

$$\text{minimize}_{\gamma_0, \dots, \gamma_P} \max_{0 \leq x \leq 1} \left( N \left( x^{1/2} - \sum_p \gamma_p x^p \right) \right)^2 + \sum_{p, p'} \gamma_p \gamma_{p'} \text{Cov}(\hat{W}_p, \hat{W}_{p'}). \quad (2.13)$$

We estimate  $\text{Cov}(\hat{W}_p, \hat{W}_{p'})$  by bootstrapping—i.e. the empirical covariance of these statistics across re-sampled datasets where  $\mathbf{z}_{1:M}$  are sampled with replacement. Given this estimate of covariance, eq. (2.13) can be cast as a convex quadratic program and the maximal bias can be bounded to a user defined limit at the expense of variance (see App. C.2). We use the maximal bias (eq. 2.13, term 1) and variance (eq. 2.13, term 2) to form approximate confidence intervals (see App. C.3).

### 2.3 Validation on synthetic data

We validate our method-of-moments estimator (section 2.2) on simulated representations jointly sampled from a multivariate normal distribution. We consider estimating the cosine shape similarity,  $\cos \theta$ , defined in eq. 2.6. Our estimator of  $\|\Sigma_{ij}\|_*$  is the principle novelty; thus, it is informative to understand its properties in isolation. To achieve this, in our experiments we use the ground truth covariance of  $\hat{W}_p$  (instead of an estimate from a bootstrap) and use the ground truth values of  $\text{Tr}[\Sigma_{ii}]$  and  $\text{Tr}[\Sigma_{jj}]$ . To draw data for our simulations, we set the eigenvalues of the  $\Sigma_{ii}$  and the singular values of  $\Sigma_{ij}$  to a ground truth nuclear norm and similarity score. To demonstrate the estimators accuracy across the space of orthogonal transformations we apply a random orthogonal rotation matrix to each population’s covariance in each new parameter setting.

We first compared the bias of the plug-in estimator to that of the moment-based estimator across a range of ground truth shape similarity values (Fig. 2A). As expected from our intuition discussed in App. B.1, the plug-in estimator (blue line) tends to grossly inflate estimated similarity when ground truth similarity is low (left side of plot). The moment-based estimator (orange line), in contrast, performs reasonably well over the full range of simulations, at the cost of modest increases in estimator variance (blue vs orange error bars).

Next, we fixed the ground truth similarity at 0.2 and studied the effect of sample size,  $M$  (Fig. 2B). The moment estimator (constrained to 5% bias) maintains small bias even with small  $M$ , at the cost of high variance (large orange error bars). Increasing  $M$  quickly reduces the variance of the estimator. A similar story emerges when we fix  $M$  and vary the ambient dimension  $N$  (Fig. 2C). As the

dimensionality increases, the the plug-in estimator bias quickly explodes. In contrast, the moment estimator (here constrained to 10% bias) has roughly constant bias; however, it’s variance grows with  $N$ . Thus our estimator bias outperforms the plug-in sample size is low and dimensionality is high.

Finally, an important property of the moment-based estimator is our ability to compute approximate confidence intervals (CI) (see App. C.3). We demonstrate 95% CIs across simulations in Figure 2D (shaded orange region). These CIs are conservative, the true shape score is not within the CI’s for only 2.3% of simulations. Results on neural data can be found in App. D.2.

### 3 Discussion

There is a vast literature of papers that utilize or develop measures of representational similarity between neural networks [see 11, for review], and recent works have shown interest in leveraging representational distances that satisfy the triangle inequality [6, 7, 16, 30]. Yet, the statistical properties of these shape distance measures appears understudied. Here, rigorously analyzed of “plug-in” estimates of shape distance in high-dimensional, noisy, and sample-limited regimes. Our analysis showed that these estimates (a) tend to over-estimate representational similarity when the true similarity is small and (b) require a large number of samples,  $M$ , to overcome this bias in high dimensional regimes. Theorems B.1 and B.2 provide precise guarantees on the worst-case performance of plug-in estimators, which should guide the design of biological experiments and analyses of their statistical power.

An equally important contribution of our work is to provide a practical method to (a) reduce the bias of plug-in estimators of shape distance, (b) quantify uncertainty in shape distance estimates, and (c) enable practitioners to explicitly trade off estimator bias and variance. When employed on a biological dataset published by Stringer et al. [25], we find that shape similarity estimates are highly uncertain, revealing the challenging nature of the problem in high dimensions and with noisy data. Importantly, this degree of uncertainty is not obvious from the procedures and plug-in estimates advertised by existing work on this subject.

Both theoretical and methodological aspects of our work may be of broader interest beyond the immediate subject of shape distance estimation. We have seen that estimating the nuclear norm of the cross-covariance,  $\|\Sigma_{ij}\|_*$ , is the key challenge in our problem. Estimating the spectrum of cross-covariance matrices is a topic of contemporary interest [2], and further exploring the connections between this problem and shape distance estimation is an intriguing direction. Similarly, the method-of-moments estimator presented in section 2.2 is broadly applicable to generalized trace estimation [1]. While others have used polynomial expansions in this context [17], a key novelty of our approach is the selection of coefficients with a tunable parameter that explicitly trades off estimator bias and variance. A more typical approach would be to choose these coefficients based on a Chebyshev polynomial expansion. While elegant, we believe our procedure for tuning these coefficients will be more relevant to scientific applications where samples are limited (such as neural data) and practitioners desire finer-scale control.

In summary, our work is one of the first to rigorously interrogate the statistical challenges of estimating shape distances in high-dimensional spaces. While shape distances can be well-behaved in certain settings (e.g. in artificial networks where a very large number of inputs can be sampled), our theoretical results and empirical observations underscore the challenging nature of this problem, suggesting the need for carefully designed biological experiments and estimation procedures.

### References

- [1] R. P. Adams, J. Pennington, M. J. Johnson, J. Smith, Y. Ovadia, B. Patton, and J. Saunderson. Estimating the spectral density of large implicit matrices, 2018.
- [2] F. Benaych-Georges, J.-P. Bouchaud, and M. Potters. Optimal cleaning for singular values of cross-covariance matrices. *The Annals of Applied Probability*, 33(2):1295 – 1326, 2023. doi: 10.1214/22-AAP1842. URL <https://doi.org/10.1214/22-AAP1842>.
- [3] S. Boyd and L. Vandenberghe. *Convex optimization*. Cambridge University Press, 2004.
- [4] M. Cai, N. W. Schuck, J. W. Pillow, and Y. Niv. A bayesian method for reducing bias in neural representational similarity analysis. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon,

- and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016. URL <https://proceedings.neurips.cc/paper/2016/file/b06f50d1f89bd8b2a0fb771c1a69c2b0-Paper.pdf>.
- [5] T. Cohen and M. Welling. Group equivariant convolutional networks. In M. F. Balcan and K. Q. Weinberger, editors, *Proceedings of The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pages 2990–2999, New York, New York, USA, 20–22 Jun 2016. PMLR. URL <https://proceedings.mlr.press/v48/cohenc16.html>.
- [6] L. R. Duong, J. Zhou, J. Nassar, J. Berman, J. Olieslagers, and A. H. Williams. Representational dissimilarity metric spaces for stochastic neural networks. In *International Conference on Learning Representations*, 2023.
- [7] H. Giaffar, C. R. Buxó, and M. Aoi. The effective number of shared dimensions: A simple method for revealing shared structure between datasets. *bioRxiv*, 2023. doi: 10.1101/2023.07.27.550815. URL <https://www.biorxiv.org/content/early/2023/07/28/2023.07.27.550815>.
- [8] J. C. Gower and G. B. Dijkstra. *Procrustes problems*, volume 30. OUP Oxford, 2004.
- [9] D. G. Kendall, D. Barden, T. K. Carne, and H. Le. *Shape and Shape Theory*. John Wiley & Sons, Sept. 2009.
- [10] T. C. Kietzmann, C. J. Spoerer, L. K. A. Sørensen, R. M. Cichy, O. Hauk, and N. Kriegeskorte. Recurrence is required to capture the representational dynamics of the human visual system. *Proceedings of the National Academy of Sciences*, 116(43):21854–21863, 2019. doi: 10.1073/pnas.1905544116. URL <https://www.pnas.org/doi/abs/10.1073/pnas.1905544116>.
- [11] M. Klabunde, T. Schumacher, M. Strohmaier, and F. Lemmerich. Similarity of neural network models: A survey of functional and representational measures, 2023.
- [12] W. Kong and G. Valiant. Spectrum estimation from samples. *The Annals of Statistics*, 45(5):2218–2247, 2017. doi: 10.1214/16-AOS1525. URL <https://doi.org/10.1214/16-AOS1525>.
- [13] S. Kornblith, M. Norouzi, H. Lee, and G. Hinton. Similarity of neural network representations revisited. In K. Chaudhuri and R. Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 3519–3529. PMLR, 09–15 Jun 2019. URL <https://proceedings.mlr.press/v97/kornblith19a.html>.
- [14] N. Kriegeskorte, M. Mur, and P. Bandettini. Representational similarity analysis - connecting the branches of systems neuroscience. *Front. Syst. Neurosci.*, 2:4, Nov. 2008.
- [15] N. Kriegeskorte, M. Mur, D. A. Ruff, R. Kiani, J. Bodurka, H. Esteky, K. Tanaka, and P. A. Bandettini. Matching categorical object representations in inferior temporal cortex of man and monkey. *Neuron*, 60(6):1126–1141, Dec. 2008.
- [16] R. D. Lange, D. Kwok, J. Matelsky, X. Wang, D. S. Rolnick, and K. P. Kording. Neural networks as paths through the space of representations. *arXiv preprint arXiv:2206.10999*, 2022.
- [17] L. Lin, Y. Saad, and C. Yang. Approximating spectral densities of large matrices. *SIAM Review*, 58(1):34–65, 2016. doi: 10.1137/130934283.
- [18] N. Maheswaranathan, A. Williams, M. Golub, S. Ganguli, and D. Sussillo. Universality and individuality in neural dynamics across large populations of recurrent networks. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL [https://proceedings.neurips.cc/paper\\_files/paper/2019/file/5f5d472067f77b5c88f69f1bcfda1e08-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2019/file/5f5d472067f77b5c88f69f1bcfda1e08-Paper.pdf).
- [19] T. Nguyen, M. Raghu, and S. Kornblith. Do wide and deep networks learn the same things? uncovering how neural network representations vary with width and depth. In *International Conference on Learning Representations*, 2021.

- [20] M. Potters and J.-P. Bouchaud. *A First Course in Random Matrix Theory: for Physicists, Engineers and Data Scientists*. Cambridge University Press, 2020. doi: 10.1017/9781108768900.
- [21] M. Raghu, J. Gilmer, J. Yosinski, and J. Sohl-Dickstein. Svcca: Singular vector canonical correlation analysis for deep learning dynamics and interpretability. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- [22] H. H. Schütt, A. D. Kipnis, J. Diedrichsen, and N. Kriegeskorte. Statistical inference on representational geometries, 2021.
- [23] J. Shi, E. Shea-Brown, and M. Buice. Comparison against task driven artificial neural networks reveals functional properties in mouse visual cortex. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL [https://proceedings.neurips.cc/paper\\_files/paper/2019/file/748d6b6ed8e13f857ceaa6cfbdca14b8-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2019/file/748d6b6ed8e13f857ceaa6cfbdca14b8-Paper.pdf).
- [24] K. R. Storrs, T. C. Kietzmann, A. Walther, J. Mehrer, and N. Kriegeskorte. Diverse Deep Neural Networks All Predict Human Inferior Temporal Cortex Well, After Training and Fitting. *Journal of Cognitive Neuroscience*, 33(10):2044–2064, 09 2021. ISSN 0898-929X. doi: 10.1162/jocn\_a\_01755. URL [https://doi.org/10.1162/jocn\\_a\\_01755](https://doi.org/10.1162/jocn_a_01755).
- [25] C. Stringer, M. Pachitariu, N. Steinmetz, M. Carandini, and K. D. Harris. High-dimensional geometry of population responses in visual cortex. *Nature*, 571(7765):361–365, 2019.
- [26] J. A. Tropp. An introduction to matrix concentration inequalities. *Foundations and Trends® in Machine Learning*, 8(1-2):1–230, 2015. ISSN 1935-8237. doi: 10.1561/22000000048. URL <http://dx.doi.org/10.1561/22000000048>.
- [27] M. J. Wainwright. *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48. Cambridge university press, 2019.
- [28] A. Walther, H. Nili, N. Ejaz, A. Alink, N. Kriegeskorte, and J. Diedrichsen. Reliability of dissimilarity measures for multi-voxel pattern analysis. *NeuroImage*, 137:188–200, 2016.
- [29] A. H. Williams and S. W. Linderman. Statistical neuroscience in the single trial limit. *Current Opinion in Neurobiology*, 70:193–205, 2021. ISSN 0959-4388. doi: <https://doi.org/10.1016/j.conb.2021.10.008>. URL <https://www.sciencedirect.com/science/article/pii/S0959438821001203>. Computational Neuroscience.
- [30] A. H. Williams, E. Kunz, S. Kornblith, and S. Linderman. Generalized shape metrics on neural representations. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P. Liang, and J. W. Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 4738–4750. Curran Associates, Inc., 2021. URL <https://proceedings.neurips.cc/paper/2021/file/252a3dbaeb32e7690242ad3b556e626b-Paper.pdf>.

## A Appendix: Background on Generalized Shape Metrics

### A.1 Definition of Generalized Shape Metrics

Intuitively, a measure of distance between neural representations should be invariant nuisance symmetries in the neural representation, such as arbitrary permutations over neuron labels [5]. Representational similarity measures are typically designed to be invariant not only to permutations, but also to rotations, reflections, translations, and isotropic scalings in neural firing rate space [13, 14].

We begin by considering a simple setting where each neural network is a deterministic map (for the stochastic setting, see appendix A.2). A collection of  $K$  neural systems can then be viewed as a set of functions, each denoted  $h_i : \mathcal{Z} \mapsto \mathbb{R}^N$  for  $i \in \{1, \dots, K\}$ . Here,  $\mathcal{Z}$  is a feature space and  $N$  can be interpreted as the number of neurons in each system (e.g. the size of a hidden layer in an artificial network, or the number of recorded neurons in a biological experiment).<sup>2</sup>

Motivated by the shape theory literature [9, 30], we consider estimating the *Procrustes size-and-shape distance*,  $\rho$ , and *Riemannian shape distance*,  $\theta$ , between neural representations. Let  $h_i$  and  $h_j$  denote neural systems that are mean-centered and bounded:

$$\mathbb{E}[h_i(\mathbf{z})] = \mathbb{E}[h_j(\mathbf{z})] = \mathbf{0} \quad \text{and} \quad \|h_i(\mathbf{z})\|_2, \|h_j(\mathbf{z})\|_2 < B\sqrt{N} \quad \text{almost surely.}$$

for some constant  $B > 0$ . Here, the expectations are taken over  $\mathbf{z} \sim P$ , for some distribution  $P$  over network inputs. Our assumption that neural population rates are bounded by  $B\sqrt{N}$  can be achieved by assuming that each neuron has a maximum firing rate equal to  $B$ . This assumption is common in the literature and reasonable in both artificial networks (since connection weights are finite) and biological networks (since neurons have a maximal firing rate).

The Procrustes and Riemannian shape distances can be defined [App. D in 30]:

$$\rho(h_i, h_j) = \min_{\mathbf{Q} \in \mathcal{O}(N)} \sqrt{\mathbb{E}\|h_i(\mathbf{z}) - \mathbf{Q}h_j(\mathbf{z})\|_2^2}$$

$$\theta(h_i, h_j) = \min_{\mathbf{Q} \in \mathcal{O}(N)} \cos^{-1} \left( \frac{\mathbb{E}[h_i(\mathbf{z})^\top \mathbf{Q}h_j(\mathbf{z})]}{\sqrt{\mathbb{E}[h_i(\mathbf{z})^\top h_i(\mathbf{z})]\mathbb{E}[h_j(\mathbf{z})^\top h_j(\mathbf{z})]}} \right)$$

where  $\mathcal{O}(N)$  denotes the set of  $N \times N$  orthogonal matrices. Again, all expectations are taken over  $\mathbf{z} \sim P$ . Note that different notions of distance arise from different choices of input distribution,  $P$ .

To simplify our analysis and exposition, we will focus on estimating the *squared Procrustes distance*,  $\rho^2$ , and what we call the *cosine shape similarity*,  $\cos \theta$ . Thus, we ignore the square root term in eq. (2.2) and the arccosine term in eq. (2.3), but it should be kept in mind that one must apply these nonlinear functions to achieve a proper metric.

**Properties of Shape Distance** It is easy to verify that shape distances are invariant to rotations and reflections: that is, if  $r : \mathbb{R}^N \mapsto \mathbb{R}^N$  is an orthogonal transformation, then for any function  $h : \mathcal{Z} \mapsto \mathbb{R}^N$  representing a neural system we have  $\rho(h, r \circ h) = \theta(h, r \circ h) = 0$ , where ‘ $\circ$ ’ denotes function composition. Furthermore,  $\rho$  and  $\theta$  are proper metrics, meaning that:

$$\rho(h_i, h_j) = \rho(h_j, h_i) \quad \text{and} \quad \rho(h_i, h_j) \leq \rho(h_i, h_k) + \rho(h_k, h_j) \quad \forall i, j, k \in \{1, \dots, K\}, \quad (\text{A.1})$$

and likewise for  $\theta$ . These properties are fundamental to rigorously establishing downstream analyses, such as for clustering networks with similar representations [30].

It is well-known that the optimal orthogonal alignment appearing in eqs. (2.2) and (2.3) can be identified in closed form. Leveraging this, we can use the covariance and cross-covariance matrices,

$$\mathbf{\Sigma}_{ii} = \mathbb{E}[h_i(\mathbf{z})h_i(\mathbf{z})^\top], \quad \mathbf{\Sigma}_{jj} = \mathbb{E}[h_j(\mathbf{z})h_j(\mathbf{z})^\top], \quad \mathbf{\Sigma}_{ij} = \mathbb{E}[h_i(\mathbf{z})h_j(\mathbf{z})^\top], \quad (\text{A.2})$$

to reformulate the squared Procrustes distance and cosine shape similarity:

$$\rho^2(h_i, h_j) = \text{Tr}[\mathbf{\Sigma}_{ii}] + \text{Tr}[\mathbf{\Sigma}_{jj}] - 2\|\mathbf{\Sigma}_{ij}\|_*$$

$$\cos \theta(h_i, h_j) = \frac{\|\mathbf{\Sigma}_{ij}\|_*}{\sqrt{\text{Tr}[\mathbf{\Sigma}_{ii}]\text{Tr}[\mathbf{\Sigma}_{jj}]}}$$

<sup>2</sup>The assumption that each layer has the same number of neurons is not necessary, and only made for convenience. For networks with dissimilar sizes, we can preprocess by zero-padding the smaller network. If necessary, one could alternatively perform PCA on the larger network to reduce to a common dimension.

where  $\|\Sigma_{ij}\|_*$  denotes the nuclear norm (or Shatten 1-norm) of the cross-covariance matrix:

$$\|\Sigma_{ij}\|_* = \sum_{n=1}^N s_n(\Sigma_{ij}) \quad (\text{A.3})$$

where  $s_1(M) \geq \dots \geq s_N(M) \geq 0$  denote the singular values of a matrix  $M$ . Equations 2.5 and 2.6 are derived in Appendix A.3 to provide the reader with a self-contained narrative.

**Plug-in Estimators** Suppose we are given  $M$  independent and identically distributed network inputs  $z_1, \dots, z_M \sim P$ . How well can we approximate the shape distances between two networks, as a function of  $M$ ? The standard approach, which was previously used in Williams et al. [30], is to use a *plug-in estimator* in which one computes eqs. (2.2) and (2.3) after identifying the optimal  $Q \in \mathcal{O}(N)$ . As we show in App. A.4, this is equivalent to estimating the squared Procrustes and cosine Riemannian distances by substituting the empirical covariances:

$$\hat{\Sigma}_{ii} = \frac{1}{M} \sum_{m=1}^M h_i(z_m)h_i(z_m)^\top, \quad \hat{\Sigma}_{jj} = \frac{1}{M} \sum_{m=1}^M h_j(z_m)h_j(z_m)^\top, \quad \hat{\Sigma}_{ij} = \frac{1}{M} \sum_{m=1}^M h_i(z_m)h_j(z_m)^\top \quad (\text{A.4})$$

to approximate the true covariances appearing in eqs. (2.5) and (2.6). Thus,

$$\hat{\rho}^2(h_i, h_j) = \text{Tr}[\hat{\Sigma}_{ii}] + \text{Tr}[\hat{\Sigma}_{jj}] - 2\|\hat{\Sigma}_{ij}\|_*$$

$$\cos \hat{\theta}(h_i, h_j) = \frac{\|\hat{\Sigma}_{ij}\|_*}{\sqrt{\text{Tr}[\hat{\Sigma}_{ii}] \text{Tr}[\hat{\Sigma}_{jj}]}}$$

define plug-in estimators for the squared Procrustes and cosine Riemannian shape distances. The empirical behavior of these estimators as a function of  $M$  was only briefly characterized by Williams et al. [30] for a pair of artificial networks trained on CIFAR-10.

## A.2 Extension to stochastic networks

Thus far, we have modeled neural networks as deterministic mappings,  $h_i : \mathcal{Z} \mapsto \mathbb{R}^N$ . This assumption is not satisfied in biological data and in many artificial networks (e.g. VAEs). Here, we briefly explain how to extend the estimators to the stochastic setting. In this setting, the response of network  $i$  can be written as  $h_i(z) + \epsilon_i(z)$ . As before,  $h_i(z)$  is a deterministic mapping conditioned on a random variable  $z \sim P$ . The “noise” term  $\epsilon_i(z)$  is a mean-zero random variable that, in addition to inheriting the randomness of  $z$ , captures the stochastic elements of each forward pass through the network (i.e. trial-to-trial variability even when the stimulus is fixed). Importantly, noise contributions are independent and identically distributed for each pass through the network.

Given a second stochastic network with same structure,  $h_j(z) + \epsilon_j(z)$ , our goal is to estimate the shape distances eqs. (2.2) and (2.3) as before, effectively ignoring contributions of the “noise” terms  $\epsilon_i(\cdot)$  and  $\epsilon_j(\cdot)$ . Ignoring these terms is not wholly justified, since it is of great interest to quantify how noise varies across networks [6]. Nonetheless, it is useful to develop metrics that isolate the “signal” component of neural representations, and a full development of methods to quantify similarity in noise structure is outside the scope of this paper.

Our basic observation is that it suffices to consider two replicates for each network input. That is, let  $z' = z$  where  $z \sim P$ . Then,  $\Sigma_{ii} = \mathbb{E}[h_i(z)h_i(z')^\top]$  which can be approximated by the slightly reformulated plug-in estimator:  $\hat{\Sigma}_{ii} = (1/M) \sum_m h_i(z_m)h_i(z'_m)^\top$ . Further, since noise is independent across networks, i.e.  $\epsilon_i(z) \perp\!\!\!\perp \epsilon_j(z)$  for all  $z \in \mathcal{Z}$ , the cross-covariance estimators, including the method-of-moments estimator described in section 2.2, do not require any modification. Here we provide several relevant derivations for generalized shape metrics. For a more thorough review, we direct the reader to [30] for the foundational results on generalized shape metrics and [6] for the extension to stochastic neural networks.

### A.3 Equivalence of eqs. (2.2) and (2.5); eqs. (2.3) and (2.6)

The squared Procrustes can be reformulated in terms of the covariance and cross-covariance matrices as follows:

$$\begin{aligned}
\rho^2(h_i, h_j) &= \min_{\mathbf{Q} \in \mathcal{O}(N)} \mathbb{E} \|h_i(\mathbf{z}) - \mathbf{Q}h_j(\mathbf{z})\|_2^2 \\
&= \min_{\mathbf{Q} \in \mathcal{O}(N)} \mathbb{E} [h_i(\mathbf{z})^\top h_i(\mathbf{z}) + h_j(\mathbf{z})^\top h_j(\mathbf{z}) - 2h_i(\mathbf{z})^\top \mathbf{Q}h_j(\mathbf{z})] \\
&= \mathbb{E} [h_i(\mathbf{z})^\top h_i(\mathbf{z})] + \mathbb{E} [h_j(\mathbf{z})^\top h_j(\mathbf{z})] - 2 \max_{\mathbf{Q} \in \mathcal{O}(N)} \mathbb{E} [h_i(\mathbf{z})^\top \mathbf{Q}h_j(\mathbf{z})] \\
&= \mathbb{E} [\text{Tr} [h_i(\mathbf{z})h_i(\mathbf{z})^\top]] + \mathbb{E} [\text{Tr} [h_j(\mathbf{z})h_j(\mathbf{z})^\top]] - 2 \max_{\mathbf{Q} \in \mathcal{O}(N)} \mathbb{E} [\text{Tr} [\mathbf{Q}h_j(\mathbf{z})h_i(\mathbf{z})^\top]] \\
&= \text{Tr} [\mathbb{E} [h_i(\mathbf{z})h_i(\mathbf{z})^\top]] + \text{Tr} [\mathbb{E} [h_j(\mathbf{z})h_j(\mathbf{z})^\top]] - 2 \max_{\mathbf{Q} \in \mathcal{O}(N)} \text{Tr} [\mathbf{Q}\mathbb{E} [h_j(\mathbf{z})h_i(\mathbf{z})^\top]] \\
&= \text{Tr} [\boldsymbol{\Sigma}_{ii}] + \text{Tr} [\boldsymbol{\Sigma}_{jj}] - 2 \max_{\mathbf{Q} \in \mathcal{O}(N)} \text{Tr} [\mathbf{Q}\boldsymbol{\Sigma}_{ij}] \\
&= \text{Tr} [\boldsymbol{\Sigma}_{ii}] + \text{Tr} [\boldsymbol{\Sigma}_{jj}] - 2\|\boldsymbol{\Sigma}_{ij}\|_*
\end{aligned}$$

Similarly for the cosine Riemannian distance:

$$\begin{aligned}
\cos \theta(h_i, h_j) &= \max_{\mathbf{Q} \in \mathcal{O}(N)} \left( \frac{\mathbb{E}[h_i(\mathbf{z})^\top \mathbf{Q}h_j(\mathbf{z})]}{\sqrt{\mathbb{E}[h_i(\mathbf{z})^\top h_i(\mathbf{z})]\mathbb{E}[h_j(\mathbf{z})^\top h_j(\mathbf{z})]}} \right) \\
&= \frac{\max_{\mathbf{Q} \in \mathcal{O}(N)} \mathbb{E} [\text{Tr} [\mathbf{Q}h_j(\mathbf{z})h_i(\mathbf{z})^\top]]}{\sqrt{\mathbb{E} [\text{Tr} [h_i(\mathbf{z})h_i(\mathbf{z})^\top]] \mathbb{E} [\text{Tr} [h_j(\mathbf{z})h_j(\mathbf{z})^\top]]}} \\
&= \frac{\max_{\mathbf{Q} \in \mathcal{O}(N)} \text{Tr} [\mathbf{Q}\mathbb{E} [h_j(\mathbf{z})h_i(\mathbf{z})^\top]]}{\sqrt{\text{Tr} [\mathbb{E} [h_i(\mathbf{z})h_i(\mathbf{z})^\top]] \text{Tr} [\mathbb{E} [h_j(\mathbf{z})h_j(\mathbf{z})^\top]]}} \\
&= \frac{\max_{\mathbf{Q} \in \mathcal{O}(N)} \text{Tr} [\mathbf{Q}\boldsymbol{\Sigma}_{ij}]}{\sqrt{\text{Tr} [\boldsymbol{\Sigma}_{ii}] \text{Tr} [\boldsymbol{\Sigma}_{jj}]}]} = \frac{\|\boldsymbol{\Sigma}_{ij}\|_*}{\sqrt{\text{Tr} [\boldsymbol{\Sigma}_{ii}] \text{Tr} [\boldsymbol{\Sigma}_{jj}]}]}
\end{aligned}$$

### A.4 Reformulations of the Plug-in Estimator of Procrustes distance

Let  $\mathbf{z}_1, \dots, \mathbf{z}_M$  denote a set of independently and identically distributed samples in the network input space. Then, stack the responses of network  $i$  row-wise into a matrix  $\mathbf{X}_i \in \mathbb{R}^{M \times N}$ . Given this set up, a common definition of Procrustes distance is [8]:

$$\min_{\mathbf{Q} \in \mathcal{O}(N)} \frac{1}{\sqrt{M}} \|\mathbf{X}_i - \mathbf{X}_j \mathbf{Q}\|_F \tag{A.5}$$

Here, we have included a multiplying factor of  $1/\sqrt{M}$  for reasons that will become clear shortly. Aside from this factor, the quantity above is how Williams et al. [30] define the Procrustes distance. Below, we show that the square of this quantity is indeed the plug-in estimator we defined in eq. (2.8) in terms of the empirical covariance matrices:

$$\begin{aligned}
\min_{\mathbf{Q} \in \mathcal{O}(N)} \frac{1}{M} \|\mathbf{X}_i - \mathbf{X}_j \mathbf{Q}\|_F^2 &= \min_{\mathbf{Q} \in \mathcal{O}(N)} \frac{1}{M} (\text{Tr}[\mathbf{X}_i^\top \mathbf{X}_i] + \text{Tr}[\mathbf{X}_j^\top \mathbf{X}_j] - 2 \text{Tr}[\mathbf{X}_i \mathbf{X}_j^\top \mathbf{Q}]) \\
&= \text{Tr} [\frac{1}{M} \mathbf{X}_i^\top \mathbf{X}_i] + \text{Tr} [\frac{1}{M} \mathbf{X}_j^\top \mathbf{X}_j] - 2 \max_{\mathbf{Q} \in \mathcal{O}(N)} \text{Tr} [\frac{1}{M} \mathbf{X}_i \mathbf{X}_j^\top \mathbf{Q}] \\
&= \text{Tr} [\hat{\boldsymbol{\Sigma}}_{ii}] + \text{Tr} [\hat{\boldsymbol{\Sigma}}_{jj}] - 2 \max_{\mathbf{Q} \in \mathcal{O}(N)} \text{Tr} [\hat{\boldsymbol{\Sigma}}_{ij} \mathbf{Q}] \\
&= \text{Tr} [\hat{\boldsymbol{\Sigma}}_{ii}] + \text{Tr} [\hat{\boldsymbol{\Sigma}}_{jj}] - 2\|\hat{\boldsymbol{\Sigma}}_{ij}\|_* \\
&= \hat{\rho}^2(h_i, h_j)
\end{aligned}$$

## B Appendix: Plug-in Estimator Theory

Here we provide a number of derivations related to the behavior of the plug-in estimator for generalized shape metrics. These results primarily rely on classic concentration inequalities and results from

random matrix theory. For readers interested in further background, we provide pointers to [27] and [26] for the concentration inequalities and [20] for the random matrix theory.

### B.1 Summary of Results: Nonasymptotic bounds on the performance of plug-in estimation

First, it is straightforward to estimate  $\text{Tr}[\Sigma_{ii}]$  and  $\text{Tr}[\Sigma_{jj}]$ . Their plug-in estimators are unbiased under our assumptions in eq. (2.1), and they rapidly converge to the correct answer. This is shown in the following lemma, whose proof relies only on classical concentration inequalities.

**Lemma B.1** (App. B.3). *Under the assumptions in eq. (2.1), with probability at least  $1 - \delta$ :*

$$\left| \text{Tr}[\Sigma_{ii}] - \text{Tr}[\hat{\Sigma}_{ii}] \right| \leq BN^{1/2}M^{-1/2}\sqrt{2\log(2/\delta)} \quad (\text{B.1})$$

In contrast, the plug-in estimator for  $\|\Sigma_{ij}\|_*$  is biased upwards (see appendix B.2) and turns out to converge more slowly. Using the Matrix Bernstein inequality [see 26], we can show:

**Lemma B.2** (App. B.4). *Under the assumptions in eq. (2.1), for any  $M$  and  $N$ :*

$$\mathbb{E} \left| \|\hat{\Sigma}_{ij}\|_* - \|\Sigma_{ij}\|_* \right| < \frac{2B^2N^2\log(2N)}{3M} + \frac{2B^2N^2\sqrt{\log(2N)}}{M^{1/2}} \quad (\text{B.2})$$

This only upper bounds the expected error. However, the fluctuations around this expectation turn out to be small (see App. B.5), and so we are able to combine lemmas B.1 and B.2 into the following:

**Theorem B.1** (App. B.5). *Under the assumptions in eq. (2.1), with probability at least  $1 - \delta$*

$$\frac{|\hat{\rho}^2 - \rho^2|}{N} \leq \frac{2B^2N\log(2N)}{3M} + \frac{2B^2N\sqrt{\log(2N)}}{M^{1/2}} + \left( \frac{B^2}{M^{1/2}} + \frac{2B}{N^{1/2}M^{1/2}} \right) \sqrt{2\log\left(\frac{6}{\delta}\right)} \quad (\text{B.3})$$

Theorem B.1 states a non-asymptotic upper bound on the plug-in estimator’s error that holds with high probability. We have expressed this bound on the squared size-and-shape Procrustes distance normalized by  $1/N$ , since the raw error,  $|\hat{\rho} - \rho|$ , will tend to increase linearly with  $N$  for an uninteresting reason—namely, since the the Procrustes shape distance is comprised of terms like  $\text{Tr}[\Sigma_{ii}]$  and  $\text{Tr}[\Sigma_{jj}]$ . The choice of normalization in theorem B.1 also makes the result more comparable to the cosine shape similarity (eq. 2.6), which is normalized by a factor,  $\sqrt{\text{Tr}[\Sigma_{ii}]\text{Tr}[\Sigma_{jj}]}$ , of order  $N$ .

We can gain intuition for theorem B.1 by ignoring logarithmic factors and noticing that the second term dominates. Then, roughly speaking, theorem B.1 says that we can guarantee the plug-in error decreases as a function of  $NM^{-1/2}$ . Thus, for any fixed  $N$ , we need to increase  $M$  by a factor of 4 to decrease estimation error by a factor of 2. Further, when comparing higher-dimensional neural representations (i.e. higher  $N$ ) we need to sample more landmarks—if  $N$  increases by a factor of 2, then  $M$  must be increased by a factor of 4 to compensate.

### B.2 Summary of Results: Failure modes of plug-in estimation and a lower bound on performance

Theorem B.1 provides a high probability upper bound on the estimation error. A natural question is whether this upper bound is tight. To investigate, we seek an example where the plug-in estimator performs badly. We intuited that when two neural representations are very far apart in shape space, the plug-in estimator of shape distance should have a large downward bias. This can be understood in two ways. First, from the definitions of  $\rho$  and  $\theta$  in eqs. (2.2) and (2.3), we see that both expressions contain a minimization over  $\mathbf{Q} \in \mathcal{O}(N)$ . For large  $N$  and small  $M$ , this high-dimensional orthogonal matrix can be “overfit” to the  $M$  observations resulting in an underestimate of distance. Second, from the alternative formulations in eqs. (2.5) and (2.6), we see that the shape distance is large if the true cross-covariance is “small” as quantified by the nuclear norm. In the extreme case where the singular values of  $\Sigma_{ij}$  are all zero, the empirical cross-covariance matrix  $(1/M)\sum_m h_i(\mathbf{z}_m)h_j(\mathbf{z}_m)^\top$  will overestimate the nuclear norm, and therefore underestimate the shape distance. This is more severe when  $M$  is small, since there are fewer terms in the sum to “average out” spurious correlations, which are particularly problematic in high dimensions (i.e. when  $N$  is large).

This intuition led us to construct an example where plug-in estimation error approaches the upper bound in theorem B.1. This is summarized in the following result.

**Theorem B.2** (Lower Bound, App. B.6). *Under the assumptions in eq. (2.1), there exist neural networks and a distribution over inputs such that in the limit that  $N \rightarrow \infty$  and  $M \gg N$ :*

$$\frac{|\hat{\rho}^2 - \rho^2|}{N} = \frac{16B^2}{3\pi} N^{1/2} M^{-1/2} \quad (\text{B.4})$$

Thus, while future work may seek to improve the upper bound in theorem B.1, we cannot hope to improve beyond the lower bound formulated in theorem B.2. If we ignore the logarithmic factors to gain intuition, we observe there is (roughly) a gap of  $N^{1/2}$  between the upper and lower bounds. Thus, it is possible that our analysis in appendix B.1 may be conservative in terms of the ambient dimension—specifically, the lower bound only suggests that  $M$  only needs to be increased two-fold to compensate for a two-fold increase in  $N$ . However, in terms of the number of sampled inputs, the rate cannot be improved beyond  $M^{-1/2}$ .

### B.3 Proof of lemma B.1

Here we show that the plug-in estimate of the total variance  $\text{Tr}[\hat{\Sigma}_{ii}]$  converges to the true variance  $\text{Tr}[\Sigma_{ii}]$  exponentially fast as  $M$  increases. We begin with some algebraic manipulations:

$$\begin{aligned} \left| \text{Tr}[\Sigma_{ii} - \hat{\Sigma}_{ii}] \right| &= \left| \text{Tr} \left[ \mathbb{E}_{z \sim P} [h_i(z_m) h_i(z_m)^\top] - \frac{1}{M} \sum_{m=1}^M h_i(z_m) h_i(z_m)^\top \right] \right| \\ &= \left| \mathbb{E}_{z \sim P} [\text{Tr}[h_i(z_m) h_i(z_m)^\top]] - \frac{1}{M} \sum_{m=1}^M \text{Tr}[h_i(z_m) h_i(z_m)^\top] \right| \\ &= \left| \mathbb{E}_{z \sim P} [\text{Tr}[h_i(z_m)^\top h_i(z_m)]] - \frac{1}{M} \sum_{m=1}^M \text{Tr}[h_i(z_m)^\top h_i(z_m)] \right| \\ &= \left| \mathbb{E}_{z \sim P} [h_i(z_m)^\top h_i(z_m)] - \frac{1}{M} \sum_{m=1}^M h_i(z_m)^\top h_i(z_m) \right| \end{aligned}$$

where we have used the property  $\text{Tr}[\mathbf{x}\mathbf{x}^\top] = \mathbf{x}^\top \mathbf{x}$  for any column vector  $\mathbf{x}$  in the last two lines.

The main assumption we are going to make is that the neural responses are constrained to an  $\ell_2$  ball of radius  $B\sqrt{N}$  or equivalently  $h_i(z_m)^\top h_i(z_m) \leq B^2 N$  for all stimuli in the support of  $P$ . Note that this is a reasonable assumption in both biological (energy constraints) and artificial neural networks (weight decay common).

**Lemma B.3** (Bounded Random Variables are Sub-Gaussian, Wainwright [27] Example 2.4). *We say that a random variable  $X$  with mean  $\mu$  is sub-Gaussian with parameter  $\sigma$  if:*

$$\mathbb{E} \left[ e^{\lambda(X-\mu)} \right] \leq e^{\sigma^2 \lambda^2 / 2} \quad \text{for all } \lambda \in \mathbb{R}$$

*Intuitively, this means that the tails of  $X$  fall off faster than a Gaussian. Furthermore, if  $X$  is mean zero and supported on the interval  $[a, b]$ , the  $X$  is sub-Gaussian with parameter  $\sigma = (b - a)/2$ .*

Thus our assumption implies that each term with  $\frac{1}{M} h_i(z_m)^\top h_i(z_m)$  is sub-Gaussian with parameter  $\sigma = B\sqrt{N}/M$ . We can then immediately apply the Hoeffding bound [27, Proposition 2.5] to obtain:

$$\mathbb{P} \left[ \left| \text{Tr}[\Sigma_{ii} - \hat{\Sigma}_{ii}] \right| \geq t \right] \leq 2 \exp \left[ -\frac{Mt^2}{2B^2N} \right] \quad (\text{B.5})$$

Analogously for term (B) we obtain:

$$\mathbb{P} \left[ \left| \text{Tr}[\Sigma_{jj} - \hat{\Sigma}_{jj}] \right| \geq t \right] \leq 2 \exp \left[ -\frac{Mt^2}{2B^2N} \right] \quad (\text{B.6})$$

### B.4 Proof of lemma B.2

Our main tool is the matrix Bernstein inequality, given as theorem 6.1.1 in Tropp [26]. We paraphrase a version of the theorem here to keep our narrative self-contained.

**Theorem B.3** (Matrix Bernstein). *Consider a finite sequence  $\{\mathbf{S}_1, \dots, \mathbf{S}_M\}$  of independent, random  $N \times N$  matrices. Assume that:*

$$\mathbb{E}[\mathbf{S}_m] = \mathbf{0} \quad \text{and} \quad \|\mathbf{S}_m\|_\infty \leq L \quad \text{for each index } m \quad (\text{B.7})$$

where  $\|\mathbf{S}_m\|_\infty = \sup\{\|\mathbf{S}_m \mathbf{v}\|_2 : \|\mathbf{v}\|_2 \leq 1\}$  is the matrix operator norm.

Further, define the variance of the sum  $\sum_m \mathbf{S}_m$  as:

$$V = \left\| \sum_m \mathbb{E} \mathbf{S}_m^\top \mathbf{S}_m \right\|_\infty = \left\| \sum_m \mathbb{E} \mathbf{S}_m \mathbf{S}_m^\top \right\|_\infty \quad (\text{B.8})$$

Then:

$$\mathbb{E} \left[ \left\| \sum_m \mathbf{S}_m \right\|_\infty \right] \leq \sqrt{2V \log(2N)} + \frac{L}{3} \log(2N) \quad (\text{B.9})$$

We now turn to the proof of theorem B.1. Define:

$$\mathbf{S}_m = \frac{1}{M} (h_i(\mathbf{z}_m) h_j(\mathbf{z}_m)^\top - \boldsymbol{\Sigma}_{ij}) \quad (\text{B.10})$$

for the sequence of network inputs  $\{\mathbf{z}_1, \dots, \mathbf{z}_M\}$ . Notice that:

$$\mathbb{E}[\mathbf{S}_m] = \frac{1}{M} (\mathbb{E} [h_i(\mathbf{z}_m) h_j(\mathbf{z}_m)^\top] - \boldsymbol{\Sigma}_{ij}) = \frac{1}{M} (\boldsymbol{\Sigma}_{ij} - \boldsymbol{\Sigma}_{ij}) = \mathbf{0} \quad (\text{B.11})$$

Next, due to triangle inequality, we have:

$$\|\mathbf{S}_m\|_\infty = \frac{1}{M} \|h_i(\mathbf{z}_m) h_j(\mathbf{z}_m)^\top - \boldsymbol{\Sigma}_{ij}\|_\infty \leq \frac{1}{M} \underbrace{\|h_i(\mathbf{z}_m) h_j(\mathbf{z}_m)^\top\|_\infty}_{(1)} + \frac{1}{M} \underbrace{\|\boldsymbol{\Sigma}_{ij}\|_\infty}_{(2)} \quad (\text{B.12})$$

Terms (1) and (2) are each upper bounded by  $B^2 N$ , since for term (1):

$$\|h_i(\mathbf{z}_m) h_j(\mathbf{z}_m)^\top\|_\infty \leq \|h_i(\mathbf{z}_m) h_j(\mathbf{z}_m)^\top \mathbf{v}\|_2 \quad (\text{for any vector } \|\mathbf{v}\|_2 \leq 1) \quad (\text{B.13})$$

$$= h_j(\mathbf{z}_m)^\top \mathbf{v} \|h_i(\mathbf{z}_m)\|_2 \quad (\text{B.14})$$

$$\leq \|h_j(\mathbf{z}_m)\|_2 \|\mathbf{v}\|_2 \|h_i(\mathbf{z}_m)\|_2 \quad (\text{Cauchy-Schwarz inequality}) \quad (\text{B.15})$$

$$\leq B\sqrt{N} \cdot 1 \cdot B\sqrt{N} = B^2 N \quad (\text{From assumptions in eq. 2.1}) \quad (\text{B.16})$$

And for term (2):

$$\|\boldsymbol{\Sigma}_{ij}\|_\infty = \|\mathbb{E} h_i(\mathbf{z}) h_j(\mathbf{z})^\top\|_\infty \quad (\text{B.17})$$

$$\leq \|\mathbb{E} h_i(\mathbf{z}_m) h_j(\mathbf{z}_m)^\top \mathbf{v}\|_2 \quad (\text{for any vector } \|\mathbf{v}\|_2 \leq 1) \quad (\text{B.18})$$

$$\leq \mathbb{E} \|h_i(\mathbf{z}_m) h_j(\mathbf{z}_m)^\top \mathbf{v}\|_2 \quad (\text{Jensen's inequality}) \quad (\text{B.19})$$

$$\leq B^2 N \quad (\text{Repeat the upper bound on term 1}) \quad (\text{B.20})$$

To summarize, we have:

$$\|\mathbf{S}_m\|_\infty \leq \frac{1}{M} \|h_i(\mathbf{z}_m) h_j(\mathbf{z}_m)^\top\|_\infty + \frac{1}{M} \|\boldsymbol{\Sigma}_{ij}\|_\infty \leq \frac{2B^2 N}{M} \quad (\text{B.21})$$

That is, we have shown that the assumptions of eq. (B.7) are satisfied with  $L = 2B^2 N/M$ .

Our next task is to determine an expression for the variance  $V$  defined in eq. (B.8). First, we have:

$$\begin{aligned} \mathbb{E} \mathbf{S}_m^\top \mathbf{S}_m &= \frac{1}{M^2} \mathbb{E} [h_j(\mathbf{z}_m) h_i(\mathbf{z}_m)^\top h_i(\mathbf{z}_m) h_j(\mathbf{z}_m)^\top + \boldsymbol{\Sigma}_{ij}^\top \boldsymbol{\Sigma}_{ij} - \boldsymbol{\Sigma}_{ij}^\top h_j(\mathbf{z}_m) h_i(\mathbf{z}_m)^\top - h_i(\mathbf{z}_m) h_j(\mathbf{z}_m)^\top \boldsymbol{\Sigma}_{ij}] \\ &= \frac{1}{M^2} \mathbb{E} [h_j(\mathbf{z}_m) h_i(\mathbf{z}_m)^\top h_i(\mathbf{z}_m) h_j(\mathbf{z}_m)^\top] + \boldsymbol{\Sigma}_{ij}^\top \boldsymbol{\Sigma}_{ij} - \boldsymbol{\Sigma}_{ij}^\top \mathbb{E} [h_j(\mathbf{z}_m) h_i(\mathbf{z}_m)^\top] - \mathbb{E} [h_i(\mathbf{z}_m) h_j(\mathbf{z}_m)^\top] \boldsymbol{\Sigma}_{ij} \\ &= \frac{1}{M^2} \mathbb{E} [h_j(\mathbf{z}_m) h_i(\mathbf{z}_m)^\top h_i(\mathbf{z}_m) h_j(\mathbf{z}_m)^\top] + \boldsymbol{\Sigma}_{ij}^\top \boldsymbol{\Sigma}_{ij} - \boldsymbol{\Sigma}_{ij}^\top \boldsymbol{\Sigma}_{ij} - \boldsymbol{\Sigma}_{ij}^\top \boldsymbol{\Sigma}_{ij} \\ &= \frac{1}{M^2} \mathbb{E} [h_j(\mathbf{z}_m) h_i(\mathbf{z}_m)^\top h_i(\mathbf{z}_m) h_j(\mathbf{z}_m)^\top] - \boldsymbol{\Sigma}_{ij}^\top \boldsymbol{\Sigma}_{ij} \end{aligned}$$

Then, by triangle inequality:

$$\begin{aligned}\|\mathbb{E} \mathbf{S}_m^\top \mathbf{S}_m\|_\infty &= \frac{1}{M^2} \|\mathbb{E}[h_j(\mathbf{z}_m)h_i(\mathbf{z}_m)^\top h_i(\mathbf{z}_m)h_j(\mathbf{z}_m)^\top] - \boldsymbol{\Sigma}_{ij}^\top \boldsymbol{\Sigma}_{ij}\|_\infty \\ &\leq \frac{1}{M^2} \underbrace{\|\mathbb{E}[h_j(\mathbf{z}_m)h_i(\mathbf{z}_m)^\top h_i(\mathbf{z}_m)h_j(\mathbf{z}_m)^\top]\|_\infty}_{(A)} + \frac{1}{M^2} \underbrace{\|\boldsymbol{\Sigma}_{ij}^\top \boldsymbol{\Sigma}_{ij}\|_\infty}_{(B)}\end{aligned}$$

Terms (A) and (B) are each upper bounded by  $N^2$ . First, taking term (A):

$$\begin{aligned}\|\mathbb{E}[h_j(\mathbf{z}_m)h_i(\mathbf{z}_m)^\top h_i(\mathbf{z}_m)h_j(\mathbf{z}_m)^\top]\|_\infty &\leq \|\mathbb{E}[h_j(\mathbf{z}_m)h_i(\mathbf{z}_m)^\top h_i(\mathbf{z}_m)h_j(\mathbf{z}_m)^\top \mathbf{v}]\|_2 \quad (\text{for } \|\mathbf{v}\| \leq 1) \\ &\leq \mathbb{E} \|h_j(\mathbf{z}_m)h_i(\mathbf{z}_m)^\top h_i(\mathbf{z}_m)h_j(\mathbf{z}_m)^\top \mathbf{v}\|_2 \quad (\text{Jensen's}) \\ &\leq \mathbb{E} [h_j(\mathbf{z}_m)^\top \mathbf{v} \|h_i(\mathbf{z}_m)\|_2^2 \|h_j(\mathbf{z}_m)\|_2] \\ &\leq \mathbb{E} [\|\mathbf{v}\|_2 \|h_i(\mathbf{z}_m)\|_2^2 \|h_j(\mathbf{z}_m)\|_2^2] \quad (\text{Cauchy-Schwarz}) \\ &\leq 1 \cdot B^2 N \cdot B^2 N = B^4 N^2 \quad (\text{from eq. 2.1})\end{aligned}$$

For term (B), we first note that  $\|\boldsymbol{\Sigma}_{ij}^\top \boldsymbol{\Sigma}_{ij}\|_\infty \leq \|\boldsymbol{\Sigma}_{ij}\|_\infty^2$  due to the fact that the operator norm is submultiplicative. Then, term (B) is upper bounded by  $B^4 N^2$  follows readily from:

$$\begin{aligned}\|\boldsymbol{\Sigma}_{ij}\|_\infty &= \|\mathbb{E} h_i(\mathbf{z})h_j(\mathbf{z})^\top\|_\infty \\ &\leq \|\mathbb{E} h_i(\mathbf{z})h_j(\mathbf{z})^\top \mathbf{v}\|_2 \quad (\text{for } \|\mathbf{v}\| \leq 1) \\ &\leq \mathbb{E} \|h_i(\mathbf{z})h_j(\mathbf{z})^\top \mathbf{v}\|_2 \quad (\text{Jensen's}) \\ &\leq \mathbb{E} \|h_i(\mathbf{z})\|_2 \|h_j(\mathbf{z})\|_2 \|\mathbf{v}\|_2 \quad (\text{Cauchy-Schwarz}) \\ &\leq B\sqrt{N} \cdot B\sqrt{N} \cdot 1 = B^2 N \quad (\text{from eq. 1})\end{aligned}$$

Taking these two bounds together, we have shown  $\|\mathbb{E} \mathbf{S}_m^\top \mathbf{S}_m\|_\infty \leq 2B^4 N^2/M^2$ . We are now ready to upper bound the variance term,  $V$ , appearing in theorem B.3. Specifically, by the triangle inequality and the bounds above, we have:

$$V = \|\sum_m \mathbb{E} \mathbf{S}_m^\top \mathbf{S}_m^\top\|_\infty \leq \sum_{m=1}^M \|\mathbb{E} \mathbf{S}_m^\top \mathbf{S}_m^\top\|_\infty \leq \frac{2B^4 N^2}{M} \quad (\text{B.22})$$

With this, we are equipped to apply the matrix Bernstein inequality to obtain an upper bound on the estimation error of the plug-in estimator. Specifically, we have:

$$\begin{aligned}\left| \|\hat{\boldsymbol{\Sigma}}_{ij}\|_* - \|\boldsymbol{\Sigma}_{ij}\|_* \right| &\leq \|\hat{\boldsymbol{\Sigma}}_{ij} - \boldsymbol{\Sigma}_{ij}\|_* \quad (\text{reverse triangle inequality}) \\ &= \left\| \sum_m \mathbf{S}_m \right\|_* \\ &\leq N \left\| \sum_m \mathbf{S}_m \right\|_\infty \\ &\leq N \sqrt{2V \log(2N)} + \frac{NL}{3} \log(2N) \quad (\text{theorem B.3}) \\ &\leq 2B^2 N^2 M^{-1/2} \sqrt{\log(2N)} + \frac{2B^2 N^2}{3M} \log(2N)\end{aligned}$$

Where we have substituted the derived quantities  $L = 2B^2 N/M$  and  $V \leq 2B^4 N^2/M$  in the final line.

## B.5 Proof of theorem B.1

Lemma B.2 provides an upper bound on the expected value on  $\left| \|\boldsymbol{\Sigma}_{ij}\|_* - \|\hat{\boldsymbol{\Sigma}}_{ij}\|_* \right|$ , which is the error of our plug-in estimate of cross-covariance nuclear norm. This bound holds for any true cross-covariance matrix  $\boldsymbol{\Sigma}_{ij}$ , provided that the constraints in eq. (2.1) are satisfied. However, this tells us nothing about how the estimation error deviates around its expectation.

Here, we use the bounded differences inequality [27, Corollary 2.21], also called McDiarmid's inequality, to show that deviations around this expectation decrease exponentially fast. Thus, the upper bound on the expected error (theorem B.1) provides accurate intuition.

**Lemma B.4** (Bounded Differences Inequality, Wainwright [27] Corollary 2.21). *Consider a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ . The function is said to have the bounded difference property for the  $k$ th coordinate if there exists an  $L_k$  for which the following holds:*

$$\max_{X_{1:n} \in \mathbb{R}^n, X'_k \in \mathbb{R}} |f(X_{1:n}) - f(X_{1:k-1}, X'_k, X_{k+1:n})| \leq L_k$$

Suppose  $f$  satisfies this property with  $L_1, \dots, L_n$  for each coordinate respectively. Then the following inequality holds:

$$\mathbb{P} \left[ \left| f(X_{1:n}) - \mathbb{E}[f(X_{1:n})] \right| \geq t \right] \leq \exp \left[ -\frac{2t^2}{\sum_{i=1}^n L_i^2} \right] \quad (\text{B.23})$$

We start by applying the reverse triangle inequality:

$$\left| \|\Sigma_{ij}\|_* - \|\hat{\Sigma}_{ij}\|_* \right| \leq \|\Sigma_{ij} - \hat{\Sigma}_{ij}\|_* = \left\| \Sigma_{ij} - \frac{1}{M} \sum_{m=1}^M h_i(\mathbf{z}_m) h_j(\mathbf{z}_m)^\top \right\|_*$$

We can bound how much this changes if we change one coordinate of the function, i.e. if  $h_i(\mathbf{z}_1)^\top h_j(\mathbf{z}_1)$  is replaced by  $h_i(\tilde{\mathbf{z}}_1)^\top h_j(\tilde{\mathbf{z}}_1)$ . The difference is then bounded by:

$$\begin{aligned} & \left\| \Sigma_{ij} - \frac{1}{M} \sum_{m=1}^M h_i(\mathbf{z}_m) h_j(\mathbf{z}_m)^\top \right\|_* - \left\| \Sigma_{ij} - \left( \frac{1}{M} \sum_{m=1}^M h_i(\mathbf{z}_m) h_j(\mathbf{z}_m)^\top - \frac{1}{M} h_i(\mathbf{z}_1) h_j(\mathbf{z}_1)^\top + \frac{1}{M} h_i(\tilde{\mathbf{z}}_1) h_j(\tilde{\mathbf{z}}_1)^\top \right) \right\|_* \\ & \leq \left\| \frac{1}{M} (h_i(\mathbf{z}_1) h_j(\mathbf{z}_1)^\top - h_i(\tilde{\mathbf{z}}_1) h_j(\tilde{\mathbf{z}}_1)^\top) \right\|_* = \frac{1}{M} \|h_i(\mathbf{z}_1) h_j(\mathbf{z}_1)^\top - h_i(\tilde{\mathbf{z}}_1) h_j(\tilde{\mathbf{z}}_1)^\top\|_* \\ & \leq \frac{1}{M} (\|h_i(\mathbf{z}_1) h_j(\mathbf{z}_1)^\top\|_* + \|h_i(\tilde{\mathbf{z}}_1) h_j(\tilde{\mathbf{z}}_1)^\top\|_*) = \frac{1}{M} (|h_i(\mathbf{z}_1)^\top h_j(\mathbf{z}_1)| + |h_i(\tilde{\mathbf{z}}_1)^\top h_j(\tilde{\mathbf{z}}_1)|) \end{aligned}$$

Finally, we can apply Cauchy-Schwartz and our assumption about the neural activations being bounded to obtain:

$$\begin{aligned} \frac{1}{M} (|h_i(\mathbf{z}_1)^\top h_j(\mathbf{z}_1)| + |h_i(\tilde{\mathbf{z}}_1)^\top h_j(\tilde{\mathbf{z}}_1)|) & \leq \frac{1}{M} (\|h_i(\mathbf{z}_1)\|_2 \|h_j(\mathbf{z}_1)\|_2 + \|h_i(\tilde{\mathbf{z}}_1)\|_2 \|h_j(\tilde{\mathbf{z}}_1)\|_2) \\ & \leq \frac{2B^2 N}{M} \end{aligned}$$

Thus we have  $\sum_{i=1}^M L_i^2 = \sum_{i=1}^M 4B^4 N^2 / M^2 = 4B^4 N^2 / M$ , and we can apply the bounded differences inequality to obtain for all  $t \geq 0$ :

$$\mathbb{P} \left[ \left| \|\Sigma_{ij}\|_* - \|\hat{\Sigma}_{ij}\|_* \right| - \mathbb{E} \left| \|\Sigma_{ij}\|_* - \|\hat{\Sigma}_{ij}\|_* \right| \geq t \right] \leq 2 \exp \left[ -\frac{Mt^2}{2B^4 N^2} \right] \quad (\text{B.24})$$

For the deviation from the expectation to be in the range  $[-t, t]$  with probability  $1 - \delta$  we require:

$$2 \exp \left[ -\frac{Mt^2}{2B^4 N^2} \right] \leq \delta$$

Solving for  $t$  gives  $t \geq B^2 N M^{-1/2} \sqrt{2 \log(2/\delta)}$ , and thus with probability  $1 - \delta$  the following holds:

$$\left| \|\Sigma_{ij}\|_* - \|\hat{\Sigma}_{ij}\|_* - \mathbb{E} \left| \|\Sigma_{ij}\|_* - \|\hat{\Sigma}_{ij}\|_* \right| \right| \leq B^2 N M^{-1/2} \sqrt{2 \log(2/\delta)}$$

To proceed we break this we use a basic identity of the absolute value: if  $|a - b| < c$  then  $a - b < c$  and also  $b - a < c$ . Thus, with probability at least  $1 - \delta$ , we have:

$$\begin{aligned} \|\Sigma_{ij}\|_* - \|\hat{\Sigma}_{ij}\|_* & \leq \mathbb{E} \left| \|\Sigma_{ij}\|_* - \|\hat{\Sigma}_{ij}\|_* \right| + \frac{B^2 N}{M^{1/2}} \sqrt{2 \log(2/\delta)} \\ & \leq \frac{2B^2 N^2}{M^{1/2}} \sqrt{\log(2N)} + \frac{2B^2 N^2}{3M} \log(2N) + \frac{B^2 N}{M^{1/2}} \sqrt{2 \log(2/\delta)} \end{aligned}$$

And we also have with probability at least  $1 - \delta$ , we have:

$$\begin{aligned} \|\hat{\Sigma}_{ij}\|_* - \|\Sigma_{ij}\|_* &\leq \mathbb{E} \left| \|\Sigma_{ij}\|_* - \|\hat{\Sigma}_{ij}\|_* \right| + \frac{B^2 N}{M^{1/2}} \sqrt{2 \log(2/\delta)} \\ &\leq \frac{2B^2 N^2}{M^{1/2}} \sqrt{\log(2N)} + \frac{2B^2 N^2}{3M} \log(2N) + \frac{B^2 N}{M^{1/2}} \sqrt{2 \log(2/\delta)} \end{aligned}$$

In the final inequalities above, we have simply plugged in our expectation bound from lemma B.2. The relations above imply that the following holds with probability  $1 - \delta$ :

$$\left| \|\Sigma_{ij}\|_* - \|\hat{\Sigma}_{ij}\|_* \right| \leq \frac{2N^2 \log(2B^2 N)}{3M} + \frac{2B^2 N^2 \sqrt{\log(2N)}}{M^{1/2}} + \frac{B^2 N}{M^{1/2}} \sqrt{2 \log\left(\frac{2}{\delta}\right)} \quad (\text{B.25})$$

To complete the proof we need to combine the above tail bound with lemma B.1. By the triangle inequality we have

$$\begin{aligned} |\hat{\rho}^2 - \rho^2| &= \left| \text{Tr}[\hat{\Sigma}_{ii}] + \text{Tr}[\hat{\Sigma}_{jj}] - 2\|\hat{\Sigma}_{ij}\|_* - \text{Tr}[\hat{\Sigma}_{ii}] - \text{Tr}[\hat{\Sigma}_{jj}] + 2\|\hat{\Sigma}_{ij}\|_* \right| \\ &= \left| \text{Tr}[\hat{\Sigma}_{ii}] - \text{Tr}[\Sigma_{ii}] + \text{Tr}[\hat{\Sigma}_{jj}] - \text{Tr}[\Sigma_{jj}] + 2\|\Sigma_{ij}\|_* - 2\|\hat{\Sigma}_{ij}\|_* \right| \\ &\leq \left| \text{Tr}[\hat{\Sigma}_{ii}] - \text{Tr}[\Sigma_{ii}] \right| + \left| \text{Tr}[\hat{\Sigma}_{jj}] - \text{Tr}[\Sigma_{jj}] \right| + 2 \left| \|\Sigma_{ij}\|_* - \|\hat{\Sigma}_{ij}\|_* \right| \end{aligned}$$

Setting  $\delta' = \delta/3$  in our results for these three terms yields that the following three inequalities independently hold with probability  $\delta/3$ :

$$\begin{aligned} \left| \text{Tr}[\Sigma_{ii}] - \text{Tr}[\hat{\Sigma}_{ii}] \right| &\geq BN^{1/2} M^{-1/2} \sqrt{2 \log(6/\delta)} \\ \left| \text{Tr}[\Sigma_{jj}] - \text{Tr}[\hat{\Sigma}_{jj}] \right| &\geq BN^{1/2} M^{-1/2} \sqrt{2 \log(6/\delta)} \\ \left| \|\Sigma_{ij}\|_* - \|\hat{\Sigma}_{ij}\|_* \right| &\geq \frac{2N^2 \log(2B^2 N)}{3M} + \frac{2B^2 N^2 \sqrt{\log(2N)}}{M^{1/2}} + \frac{B^2 N}{M^{1/2}} \sqrt{2 \log\left(\frac{6}{\delta}\right)} \end{aligned}$$

By applying the union bound, we obtain that all three inequalities hold simultaneously with probability  $\leq \delta/3 + \delta/3 + \delta/3 = \delta$ . The three reverse inequalities then hold simultaneously with probability greater than or equal to  $1 - \delta$ . Thus with probability at least  $1 - \delta$ , the following holds:

$$|\hat{\rho}^2 - \rho^2| \leq \frac{2B^2 N^2 \log(2N)}{3M} + \frac{2B^2 N^2 \sqrt{\log(2N)}}{M^{1/2}} + \left( \frac{NB^2}{M^{1/2}} + \frac{2N^{1/2}B}{M^{1/2}} \right) \sqrt{2 \log\left(\frac{6}{\delta}\right)}$$

as claimed in theorem B.1.

## B.6 Proof of theorem B.2 (Lower Bound on Plug-In Estimator Error)

We derive a lower bound by constructing an explicit example where the plug-in estimator performs badly. Specifically, we consider a scenario where two networks have entirely decorrelated, high-variance representations. To do this, we use *Rademacher random variables*—a random variable  $R$  is called a Rademacher variable if it behaves as follows:

$$R = \begin{cases} +1 & \text{with probability } 1/2 \\ -1 & \text{with probability } 1/2 \end{cases} \quad (\text{B.26})$$

Now, suppose we sample  $M$  network inputs,  $z_1, \dots, z_M \sim P$ , independently. Further, let  $B > 0$  be the constant appearing in eq. (2.1). For  $m \in \{1, \dots, M\}$  define

$$X_m = \frac{1}{B} h_i(z_m) \quad \text{and} \quad Y_m = \frac{1}{B} h_j(z_m) \quad (\text{B.27})$$

Note that  $X_m$  and  $Y_m$  are  $N$ -dimensional random vectors. Due to eq. (2.1), we have  $\|h_i(\mathbf{z})\|_2 \leq B\sqrt{N}$  and  $\|h_j(\mathbf{z})\|_2 \leq B\sqrt{N}$  almost surely. Thus,  $\|X_m\| \leq \sqrt{N}$  and  $\|Y_m\| \leq \sqrt{N}$  almost surely.

Define  $X = (1/B)h_i(\mathbf{z})$  and  $Y = (1/B)h_j(\mathbf{z})$  for randomly sampled  $\mathbf{z} \sim P$ . The case we will consider is that  $X$  and  $Y$  are each composed of  $N$  independent Rademacher variables. One trivial way to construct this is to suppose each  $\mathbf{z} \sim P$  is a random vector with  $2N$  elements, all of which are independent Rademacher variables scaled by a factor  $B > 0$ . Then, let  $h_i : \mathbb{R}^{2N} \mapsto \mathbb{R}^N$  be the function which extracts the first  $N$  elements of  $\mathbf{z}$  and let  $h_j : \mathbb{R}^{2N} \mapsto \mathbb{R}^N$  be the function which extracts the final  $N$  elements.

Thus, we have constructed a setting where  $X_1, \dots, X_M, Y_1, \dots, Y_M$  are all composed of independent Rademacher variables. In this setting, the squared Procrustes distance is given by:

$$\rho^2 = \text{Tr}[\mathbf{\Sigma}_{ii}] + \text{Tr}[\mathbf{\Sigma}_{jj}] - 2\|\mathbf{\Sigma}_{ij}\|_* \quad (\text{B.28})$$

$$= \text{Tr}[\mathbb{E}[h_i(\mathbf{z})h_i(\mathbf{z})^\top]] + \text{Tr}[\mathbb{E}[h_j(\mathbf{z})h_j(\mathbf{z})^\top]] - 2\|\mathbb{E}[h_i(\mathbf{z})h_j(\mathbf{z})^\top]\|_* \quad (\text{B.29})$$

$$= B^2 \cdot (\text{Tr}[\mathbb{E}[XX^\top]] + \text{Tr}[\mathbb{E}[YY^\top]] - 2\|\mathbb{E}[XY^\top]\|_*) \quad (\text{B.30})$$

$$= B^2 \cdot (\mathbb{E}[X^\top X] + \mathbb{E}[Y^\top Y] - 2\|\mathbb{E}[X]\mathbb{E}[Y^\top]\|_*) \quad (\text{B.31})$$

$$= B^2 \cdot (N + N - 0) \quad (\text{B.32})$$

$$= 2B^2N \quad (\text{B.33})$$

where we have used the fact that  $X$  and  $Y$  are independent, mean zero, random vectors to conclude that the cross covariance is an  $N \times N$  matrix filled with zeros. Furthermore, note that  $X_m^\top X_m = N$  and  $Y_m^\top Y_m = N$  almost surely for all  $m \in 1, \dots, M$  since they are comprised of  $N$  Rademacher variables. Thus, the plug-in estimate of the squared Procrustes distance takes the form:

$$\hat{\rho}^2 = B^2 \cdot (\text{Tr}[\frac{1}{M} \sum_m X_m X_m^\top] + \text{Tr}[\frac{1}{M} \sum_m Y_m Y_m^\top] - 2\|\frac{1}{M} \sum_m X_m Y_m^\top\|_*) \quad (\text{B.34})$$

$$= B^2 \cdot (\frac{1}{M} \sum_m X_m^\top X_m + \frac{1}{M} \sum_m Y_m^\top Y_m - 2\|\frac{1}{M} \sum_m X_m Y_m^\top\|_*) \quad (\text{B.35})$$

$$= B^2 \cdot (N + N - 2\|\frac{1}{M} \sum_m X_m Y_m^\top\|_*) \quad (\text{B.36})$$

$$= 2B^2N - 2B^2\|\frac{1}{M} \sum_m X_m Y_m^\top\|_* \quad (\text{B.37})$$

Putting these two results together, we conclude that the absolute error of the plug-in estimator is:

$$|\rho^2 - \hat{\rho}^2| = 2B^2\|\frac{1}{M} \sum_m X_m Y_m^\top\|_* \quad (\text{B.38})$$

Now, the product of two independent Rademacher variables is also a standard Rademacher variable. Thus, each element inside the matrix  $(1/M) \sum_m X_m Y_m^\top$ , is the empirical average of  $M$  independent Rademacher variables. These matrix elements are asymptotically independent in the limit that  $M \rightarrow \infty$ . Further, the central limit theorem applies in this limit, and thus the distribution of each matrix element approaches a Gaussian distribution  $\mathcal{N}(0, 1/M)$ .

Such random matrices are well-studied under the name of Ginibre ensembles. In the limit that  $N \rightarrow \infty$  and the variance of each matrix element is taken to be  $\sigma^2/N$ , the density of the singular values takes the following form [see e.g. 20, sec. 3.1.3]:

$$\rho(s) = \frac{\sqrt{4\sigma^2 - s^2}}{\pi\sigma^2} \quad s \in (0, 2\sigma) \quad (\text{B.39})$$

This is called the quarter circle law since if we look at the density of  $s$  it forms a quarter circle. The nuclear norm of the matrix is  $N$  times the expected value of  $s$  with respect to the density  $\rho(s)$ . Integrating this density, we obtain:

$$\lim_{\substack{N \rightarrow \infty \\ M \gg N}} \|\frac{1}{M} \sum_m X_m Y_m^\top\|_* = \frac{N}{\pi\sigma^2} \int_0^{2\sigma} s \sqrt{4\sigma^2 - s^2} ds \quad (\text{B.40})$$

$$= \frac{N}{4\pi\sigma^2} \left[ -\frac{1}{3}(4\sigma^2 - s^2)^{3/2} \right]_0^{2\sigma} \quad (\text{B.41})$$

$$= \frac{N}{\pi\sigma^2} \left[ \frac{1}{3}(4\sigma^2)^{3/2} \right] = \frac{N}{\pi\sigma^2} \left[ \frac{8}{3}\sigma^3 \right] \quad (\text{B.42})$$

$$= \frac{8\sigma}{3\pi} N = \frac{8}{3\pi} N^{3/2} M^{-1/2} \quad (\text{B.43})$$

Where in the last line we have substituted  $\sigma = \sqrt{N/M}$ , which comes from equating  $\sigma^2/N$  (the variance in of each matrix element in eq. B.39) with  $1/M$  (the variance given by the average of  $M$  Rademacher variables under the central limit theorem). Note that the analysis above holds asymptotically as  $M, N \rightarrow \infty$  and we keep  $M \gg N$  so that the central limit theorem continues to hold.

Plugging eq. (B.43) into eq. (B.38) and dividing both sides by  $N$  we arrive at the expression appearing in theorem B.2.

## C Appendix: Method-of-Moments Estimator

### C.1 Derivation of method-of-moment estimator

We now turn to constructing our method-of-moments estimator of  $\|\Sigma_{ij}\|_* = \sum_{n=1}^N s_n(\Sigma_{ij})$ , which is required for our novel estimator of the Riemannian shape distance. We can form an unbiased estimator of the matrix  $\Sigma_{ij}$  by observing a single random stimuli in the two networks:

$$\hat{\Sigma}_{ijm} := h_i(\mathbf{z}_m)h_j(\mathbf{z}_m)^\top \in \mathbb{R}^{N \times N}, \quad \mathbb{E}[\hat{\Sigma}_{ijm}] = \Sigma_{ij}$$

Note that here the randomness comes from the selection of the stimuli, i.e.  $\mathbf{z}_m \sim P$ ; the output of the network is deterministic. Assuming  $m, m'$  are distinct stimuli drawn independently from the distribution  $P$ , we then have:

$$\mathbb{E} \left[ \hat{\Sigma}_{ijm} \hat{\Sigma}_{ijm'}^\top \right] = \Sigma_{ij} \Sigma_{ij}^\top$$

This means we can estimate  $\Sigma_{ij} \Sigma_{ij}^\top$  by observing a pair of stimuli in both networks.

$$\begin{aligned} \text{Tr} \left[ f(\Sigma_{ij} \Sigma_{ij}^\top) \right] &= \sum_{n=1}^N f(s_n^2(\Sigma_{ij})) = \sum_{n=1}^N \sum_{p=0}^{\infty} \gamma_p s_n^{2p}(\Sigma_{ij}) && \text{Taylor expansion of } f(\cdot) \\ &= \sum_{p=0}^{\infty} \gamma_p \sum_{n=1}^N s_n^{2p}(\Sigma_{ij}) = \sum_{p=0}^{\infty} \gamma_p \text{Tr} \left[ \left( \Sigma_{ij} \Sigma_{ij}^\top \right)^p \right] && \text{Tr} \left[ \left( \Sigma_{ij} \Sigma_{ij}^\top \right)^p \right] = \sum_{n=1}^N s_n^{2p}(\Sigma_{ij}) \\ &= \sum_{p=0}^{\infty} \gamma_p \mathbb{E} \left[ \text{Tr} \left[ \prod_{\sigma=1}^p \hat{\Sigma}_{ij(2\sigma-1)} \hat{\Sigma}_{ij(2\sigma)}^\top \right] \right] && \text{Substitute unbiased estimator for } \left( \Sigma_{ij} \Sigma_{ij}^\top \right)^p \\ &\approx \sum_{p=0}^P \gamma_p \mathbb{E} \left[ \text{Tr} \left[ \prod_{\sigma=1}^p \hat{\Sigma}_{ij(2\sigma-1)} \hat{\Sigma}_{ij(2\sigma)}^\top \right] \right] && \text{Approximate with truncated power series} \end{aligned}$$

Our estimator for the nuclear norm of  $\Sigma_{ij}$  is thus:

$$\widehat{\|\Sigma_{ij}\|_*} = \sum_{p=0}^P \gamma_p \text{Tr} \left[ \prod_{\sigma=1}^p \hat{\Sigma}_{ij(2\sigma-1)} \hat{\Sigma}_{ij(2\sigma)}^\top \right] \quad (\text{C.1})$$

Note that for each element of the product we are considering the estimator based on stimuli  $(2\sigma - 1)$  and  $(2\sigma)$ ; in total this estimator will use  $2P$  unique stimuli.

### C.2 Deriving the Quadratic Program

The optimization problem in eq. (2.13) takes the form:

$$\underset{\gamma}{\text{minimize}} \quad \gamma^\top \mathbf{A} \gamma + N^2 \left( \max_x f^2(\gamma, x) \right) \quad (\text{C.2})$$

where  $f(\gamma, x) = x^{1/2} - \sum_p \gamma_p x^p$ ,

$$\gamma = \begin{bmatrix} \gamma_1 \\ \vdots \\ \gamma_P \end{bmatrix} \in \mathbb{R}^P, \quad \mathbf{A} = \begin{bmatrix} \text{Cov}(\hat{W}_1, \hat{W}_1) & \dots & \text{Cov}(\hat{W}_1, \hat{W}_P) \\ \vdots & & \vdots \\ \text{Cov}(\hat{W}_P, \hat{W}_1) & \dots & \text{Cov}(\hat{W}_P, \hat{W}_P) \end{bmatrix} \in \mathbb{R}^{P \times P}, \quad (\text{C.3})$$

Notice that  $f$  is linear in  $\gamma$ , and that  $\mathbf{A}$  is symmetric, positive-definite.

We will reformulate eq. (C.2) in several steps, and ultimately obtain a quadratic program that can be efficiently solved. First, we introduce a new optimization variable  $u \in \mathbb{R}$  whose square is an upper bound on  $f^2(\gamma, x)$  for all  $x \in [0, 1]$ . Thus, the optimal  $\gamma$  for the problem:

$$\begin{aligned} & \underset{\gamma, u}{\text{minimize}} && \gamma^\top \mathbf{A} \gamma + N^2 u^2 \\ & \text{subject to} && u^2 \geq f^2(\gamma, x) \quad \text{for all } x \in [0, 1] \end{aligned} \quad (\text{C.4})$$

coincides to the optimal  $\gamma$  solving eq. (C.2). This is essentially an *epigraph reformulation* of the original problem [see 3, equation 4.11]. Notice that the objective function is quadratic in this reformulation.

Next, we lay down a fine grid of linearly spaced test points  $x_1, \dots, x_T \in [0, 1]$ . We can then obtain a good approximation to the solution in eq. (C.4) by solving:

$$\begin{aligned} & \underset{\gamma, u}{\text{minimize}} && \gamma^\top \mathbf{A} \gamma + N^2 u^2 \\ & \text{subject to} && u^2 \geq f^2(\gamma, x_t) \quad \text{for all } t \in 1, \dots, T \end{aligned} \quad (\text{C.5})$$

Of course, increasing  $T$  (the number of test points) improves the approximation arbitrarily well.

Finally, the constraints of the problem can be put into a form that is jointly linear in  $\gamma$  and  $u$ . First, constraining  $u^2 \geq f^2(\gamma, x_t)$  is equivalent to simultaneously constraining  $u \geq f(\gamma, x_t)$  and  $u \geq -f(\gamma, x_t)$ . Then, plugging in the definition of  $f(\gamma, x_t)$ , and rearranging we have:

$$\begin{aligned} & \underset{\gamma, u}{\text{minimize}} && \gamma^\top \mathbf{A} \gamma + N^2 u^2 \\ & \text{subject to} && u + \sum_p \gamma_p x_t^p \geq x_t^{1/2} \quad \text{for all } t \in 1, \dots, T \\ & && u - \sum_p \gamma_p x_t^p \geq -x_t^{1/2} \quad \text{for all } t \in 1, \dots, T \end{aligned} \quad (\text{C.6})$$

This objective is quadratic and the constraints are linear with respect to the optimized quantities. Thus, a solution (approximated to high accuracy) can be achieved efficiently using off-the-shelf quadratic programming solvers. To enforce the user defined bound on the bias a final two constraints are appended to eq. (C.6):  $-Nu \geq -c$  and  $Nu \geq -c$ , where  $c$  is the upper bound on the absolute bias.

### C.3 Confidence intervals

To form approximate  $\alpha$  level confidence intervals around  $\|\widehat{\Sigma}_{ij}\|_*$  we use the maximal bias (eq. 2.13, term 1) and variance (eq. 2.13, term 2) from the quadratic program's solution:

$$\left[ \|\widehat{\Sigma}_{ij}\|_* - z^* \sqrt{\gamma^\top \mathbf{A} \gamma} - Nu, \quad \|\widehat{\Sigma}_{ij}\|_* + z^* \sqrt{\gamma^\top \mathbf{A} \gamma} + Nu \right],$$

where  $z^*$  is the critical value of the standard normal. For confidence intervals of the similarity score we scale this interval by the denominator of the similarity score.

## D Appendix: Extended Experiments

### D.1 Control of estimator bias

Here we demonstrate the bias-variance tradeoff controlled by the upper-bound on bias defined by the user. The quadratic program in eq. (2.13) can be constrained to keep the maximal absolute bias below

a chosen constant (Fig. 3A, extent of blue shaded area centered around true similarity score). The actual maximal bias for a given solution to the program will then be less than or equal to the user defined bias (cyan shaded area within blue). The expected value of the moment estimator stays within the maximal bias, in this case on its bound (orange trace mean and SD across 5,000 simulations). The user defined bias bound remains inactive until it is less than the MSE minimizing solution’s bias (blue shaded area completely overlapped by cyan when user defined bias bound is less than 0.2). Variance then begins to increase as higher order  $\hat{W}_p$  terms are weighted more heavily to reduce bias (orange standard deviation bars from simulation increase as cyan region narrows). The expected value of the estimator converges to ground truth as it is constrained by the bias bound (dotted orange line converges to dashed black). The plug-in estimator exceeds the maximal bias of the moment estimator (blue trace outside of cyan shaded area).

Intuition for the method-of-moments estimator can be drawn from example plots of solutions to the power series approximation to the square root (eq. 2.11, Fig. 3B, orange trace approximates black dashed trace) of the squared singular values of  $\Sigma_{1,2}$  (black points all overlapping). Here we have re-scaled the singular values on the vertical axis so that the deviation between the square root and power series approximation is exactly the bias of the moment estimator. In the case where bias is not constrained (associated with left most estimates in panel B) the approximation is poor (dashed-dot orange trace does not match dashed black trace). For these eigenvalues the the deviation is near the worst possible bias (distance from black point dashed dot orange line is nearly as far as any other vertical deviation between the traces), this is why the estimator in panel B sits at the bound of maximal possible bias. On the other hand when the upper bound on bias is very small (far right of B) the approximation is very good (dashed orange overlaps dashed black) because higher order terms are used. Yet this results in very high variance (Fig. 3B).

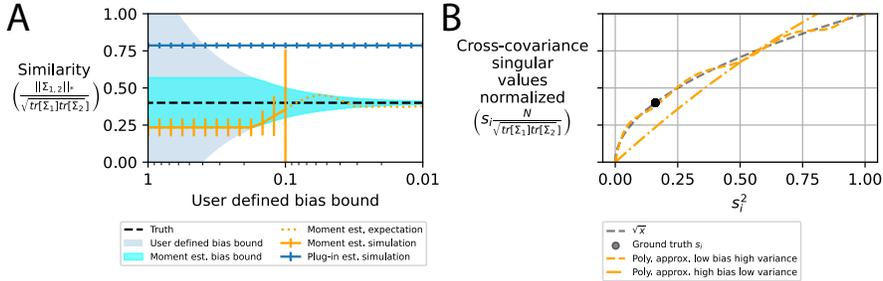


Figure 3: Control of bias-variance tradeoff with user defined bound on bias. **(A)** Here the moment based estimator is constrained to be within the user defined bias bound (blue region) and to minimize worst case MSE (eq. (2.13)). Maximal bias can be less than the user defined bias (cyan region within blue). As the estimator is constrained to have less bias variance increases (orange trace converges to black dashed as SD bars widen as ). Where simulations become unstable we plot the theoretical expected value (dotted orange). Plug-in estimator is well outside bias bounds of moment estimator thus is more biased than moment estimator (blue trace outside cyan line). **(B)** Example plots of solutions to the quadratic program’s approximation (orange traces) to square root (black dashed trace) of the eigenvalues of  $\Sigma_{1,2}$  (black points). Re-scaling of singular values on vertical axis results in the deviation between the polynomial and the true square root evaluated at the true eigenvalues being exactly the bias of the associated estimates in panel A.

## D.2 Validation on neural data

Here we demonstrate that the estimator performs as expected when applied to noisy non-normal data where covariance of the  $\hat{W}_p$  and the denominator of the similarity score must be estimated from data. (Experiments on synthetic data verifying the estimators behave as expected can be found in App. 2.3.) We do so by applying our estimator to neural data: calcium recordings from mouse primary visual cortex in responses to a set of 2,800 natural images repeated twice [25]. We found that our estimator became highly variable when applied to this data in part because of its low SNR and low number of repeat (average SNR  $\approx 0.1$ ). We thus select neurons with the highest levels of SNR in each recording to perform our analyses on. To assess variability of the estimates we ran independent simulations from the same distribution by randomly sub sampling stimuli presentations within a recording into 3 disjoint sets.

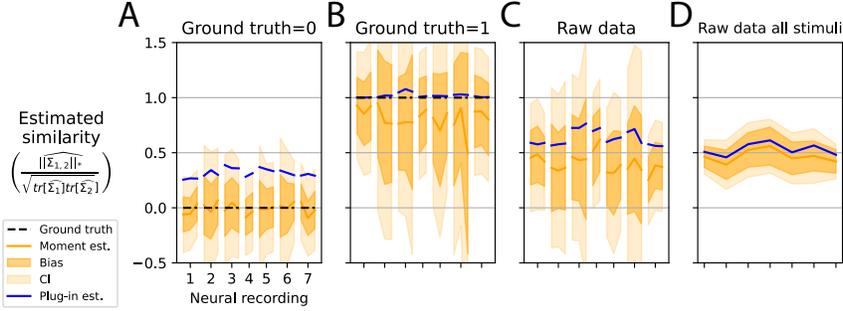


Figure 4: Validation of estimator on neural data [25]. **(A)** Comparison of estimators when ground truth similarity of neural data is set to 0. The estimator is applied to three disjoint sets of random stimuli for each recording ( $n = 7$ ). The estimated maximal bias is plotted in dark orange area and the confidence interval, which includes bias, is plotted in light orange. **(B)** Same simulation as (A) except ground truth similarity is 1. **(C)** Same as (B) except estimation of true similarity. **(D)** Estimation of true similarity on all stimuli. ( $M \approx 2,800$ ).

To determine the properties of the bias of our estimator requires comparison to the ground truth value of the similarity score. In the neural data ground truth is unknown. We thus developed two sampling schemes to set the ground truth similarity in the neural data. To set similarity to 0 we measured similarity between different populations of neurons shown different stimuli, thus the two populations responses are independent, thus their cross covariance is 0 so that the similarity score is 0. To set the similarity to 1 we measured similarity between the same population of neurons shown the same stimuli but on different trials, thus the only deviation in their responses is owing to trial-to-trial variability, thus their tuning similarity is 1.

We applied our estimator to populations of neurons ( $N = 40$  each) where the ground truth was zero. We found that across recordings the moment estimator correctly indicated the similarity was near 0 (Fig. 4A, orange trace overlaps black dashed) and the confidence intervals always contained the true similarity (light orange contains black dashed). On the other hand the plug-in estimator was upwardly biased (blue above black dashed). Thus the moment based estimator can accurately determine when the similarity is low in noisy neural data whereas the plug-in estimator cannot.

When ground truth similarity was 1, we found the bias of the moment estimator was worse than that of the plug-in (Fig. 4B, blue overlaps black dashed, orange below). This is consistent with our synthetic simulations (see Fig. 2A far right). The CIs always contained the true value but contained nearly the entire possible range of similarity values. Thus while the average estimate is high our confidence intervals are so wide that we do not have much information about the true similarity.

Finally, we assessed the estimators' performance measuring the true similarity between these populations of high SNR neurons (Fig. 4C). Across recordings the moment estimator was near 0.5 but confidence intervals were wide so there is little information about similarity even for the highest SNR neurons (light orange extends from 0 to 1 on vertical axis). The plug-in estimator reports a higher degree of similarity, that we heavily discount given its upward bias. When we included all stimuli ( $M \approx 2800$ ) we obtained more accurate estimates, learning that the true similarity is most likely between 0.25 and 0.75 (Fig. 4D). Thus small populations of well-tuned neurons in the same brain region have only intermediate levels of representational similarity. Overall, we find noisy data is a challenging setting for reducing the bias of shape similarity estimates.

### D.3 Experimental data from Stringer et al. [25]

Neural activity in mouse primary visual cortex was recorded using a two-photon microscope while mice were free to run on an air-floating ball. Recordings were collected across multiple depth planes at a frequency of 2.5 or 3 Hz, with planes 30-35  $\mu m$  apart. The field of view of the microscope was selected such that 10,000 neurons could be observed within a retinotopic location on the stimulus display.

All stimuli were presented for 0.5s with a random inter-stimulus interval between 0.3 and 1.1s consisting of a grey-screen. The images used in the experiment were taken from the ImageNet database, which includes categories such as birds, cats, and insects. The researchers manually

selected images that had a mix of low and high spatial frequencies and that did not consist of more than 50 % uniform background. All images were uniformly contrast-normalized by subtracting the local mean brightness and dividing by the local mean contrast. Each stimulus consisted of a different normalized image from the ImageNet database, with 2,800 different images used in total. The same image was displayed on all three screens, but each screen showed the image at a different rotation. Each of the 2,800 natural image stimuli were displayed twice in a recording in two blocks of the same randomized order.

Calcium movie data was processed using the Suite2p toolbox to estimate spike rates of neurons. Underlying neural activity was estimated using non-negative spike deconvolution (Frierich et. al., 2017). These deconvolved traces were normalized to the mean and standard deviation of their activity during a 30-minute period of grey-screen spontaneous activity. For further detail please see the original study [25]. All analyses done in this paper were performed on the pre-processed data available on figshare ([https://figshare.com/articles/Recordings\\_of\\_ten\\_thousand\\_neurons\\_in\\_visual\\_cortex\\_in\\_response\\_to\\_2\\_800\\_natural\\_images/6845348](https://figshare.com/articles/Recordings_of_ten_thousand_neurons_in_visual_cortex_in_response_to_2_800_natural_images/6845348)).