
Spectral Analysis of Representational Similarity with Limited Neurons

Hyunmo Kang^{*1} Abdulkadir Canatar^{*1,2} SueYeon Chung^{1,2}

Abstract

Measuring representational similarity between neural recordings and computational models is challenging due to constraints on the number of neurons that can be recorded simultaneously. In this work, we investigate how such limitations affect similarity measures, focusing on Canonical Correlation Analysis (CCA) and Centered Kernel Alignment (CKA). Leveraging tools from Random Matrix Theory, we develop a predictive spectral framework for these measures and demonstrate that finite neuron sampling systematically underestimates similarity due to eigenvector delocalization. To overcome this, we introduce a denoising method to infer population-level similarity, enabling accurate analysis even with small neuron samples. Our theory is validated on synthetic and real datasets, offering practical strategies for interpreting neural data under finite sampling constraints.

1. Introduction

Understanding how artificial neural networks relate to biological neural activity remains one of the central challenges in computational neuroscience (Carandini et al., 2005; van Gerven, 2017; Naselaris et al., 2011). As deep learning models become increasingly sophisticated at matching human-level performance on complex tasks, there is growing interest in whether these models actually learn representations that mirror those found in the brain (Yamins et al., 2014; Khaligh-Razavi & Kriegeskorte, 2014; Kell et al., 2018; Richards et al., 2019; Lindsay, 2021). However, a fundamental obstacle stands in the way of making this comparison: while artificial networks can be analyzed in their entirety, neuroscientists can only record from a small subset of neurons in any given brain region (Cai et al., 2016; Walther et al., 2016; Schütt et al., 2023). This sampling

limitation poses a critical challenge for the field. When we measure the similarity between model and neural representations using standard techniques like Canonical Correlation Analysis (CCA) or Centered Kernel Alignment (CKA), how much does our limited neural sample size distort the true relationship? The stakes for answering this question are high - these similarity metrics are increasingly used to evaluate competing neural network architectures and training approaches based on their match to brain data.

Our work provides the first rigorous theoretical framework for understanding how neuron sampling affects representational similarity measures. Our analysis reveals that measuring CCA and CKA with a limited number of recorded neurons systematically underestimates the true population-level similarity. This underestimation stems primarily from eigenvector delocalization (Aggarwal et al., 2023; Cizeau & Bouchaud, 1994; Baik et al., 2004)—a phenomenon where sample eigenvectors become increasingly misaligned with their population counterparts as the number of recorded neurons decreases. Understanding and accounting for this effect is crucial for accurate interpretation of neural representational similarities, particularly in experimental settings where only a subset of neurons can be recorded.

Our analysis proceeds in two parts. First, in the forward problem, we investigate how neuron sub-sampling from the full underlying population distorts the population eigencomponents and how this distortion affects the computed similarity measures. Second, in the backward problem, we ask whether observations from a finite number of neurons can be used to reliably infer the population representational similarity.

1.1. Our Contributions

- **Eigenvector-wise Analysis of Representation Similarity:** We show how neuron sub-sampling alters the eigenvalues and eigenvectors of the Gram matrix, leading to a systematic underestimation of CCA/CKA due to eigenvector delocalization.
- **Backward Inference via Denoising Eigenvectors:** We introduce a denoising method that leverages population eigenvalue priors (e.g., power law) to infer the true population similarity from limited data, substantially correcting the sampling bias.

^{*}Equal contribution ¹Center for Computational Neuroscience, Flatiron Institute, New York, NY, 10010 ²Center for Neural Science, New York University, New York, NY, 10003, USA. Correspondence to: SueYeon Chung <schung@nyu.edu>.

- **Validation on Real Neural Data:** Applying our framework to primate visual cortex recordings confirms that even modest neuron counts can lead to severe underestimation of model–brain similarity and that our method effectively recovers the missing signal.

1.2. Related Works

Representation similarity measures expressed in terms of eigencomponents were presented in detail by (Kornblith et al., 2019), who showed that CCA, CKA, and linear regression scores can all be written in terms of the eigenvalues and eigenvectors of the Gram matrices.

A key question is how these similarity measures behave under different kinds of noise. Broadly, there are two primary noise sources:

1. *Additive noise*, which arises from trial-to-trial variability and measurement error. In many studies, repeated trials and averaging can substantially mitigate this type of noise.
2. *Sampling noise*, which occurs because we can only record from a limited subset of neurons rather than the entire population. Consequently, the sample eigenvectors and eigenvalues differ from their population counterparts.

In this work, we focus on the latter issue—sampling noise—since we assume trial averaging already reduces the additive noise to a manageable level.

One approach to address sampling noise is by studying the *moments* of the Gram matrix (Kong & Valiant, 2017; Chun et al., 2024). While these methods provide a way to approximate the effect of sampling on the scalar values of certain similarity measures, they do not directly offer an interpretable description of what happens to the underlying eigencomponents. Recent work by (Pospisil & Pillow, 2024) provides bounds on representation similarity measures when the number of sampled neurons is limited. However, these bounds are tight only under the assumption of a white Wishart model (i.e., all population eigenvalues are 1). For more realistic data, where eigenvalues often decay according to a power-law, these bounds can become too loose to be practically informative.

Instead, we directly investigate how sampling noise affects both the eigenvalues and eigenvectors of the sample Gram matrix using random matrix theory (Potters & Bouchaud, 2020; Bun et al., 2018; 2017). Extensive results exist for white Wishart matrices and low-rank “spiked” models, including the Baik–Ben Arous–Péché (BBP) phase transition (Baik et al., 2004), which reveals that sample eigenvectors often serve as poor estimators of their population counter-

parts. These ideas have been extended to canonical correlation analysis (CCA) (Ma & Yang, 2022; Bykhovskaya & Gorin, 2025). However, the power-law-like spectra observed in neural data have not yet received comparable attention. Our work attempts to bridge this gap by studying sampling noise in representations with strongly decaying eigenvalues, which are ubiquitous in neural datasets.

2. Notation & Problem Setup

We use bold fonts for matrices and bracket notation for vectors. We use a tilde to denote quantities related to their population values.

We consider two centered population activations $\tilde{\mathbf{X}} \in \mathbb{R}^{P \times \tilde{N}_x}$ and $\tilde{\mathbf{Y}} \in \mathbb{R}^{P \times \tilde{N}_y}$ with \tilde{N}_x and \tilde{N}_y neurons, recorded in response to a fixed set of stimuli of size P . Centered means that we subtracted column-wise mean. Their corresponding population Gram matrices are given by $\tilde{\Sigma}_x = \tilde{\mathbf{X}}\tilde{\mathbf{X}}^\top$ and $\tilde{\Sigma}_y = \tilde{\mathbf{Y}}\tilde{\mathbf{Y}}^\top$ with eigendecomposition:

$$\tilde{\Sigma}_x = \sum_{i=1}^P \tilde{\lambda}_i |\tilde{u}_i\rangle\langle\tilde{u}_i|, \quad \tilde{\Sigma}_y = \sum_{a=1}^P \tilde{\mu}_a |\tilde{w}_a\rangle\langle\tilde{w}_a|. \quad (1)$$

The sample activations $\mathbf{X} \in \mathbb{R}^{P \times N_x}$ and $\mathbf{Y} \in \mathbb{R}^{P \times N_y}$ are assumed to be generated from the population ones by a random projection $\mathbf{X} = \tilde{\mathbf{X}}\mathbf{R}$ where $\mathbf{R} \in \mathbb{R}^{\tilde{N}_x \times N_x}$ is a random matrix with Gaussian i.i.d entries. Their Gram matrices are defined as $\Sigma_x = \mathbf{X}\mathbf{X}^\top$ and $\Sigma_y = \mathbf{Y}\mathbf{Y}^\top$ with eigendecomposition:

$$\Sigma_x = \sum_{i=1}^P \lambda_i |u_i\rangle\langle u_i|, \quad \Sigma_y = \sum_{a=1}^P \mu_a |w_a\rangle\langle w_a|. \quad (2)$$

Random projections serve as an effective approach for sampling high-dimensional data due to their geometry-preserving properties (Lahiri et al., 2016) and a popular method in analyzing neural dynamics from limited recordings (Gao et al., 2017). It also reduces our problem when we consider sample Gram matrices as structured random Wishart matrices (see SI.A.1).

We define the overlap as the squared inner product of two unit vectors. The overlap matrices are:

$$\begin{aligned} Q_{ij}^x &:= \mathbb{E}[\langle u_i | \tilde{u}_j \rangle^2] && \text{(brain_1 sample vs population)} \\ Q_{ab}^y &:= \mathbb{E}[\langle w_a | \tilde{w}_b \rangle^2] && \text{(brain_2 sample vs population)} \end{aligned} \quad (3)$$

$$\begin{aligned} M_{ia} &:= \mathbb{E}[\langle u_i | w_a \rangle^2] && \text{(brain_1 vs brain_2 sample)} \\ \tilde{M}_{ia} &:= \langle \tilde{u}_i | \tilde{w}_a \rangle^2 && \text{(brain_1 vs brain_2 population)} \end{aligned} \quad (4)$$

where the expectations are over different instances of neuron samplings. Here, the matrices \mathbf{Q} represent the self-overlap

between sample and population eigenvectors, \mathbf{M} represents the cross-overlap between two sample eigenvectors, and $\tilde{\mathbf{M}}$ between two population eigenvectors.

2.1. Common Representational Similarity Measures

Here, we review common representational similarity measures and show that these measures can be expressed in terms of the average quantities presented above.

Canonical Correlation Analysis (CCA) is an algorithm that sequentially finds a set of orthonormal vectors $\{\mathbf{v}_\alpha\}$ for which the correlation coefficients $\rho_\alpha = \text{corr}(\mathbf{X}\mathbf{v}_\alpha, \mathbf{Y}\mathbf{v}_\alpha)$ for two matrices \mathbf{X}, \mathbf{Y} are maximized (Hotelling, 1936). The squared sum of these coefficients gives the CCA similarity $\text{CCA} = \sum_\alpha \rho_\alpha^2$, and can be expressed in terms of the overlap matrix M_{ia} (Bjorck & Golub, 1973; Kornblith et al., 2019)

$$\text{CCA} = \sum_{a=1}^{N_y} \frac{\langle u_i | w_a \rangle^2}{\min(N_x, N_y)} = \sum_{i=1}^{N_x} \sum_{a=1}^{N_y} \frac{M_{ia}}{\min(N_x, N_y)}. \quad (5)$$

CCA has emerged as a popular tool in deep learning to compare neural representations (Raghu et al., 2017).

Canonical Correlation Analysis (CCA) is sensitive to perturbations when the condition number of \mathbf{X} or \mathbf{Y} is large (Golub & Zha, 1995). To enhance robustness, Singular Value CCA (SVCCA) performs CCA on truncated singular values of \mathbf{X} and \mathbf{Y} . In this approach, the sum of the overlap matrix \mathbf{M} is truncated to include only the first few components. To avoid confusion, from now on, we will refer to SVCCA truncated to the top ten components for both \mathbf{X} and \mathbf{Y} as CCA, i.e. (SV)CCA = $\frac{1}{10} \sum_{i=1}^{10} \sum_{a=1}^{10} M_{ia}$.

Centered Kernel Alignment (CKA) is a summary statistic of whether two representations agree on the (dis)similarity between a pair of examples based on their dot products (Cristianini et al., 2001). CKA is defined as $\frac{\text{Tr} \Sigma_x \Sigma_y}{\sqrt{\text{Tr} \Sigma_x^2 \text{Tr} \Sigma_y^2}}$ and essentially measures the angle between two Gram matrices. In terms of spectral components, it can be expressed as:

$$\text{CKA} = \sum_{i=1}^P \sum_{a=1}^P \frac{\lambda_i}{\sqrt{\sum_{j=1}^P \lambda_j^2}} \frac{\mu_a}{\sqrt{\sum_{b=1}^P \mu_b^2}} M_{ia}. \quad (6)$$

Note that CKA is very similar to CCA but with additional (normalized) eigenvalue terms. CKA will be the main focus of our work.

Representational Similarity Analysis (RSA) is a popular method in neuroscience used to compare different brain regions in response to the same set of stimuli (Kriegeskorte et al., 2008). It is similar to CKA, except RSA compares pair-wise Euclidean distances instead of dot products. Recent work has established its equivalence to CKA when

RSA is combined with an extra centering step (Williams, 2024). Therefore, our analyses are directly applicable to (centered-)RSA.

3. Theoretical Background

Treating Σ_x and Σ_y as random matrices described in Sec.2, we leverage results from random matrix theory (Potters & Bouchaud, 2020) to compute deterministic equivalents of average CCA and CKA in the asymptotic limit. Defining $q_x = P/N_x$ and $q_y = P/N_y$, we consider the limit $P, N_x, N_y \rightarrow \infty$ by keeping $q_x, q_y \sim \mathcal{O}(1)$.

Both similarity measures depend on the cross-overlap between sample eigenvectors M_{ia} defined in Eq. (4). Asymptotically M_{ia} decouples as (Bun et al., 2018)

$$M_{ia} = \sum_{j,b} Q_{ij}^x \tilde{M}_{jb} Q_{ba}^y, \quad (7)$$

where the self-overlaps Q_{ij}^x and Q_{ab}^y can be computed analytically (Ledoit & P  ch  , 2011). The self-overlap matrix for \mathbf{X} can be expressed in terms of the resolvent matrix $\mathbf{G}(z) = (z - \Sigma)^{-1}$ given by:

$$Q_{ij} = \text{const.} \lim_{\eta \rightarrow 0^+} \text{Im} \mathbf{G}_{jj}(\lambda_i - i\eta), \quad (8)$$

where the resolvent $\mathbf{G}(z)$ has a deterministic equivalent given by the following self-consistent equation

$$\mathbf{G}_{ij}(z) = \frac{\delta_{ij}}{z - \tilde{\lambda}_j(1 + q(z\mathbf{g}(z) - 1))}, \quad \mathbf{g}(z) = \frac{1}{P} \text{Tr} \mathbf{G}(z). \quad (9)$$

We provide a detailed derivation of these results in SI.A. Here, we note that the complex function $\mathbf{g}(z)$ and Eq. (8) can be solved numerically (see SI.D for details). Plugging in the formula for expected M_{ia} in Eq. (5) and Eq. (6), we get an analytical formula for CCA and CKA, respectively. Several remarks are in order:

- While our theory is applicable to the general cases where observations from both models are sampled, henceforth, we fix one of the models to be deterministic for practical reasons. Often, neural similarity measures are applied to compare biological data with limited neuron recordings to an artificial model where the entire population is available. For example fixing model \mathbf{Y} implies that its self-overlap \mathbf{Q}^y is just an identity matrix, hence simplifying Eq. (7) to $\mathbf{M} = \mathbf{Q}^x \tilde{\mathbf{M}}$.

- The analytical formula for CCA and CKA depends only on the population quantities. However, since the self-overlap matrix Q_{ij} in Eq. (8) explicitly depends on individual eigencomponents, its deterministic equivalent specifically depends on the population eigenvalue for the j^{th} component

($\tilde{\lambda}_j$), and the expected value of the sample eigenvalue for the i^{th} component ($\mathbb{E}[\lambda_i]$).

Sample Eigenvalues: Theoretical values of sample eigenvalues $\mathbb{E}[\lambda_i]$, in principle, can be computed by solving the following integral equation (Potters & Bouchaud, 2020)

$$\int_{\mathbb{E}[\lambda_i]}^{\infty} \rho(\lambda) d\lambda = \frac{i}{P}, \quad \rho(\lambda) = \frac{1}{\pi} \lim_{\eta \rightarrow 0^+} \text{Im} \mathfrak{g}(\lambda - i\eta). \quad (10)$$

Here, $\rho(\lambda)$ is the deterministic equivalent of the empirical eigenvalue density and depends only on the population eigenvalues (see SI.A.2). This may be problematic due to numerical instabilities.

However, we know that each single-trial eigenvalue concentrates around this mean with trial-to-trial fluctuations of $\mathcal{O}(1/\sqrt{P})$ (Potters & Bouchaud, 2020). In large P limit, we can neglect these fluctuations and replace $\mathbb{E}[\lambda_i]$ with a single-trial observation. We provide a detailed account of this approximation with supporting numerical experiments in SI.A.5.

Sample Eigenvectors: Unlike eigenvalues, the sample eigenvectors $\langle u_i | \tilde{u}_j \rangle^2$ exhibit trial-to-trial fluctuations that does not vanish even as $P \rightarrow \infty$ (see SI.A.6). Instead, we need to compute the mean value of the overlap represented by the squared overlap Q_{ij} in Eq. (4) (Bun et al., 2018).

BBP Phase Transition: Despite the inevitable fluctuation in the sample eigenvectors, their mean behavior can still differ markedly from that of the population eigenvectors. A classic example is the Baik–Ben Arous–Péché (BBP) phase transition (Baik et al., 2004). Consider a population Gram matrix with one large “spike” eigenvalue and the rest equal to 1. If the spike strength exceeds a critical threshold determined by P/N , then the sample eigenvector associated with that spike has an overlap on the order of $\mathcal{O}_P(1)$ with the population spike eigenvector. If the spike strength is below that threshold, however, the corresponding sample eigenvector becomes delocalized, and its overlap with the spike is of order $\mathcal{O}(1/P)$. We provide a detailed discussion of this in SI.A.6).

Numerical Confirmation: Finally, we numerically test the theoretical prediction for self-overlap given by Eq. (8) on the eigenvectors of deep neural network activations. We extract layer activations from a pre-trained ResNet18 on CIFAR-10 images and subsample N neurons through random projection. In Fig. 2a, we show the self-overlap Q_{ii} for the first few eigenvectors of the layer activations and demonstrate a perfect match with theory. As the number of neurons decreases, the number of delocalized eigenvectors increases since fewer eigenvectors have self-overlap $Q_{ii} \approx 1$.

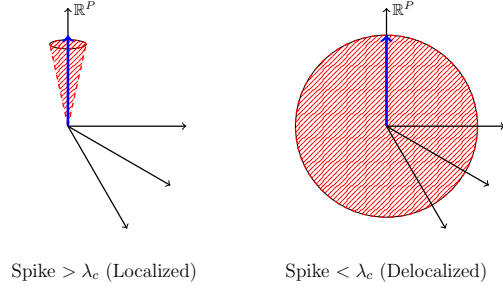


Figure 1: Illustration of BBP phase transition. On the left is the case where the spike is bigger than the critical value, and thus, the sample eigenvector related to the spike is close to the actual spike, lying at the cone with a small angle. On the right is the case where the spike is smaller than the critical value; in this case, sample eigenvector related to the spike mixes with bulk eigenvectors, ending up completely delocalized.

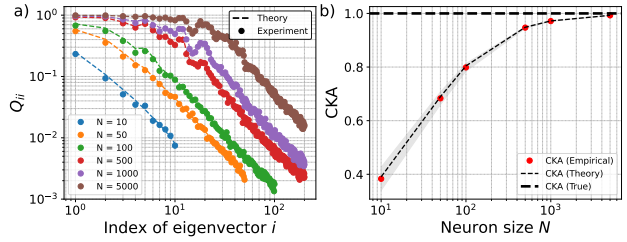


Figure 2: **a)** Self-overlap Q_{ii} between sample and population eigenvectors for ResNet18 activations. **b)** CKA between population and sample activations when N neurons are sampled. The gray-shaded region represents the standard deviation of empirical CKA across different random samplings.

The effect of eigenvector delocalization is reflected in the CKA between the sampled and population layer activations as shown in Fig. 2b. The alignment is completely misleading when small amounts of neurons are sampled and pose a significant problem for practical purposes.

4. Applying Theory to Representation Similarity

4.1. Forward Problem: Impact of Neuron Sampling on Similarity

In the forward problem, we assume that the population eigenvalues and eigenvectors are known. The first step is to obtain the typical sample eigenvalues by running a single-trial numerical simulation. We then move on to the eigenvectors by computing $\mathbf{Q}^{(x)}$ using Eq. (8). Finally, we calculate the overlap between the two systems, \mathbf{M} , using Eq. (7). Having these components allows us to evaluate

both CCA and CKA as functions of the number of neurons N .

As illustrated in Fig. 3, the theoretical predictions obtained from this eigen-decomposition match the observed CCA and CKA across different values of N . Notice that CKA decreases when the number of neurons N is reduced. As discussed above, both of these effects can be explained by the delocalization of eigenvectors.

4.2. Backward Problem: Inferring Population Similarity from Limited Neurons

Just like in our earlier analysis, inferring the population representational similarity begins with estimating the eigenvalues of the underlying population. In general, this is difficult because sample eigenvalues can deviate substantially from their population counterparts. Moreover, if $N < P$, there are $P - N$ zero eigenvalues in the sample covariance matrix, further complicating the problem.

However, if we adopt a parametric form, we can often achieve significant improvements in accuracy (Pospisil & Pillow, 2024). Here, we assume a power-law spectrum of the form $\tilde{\lambda}_i = i^{-1-\delta}$ (Stringer et al., 2019). We develop a numerical method based on random matrix theory that reliably infers the true decay rate of population eigenvalues based on only the sample eigenvalues (see SI.C for detailed analysis). More sophisticated approaches—such as allowing a broken power law and minimizing the error using unbiased moment estimators—are also possible (Pospisil & Pillow, 2024).

After estimating the population eigenvalues $\tilde{\lambda}_i$, we address the eigenvectors by computing $\mathbf{Q}^{(x)}$ using Eq. (8). Since every population eigenbasis produces the same mean self-overlap, having $\tilde{\lambda}_i$ is sufficient to find $\mathbf{Q}^{(x)}$.

Our final goal is to estimate the population cross-overlaps $\tilde{\mathbf{M}}$, which are required to infer the true population similarity between two systems. A straightforward way to do this is to invert the forward relationship $\mathbf{M} = \mathbf{Q}^{(x)} \tilde{\mathbf{M}}$. Two challenges arise in this naive approach. First, eigenvector statistics do not self-average (Potters & Bouchaud, 2020), so the empirical cross-overlap differs from \mathbf{M} , its expected value. This discrepancy can be partially mitigated by trial averaging or statistical bootstrapping.

Even if we obtain the mean cross-overlap, a second issue emerges: $\mathbf{Q}^{(x)}$ is not invertible unless $P \ll N$. As a result, it is impossible to recover the entire matrix $\tilde{\mathbf{M}}$. Intuitively, only the first few eigenvectors are well-localized; the rest delocalize and lose information, so we can only reliably retrieve the corresponding columns of $\tilde{\mathbf{M}}$.

To handle this practically, we use a constrained optimization:

$$\min_{\tilde{\mathbf{M}}} \|\mathbf{M} - \mathbf{Q}^{(x)} \tilde{\mathbf{M}}\|_F \quad (11)$$

where each element of $\tilde{\mathbf{M}}$ is restricted to lie in $[0, 1]$, as these entries represent squared inner products of eigenvectors. We provide the pseudo-code for this denoising algorithm in Alg.1

Algorithm 1 Inference of Population Cross-Overlap $\tilde{\mathbf{M}}$

Require: $\{\lambda_i\}_{i=1}^P$: Sample eigenvalues
 P : Number of stimuli
 N : Number of sampled neurons
 $\mathbf{M} \in \mathbb{R}^{P \times P}$: sample cross-overlap matrix

- 1: **Step 1: Estimate Population Eigenvalues**
 - 2: Assume power-law ansatz: $\lambda_i \propto i^{-1-\delta}$
 - 3: Find δ that best explains observed $\{\lambda_i\}_{i=1}^P$
 - 4: **Step 2: Compute Self-overlap matrix**
 - 5: $\mathbf{Q}^{(x)} \leftarrow \text{function}(\{\tilde{\lambda}_i\}, P, N)$
 - 6: **Step 3: Optimize Population Similarity**
 - 7: Solve constrained optimization problem:
 - 8: $\min_{\tilde{\mathbf{M}}} \|\mathbf{M} - \mathbf{Q}^{(x)} \cdot \tilde{\mathbf{M}}\|_F$
 - 9: subject to $\tilde{M}_{ij} \in [0, 1]$ for $\forall i, j$
 - 10: **return** $\tilde{\mathbf{M}}$
-

4.2.1. UP TO HOW MANY EIGENVECTORS CAN WE RESOLVE FOR GIVEN N, P ?

Consider a power-law spectrum, which decays relatively quickly. Under such a spectrum, only the leading sample eigenvectors tend to be well-localized, as shown in Fig. 4. If we run the backward algorithm, we observe that for a given N, P , we can reliably recover only those initial components that remain localized.

We can explicitly truncate these eigenvectors by taking a partial inverse of $\mathbf{Q}^{(x)}$ (see SI. E). However, this approach can be numerically unstable and might produce values of M_{ij} outside the $[0, 1]$ range.

Additionally, Fig. 5 demonstrates that, under a power-law of the same exponent, varying P has a subtler effect on these leading indices than varying N , which significantly affects localization.

4.2.2. WHY THIS IS SUFFICIENT FOR INFERRING POPULATION SIMILARITY

Although our denoising approach only manages to recover the leading few eigencomponents (those that remain localized), it is precisely these components that matter most for similarity measures like CKA and (SV)CCA. As shown in

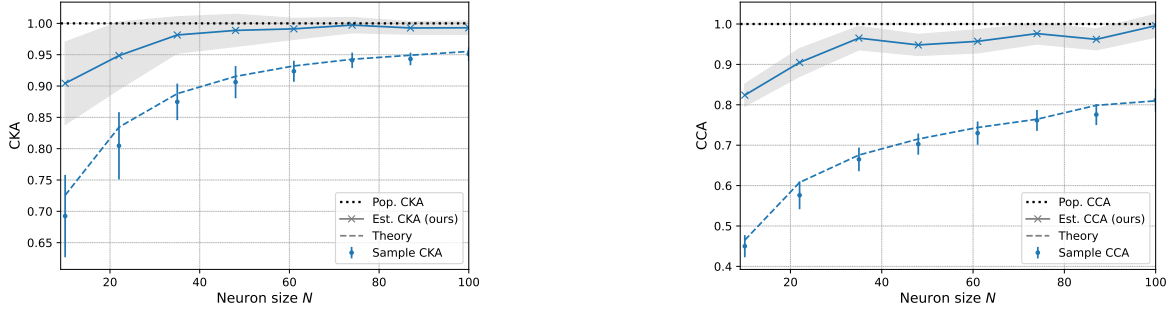


Figure 3: **Comparison of sample vs population measures for CKA and CCA:** Error bars represent empirical sample similarity and dotted lines the theoretical predictions. The black dotted line marks the true population similarity which is set to 1 for both measures. Solid lines indicate inferred true similarity from samples. Sample similarity is lower due to eigenvector delocalization, while our method consistently provides a closer estimate of the true value.

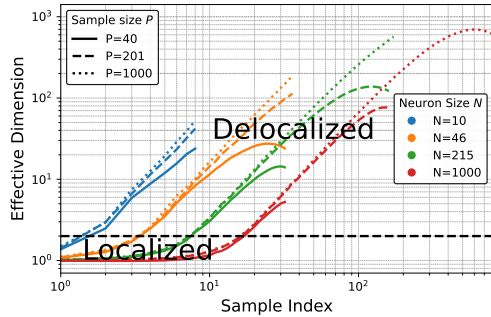


Figure 4: Participation ratio (P.R.) of self-overlap ($1/\sum_j Q_{ij}^2$), indicating the onset of eigenvector delocalization, for a power-law spectrum $\tilde{\lambda}_i \sim i^{-1.2}$. For fixed N , increasing P marginally affects the leading eigenvectors. By contrast, for fixed P , increasing N makes more eigenvectors localized. Only sample eigenvectors below the black horizontal line are localized (P.R. ≈ 1). Heuristically, \tilde{M}_{ia} can be recovered reliably for only indices below this line.

Fig. 3, these metrics are governed primarily by the initial eigenvalues and eigenvectors. Thus, even with a very limited number of neurons, estimating those leading components is sufficient for practical purposes.

To estimate population CKA from sample observations, one could focus solely on the denominator, as the sample numerator’s expected value matches the population value (Gretton et al., 2005). This allows applying methods from (Kong & Valiant, 2017; Chun et al., 2024), exploiting the fact that the changes in eigenvalues and eigenvectors offset each other. However, this approach estimates CKA using only eigenvalues, ignoring eigenvector statistics. In contrast, our method incorporates both inferred eigenvalues and eigenvectors for a more complete estimation.

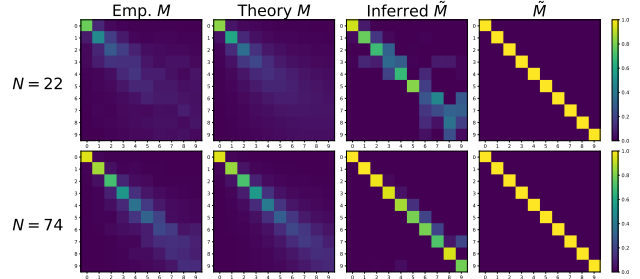


Figure 5: **Recovering population overlaps:** Each column shows the single-trial empirical M , the theoretical prediction of M , the inferred population overlap \tilde{M} , and the actual population overlap \tilde{M} . With fewer neurons N , sample eigenvectors become delocalized, causing large discrepancies. Nevertheless, our inference method successfully recovers the dominant overlaps, which are enough for global similarity measures such as CKA and CCA.

5. Experiments

5.1. Synthetic Data with a Known Population Gram Matrix

We first evaluate our approach on a synthetic dataset where the population Gram matrix is fully specified, allowing us to directly compare our estimated similarity measures against the ground-truth population values. For simplicity, we set the two population Gram matrices to be identical, i.e., $\tilde{\Sigma}_x = \tilde{\Sigma}_y$. Under this setup, the population CKA and CCA should both be 1.

Fig. 3 illustrates that our forward and backward procedures work well. In the forward approach, we show that the eigenvector-based analysis matches the empirical results closely. In the backward approach, even with an extremely limited number of neurons ($N \approx 20$), our method

infers a population similarity close to 1, despite the observed sample similarity being substantially lower.

Since the population eigenvectors (and hence the population cross-overlaps \tilde{M}) are known, we can also verify how well the inferred overlaps \hat{M} match the true overlaps \tilde{M} . Specifically, Fig. 5 displays the top-left 10×10 block of each matrix: the empirical \mathbf{M} , the theoretical \mathbf{M} (second column), the inferred population overlap $\hat{\mathbf{M}}$ (third column), and the actual population overlap $\tilde{\mathbf{M}}$ (fourth column). Because we set $\tilde{\Sigma}_x = \tilde{\Sigma}_y$, the actual population cross-overlap $\tilde{\mathbf{M}}$ should be the identity matrix. However, with fewer neurons, the sample eigenvectors become more delocalized, as evident in the first column. Our theoretical prediction of this phenomenon (second column) aligns closely with the empirical observation. Notably, even with severely limited neurons, our backward-inference method recovers a cross-overlap matrix $\hat{\mathbf{M}}$ (third column) much closer to the true identity than the naive observed \mathbf{M} .

5.1.1. SAMPLING NEURONS CAN CHANGE REPRESENTATION SIMILARITY RANKING

Next, we showcase a synthetic example in which *sampling* can lead to a reversal in the similarity rankings of models. Specifically, we construct two models:

- **Model 1** has significant overlap with the “Brain” representation on its first 3 population eigenvectors.
- **Model 2** has significant overlap with the Brain on the next 3 eigenvectors.

We set the total population (SV)CCA of Model 2 to be higher than that of Model 1. However, as neurons are sampled, eigenvectors corresponding to larger indices (smaller eigenvalues) tend to delocalize more. Hence, the empirical cross-overlap \mathbf{M} for Model 2 deteriorates faster, causing its (SV)CCA to drop more than that of Model 1. Eventually, Model 1 overtakes Model 2 in the sample-based (SV)CCA ranking, as illustrated in Fig. 7.

Fig. 8 presents the empirical and population cross-overlaps of the two models (each compared to the Brain). We set $P = 200$ and $N = 30$, and all population eigenvalues follow a power-law with exponent -1.2 . Model 2’s higher-dimensional overlaps delocalize more strongly, producing an apparent discrepancy that flips their observed ranking once neuron sampling is taken into account.

5.2. Brain Data

Finally, we apply our denoising framework to real neural recordings in the primate visual cortex, comparing them against various computational model predictions. (for experimental details see SI.D)

In Fig. 6, we illustrate a scatter plot of the representation similarity for different models compared to neural responses from V2 cortex (Freeman et al., 2013; Schrimpf et al., 2018), given an artificially limited neuron count of $N = 20$ out of 103 neurons. The x -axis corresponds to the observed *sample* CKA or CCA, while the y -axis is our *inferred population* measure. Observe that our inference method consistently produces higher population similarity estimates than the naive sample estimates. In particular, certain models that appear to have lower similarity (when judged by the raw, sample-based metric) can actually exhibit higher *true* similarity to the brain once sampling effects are taken into account.

6. Conclusion and Outlook

We have presented an eigenvector-based analysis of how sampling a finite number of neurons affects representational similarity measures, including CCA and CKA. By applying methods from Random Matrix Theory, we established that this limited sampling systematically underestimates similarity because of eigenvector delocalization in the sample Gram matrices. Our framework provides:

- **Forward Analysis:** A procedure to predict how population eigenvalues and eigenvectors will manifest under neuron sampling, thus explaining the observed drop in representation similarity.
- **Backward Inference:** A denoising algorithm capable of inferring the *population* representation similarity from limited data, overcoming the biases introduced by sampling noise.

We validated our approach on both synthetic and real datasets. In the synthetic experiments, where the population Gram matrices were fully known, we showed that our method reliably recovers the true population overlaps and similarity values, even in regimes with very few neurons. Importantly, we highlighted a striking effect of sampling: under certain configurations, the ranking of two models with respect to the brain can be inverted when only a limited set of neurons is recorded. In real datasets from primate visual cortex, our method consistently produced higher *population* similarity estimates than naive sample-based methods, underscoring that the observed decrease in similarity is largely a sampling artifact.

Future Directions. There are several promising avenues for extending our work. First, it would be valuable to explore more sophisticated spectral priors—such as broken power-law spectra—to account for multiple functional subpopulations in the data, each contributing a distinct spectral structure. Second, while we have focused on sampling noise,

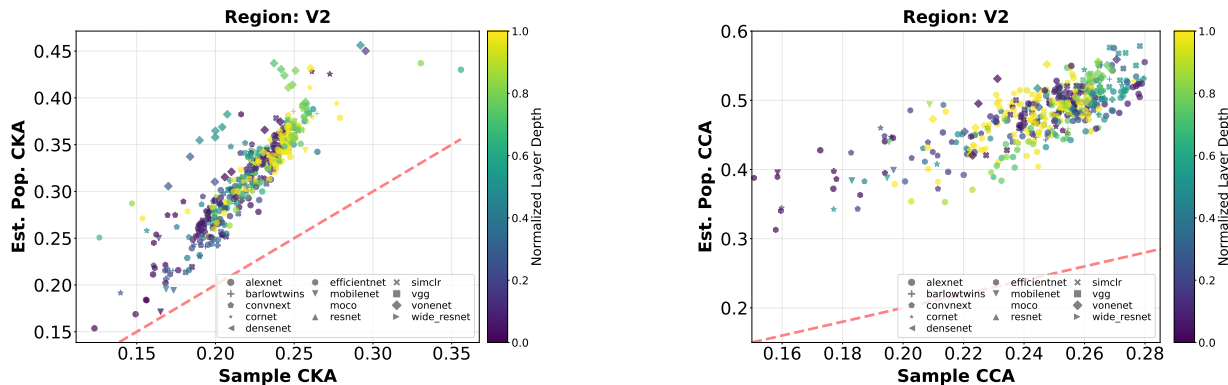


Figure 6: Scatter plots of observed sample similarity vs. inferred population similarity for multiple models compared to V2 cortex, using only $N = 20$ neurons (out of a larger set). **(Left)** CKA results; **(Right)** CCA results. The dotted line $y = x$ indicates equality. Notice that the inferred population similarity is consistently higher than the naive sample-based measure, demonstrating how limited neuron sampling can lead to underestimation of the true model-brain correspondence.

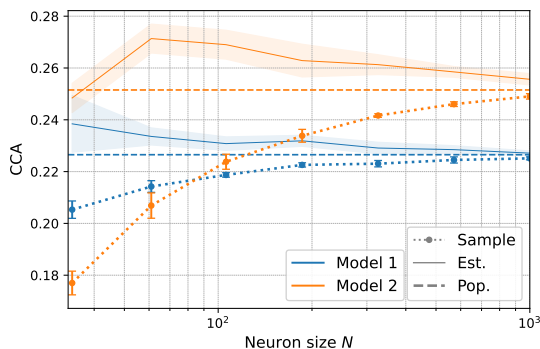


Figure 7: Sample-based CCA ranking flips despite Model 2 having a larger *population* CCA than Model 1. The decrease in Model 2’s CCA is more pronounced due to its stronger reliance on higher-indexed eigenvectors, which become more delocalized with limited neuron sampling.

future work should incorporate explicit models of additive noise that arises in real-time neurophysiological recordings, relaxing the assumption that trial averaging eliminates most of it. Third, improved denoising methods could be developed by adopting Bayesian approaches to model the joint distribution of sample eigenvectors and population eigenvectors (Monasson & Villamaina, 2015), thus allowing more accurate recovery of the population eigenspaces. Finally, as we outline in SI.B, our framework naturally extends to regression settings, where sampling-induced distortions in eigencomponents can adversely affect regression scores, much like their impact on representational similarity measures.

Overall, our results suggest that practical neuroscience studies must account for sampling-induced eigenvector delo-

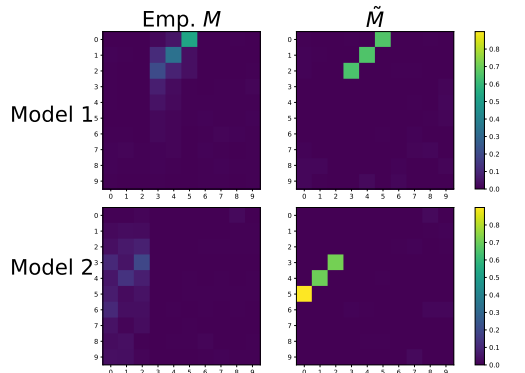


Figure 8: Empirical vs. population cross-overlaps for Model 1 vs. Brain and Model 2 vs. Brain. Here, $P = 200$ and $N = 30$. All three population eigenvalue spectra follow a power law with exponent -1.2 . Although Model 2’s true overlap is higher at the population level, it relies on higher-indexed (smaller eigenvalue) components, which delocalize more severely in the sample.

calization when interpreting representational similarity. By unveiling the intrinsic biases introduced by limited neuron sampling and proposing a systematic solution, we aim to provide neuroscientists and machine learning researchers with more reliable tools for comparing computational models and neural data.

7. Impact Statement

Our work may be highly impactful in providing reliable, robust similarity measures and check the reliability of existing studies based on neural recordings. Furthermore our method may help comparing similarities of artificial networks and may potentially be used in distilling large models.

References

- Aggarwal, A., Bordenave, C., and Lopatto, P. Mobility edge for lévy matrices, 2023. URL <https://arxiv.org/abs/2210.09458>.
- Atanasov, A., Zavatone-Veth, J. A., and Pehlevan, C. Scaling and renormalization in high-dimensional regression, 2024. URL <https://arxiv.org/abs/2405.00592>.
- Bahri, Y., Dyer, E., Kaplan, J., Lee, J., and Sharma, U. Explaining neural scaling laws. *Proceedings of the National Academy of Sciences*, 121(27):e2311878121, 2024.
- Baik, J., Arous, G. B., and Peche, S. Phase transition of the largest eigenvalue for non-null complex sample covariance matrices, 2004. URL <https://arxiv.org/abs/math/0403022>.
- Bjorck, A. and Golub, G. H. Numerical methods for computing angles between linear subspaces. *Mathematics of computation*, 27(123):579–594, 1973.
- Bordelon, B., Canatar, A., and Pehlevan, C. Spectrum dependent learning curves in kernel regression and wide neural networks. In *International Conference on Machine Learning*, pp. 1024–1034. PMLR, 2020.
- Bun, J., Bouchaud, J.-P., and Potters, M. Cleaning large correlation matrices: Tools from random matrix theory. *Physics Reports*, 666:1–109, January 2017. ISSN 0370-1573. doi: 10.1016/j.physrep.2016.10.005. URL <http://dx.doi.org/10.1016/j.physrep.2016.10.005>.
- Bun, J., Bouchaud, J.-P., and Potters, M. Overlaps between eigenvectors of correlated random matrices. *Physical Review E*, 98(5):052145, 2018.
- Bykhovskaya, A. and Gorin, V. High-dimensional canonical correlation analysis, 2025. URL <https://arxiv.org/abs/2306.16393>.
- Cai, M., Schuck, N. W., Pillow, J. W., and Niv, Y. A bayesian method for reducing bias in neural representational similarity analysis. In Lee, D., Sugiyama, M., Luxburg, U., Guyon, I., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016. URL <https://proceedings.neurips.cc/paper/2016/file/b06f50d1f89bd8b2a0fb771c1a69c2b0-Paper.pdf>.
- Canatar, A., Bordelon, B., and Pehlevan, C. Spectral bias and task-model alignment explain generalization in kernel regression and infinitely wide neural networks. *Nature communications*, 12(1):2914, 2021.
- Canatar, A., Feather, J., Wakhloo, A., and Chung, S. A spectral theory of neural prediction and alignment. *Advances in Neural Information Processing Systems*, 36, 2024.
- Carandini, M., Demb, J. B., Mante, V., Tolhurst, D. J., Dan, Y., Olshausen, B. A., Gallant, J. L., and Rust, N. C. Do we know what the early visual system does? *Journal of Neuroscience*, 25(46):10577–10597, 2005.
- Chun, C., Chung, S., and Lee, D. D. Estimating the spectral moments of the kernel integral operator from finite sample matrices, 2024. URL <https://arxiv.org/abs/2410.17998>.
- Cizeau, P. and Bouchaud, J. P. Theory of lévy matrices. *Phys. Rev. E*, 50:1810–1822, Sep 1994. doi: 10.1103/PhysRevE.50.1810. URL <https://link.aps.org/doi/10.1103/PhysRevE.50.1810>.
- Cristianini, N., Shawe-Taylor, J., Elisseeff, A., and Kandola, J. On kernel-target alignment. *Advances in neural information processing systems*, 14, 2001.
- Freeman, J., Ziemba, C. M., Heeger, D. J., Simoncelli, E. P., and Movshon, J. A. A functional and perceptual signature of the second visual area in primates. *Nature Neuroscience*, 16(7):974–981, Jul 2013. ISSN 1546-1726. doi: 10.1038/nn.3402. URL <https://doi.org/10.1038/nn.3402>.
- Gao, P., Trautmann, E., Yu, B., Santhanam, G., Ryu, S., Shenoy, K., and Ganguli, S. A theory of multineuronal dimensionality, dynamics and measurement. *BioRxiv*, pp. 214262, 2017.
- Golub, G. H. and Zha, H. *The canonical correlations of matrix pairs and their numerical computation*. Springer, 1995.
- Gretton, A., Bousquet, O., Smola, A., and Schölkopf, B. Measuring statistical dependence with hilbert-schmidt norms. In Jain, S., Simon, H. U., and Tomita, E. (eds.), *Algorithmic Learning Theory*, pp. 63–77, Berlin, Heidelberg, 2005. Springer Berlin Heidelberg. ISBN 978-3-540-31696-1.
- Hotelling, H. Relations between two sets of variates. *Biometrika*, 1936.
- Jacot, A., Simsek, B., Spadaro, F., Hongler, C., and Gabriel, F. Implicit regularization of random feature models. In *International Conference on Machine Learning*, pp. 4631–4640. PMLR, 2020.
- Kell, A. J., Yamins, D. L., Shook, E. N., Norman-Haignere, S. V., and McDermott, J. H. A task-optimized neural network replicates human auditory behavior, predicts brain responses, and reveals a cortical processing hierarchy. *Neuron*, 98(3):630–644, 2018.

- Khaligh-Razavi, S.-M. and Kriegeskorte, N. Deep supervised, but not unsupervised, models may explain it cortical representation. *PLoS computational biology*, 10(11): e1003915, 2014.
- Knowles, A. and Yin, J. Anisotropic local laws for random matrices. *Probability Theory and Related Fields*, 169: 257–352, 2017.
- Kong, W. and Valiant, G. Spectrum estimation from samples, 2017. URL <https://arxiv.org/abs/1602.00061>.
- Kornblith, S., Norouzi, M., Lee, H., and Hinton, G. Similarity of neural network representations revisited. In *International Conference on Machine Learning*, pp. 3519–3529. PMLR, 2019.
- Kriegeskorte, N., Mur, M., and Bandettini, P. A. Representational similarity analysis-connecting the branches of systems neuroscience. *Frontiers in systems neuroscience*, pp. 4, 2008.
- Lahiri, S., Gao, P., and Ganguli, S. Random projections of random manifolds. *arXiv preprint arXiv:1607.04331*, 2016.
- Ledoit, O. and Péché, S. Eigenvectors of some large sample covariance matrix ensembles. *Probability Theory and Related Fields*, 151(1):233–264, 2011.
- Ledoit, O. and Wolf, M. Numerical implementation of the quest function, 2016. URL <https://arxiv.org/abs/1601.05870>.
- Lindsay, G. W. Convolutional neural networks as a model of the visual system: Past, present, and future. *Journal of cognitive neuroscience*, 33(10):2017–2031, 2021.
- Ma, Z. and Yang, F. Sample canonical correlation coefficients of high-dimensional random vectors with finite rank correlations, 2022. URL <https://arxiv.org/abs/2102.03297>.
- Marchenko, V. A. and Pastur, L. A. Distribution of eigenvalues for some sets of random matrices. *Matematicheskii Sbornik*, 114(4):507–536, 1967.
- Monasson, R. and Villamaina, D. Estimating the principal components of correlation matrices from all their empirical eigenvectors. *EPL (Europhysics Letters)*, 112(5): 50001, December 2015. ISSN 1286-4854. doi: 10.1209/0295-5075/112/50001. URL <http://dx.doi.org/10.1209/0295-5075/112/50001>.
- Naselaris, T., Kay, K. N., Nishimoto, S., and Gallant, J. L. Encoding and decoding in fmri. *Neuroimage*, 56(2):400–410, 2011.
- Pospisil, D. A. and Pillow, J. W. Revisiting the high-dimensional geometry of population responses in visual cortex. *bioRxiv*, 2024. doi: 10.1101/2024.02.16.580726. URL <https://www.biorxiv.org/content/early/2024/02/21/2024.02.16.580726>.
- Potters, M. and Bouchaud, J.-P. *A First Course in Random Matrix Theory: for Physicists, Engineers and Data Scientists*. Cambridge University Press, 2020.
- Raghu, M., Gilmer, J., Yosinski, J., and Sohl-Dickstein, J. Svcca: Singular vector canonical correlation analysis for deep learning dynamics and interpretability, 2017. URL <https://arxiv.org/abs/1706.05806>.
- Richards, B. A., Lillicrap, T. P., Beaudoin, P., Bengio, Y., Bogacz, R., Christensen, A., Clopath, C., Costa, R. P., de Berker, A., Ganguli, S., et al. A deep learning framework for neuroscience. *Nature neuroscience*, 22(11): 1761–1770, 2019.
- Schrimpf, M., Kubilius, J., Hong, H., Majaj, N. J., Rajalingham, R., Issa, E. B., Kar, K., Bashivan, P., Prescott-Roy, J., Geiger, F., et al. Brain-score: Which artificial neural network for object recognition is most brain-like? *BioRxiv*, pp. 407007, 2018.
- Schütt, H. H., Kipnis, A. D., Diedrichsen, J., and Kriegeskorte, N. Statistical inference on representational geometries. *eLife*, 12:e82566, aug 2023. ISSN 2050-084X. doi: 10.7554/eLife.82566. URL <https://doi.org/10.7554/eLife.82566>.
- Silverstein, J. W. and Choi, S.-I. Analysis of the limiting spectral distribution of large dimensional random matrices. *Journal of Multivariate Analysis*, 54(2):295–309, 1995.
- Stringer, C., Pachitariu, M., Steinmetz, N., Carandini, M., and Harris, K. D. High-dimensional geometry of population responses in visual cortex. *Nature*, 571(7765): 361–365, 2019.
- van Gerven, M. A. A primer on encoding models in sensory neuroscience. *Journal of Mathematical Psychology*, 76: 172–183, 2017.
- Walther, A., Nili, H., Ejaz, N., Alink, A., Kriegeskorte, N., and Diedrichsen, J. Reliability of dissimilarity measures for multi-voxel pattern analysis. *NeuroImage*, 137:188–200, 2016.
- Williams, A. H. Equivalence between representational similarity analysis, centered kernel alignment, and canonical correlations analysis. In *UniReps: 2nd Edition of the Workshop on Unifying Representations in Neural Models*, 2024. URL <https://openreview.net/forum?id=zMdnFasgC>.

Yamins, D. L., Hong, H., Cadieu, C. F., Solomon, E. A., Seibert, D., and DiCarlo, J. J. Performance-optimized hierarchical models predict neural responses in higher visual cortex. *Proceedings of the national academy of sciences*, 111(23):8619–8624, 2014.

A. Detailed Derivation of the Main Result

A.1. The Sample Gram Matrix

Let $\tilde{\mathbf{X}} \in \mathbb{R}^{P \times \tilde{N}_x}$ denote the true population matrix with P samples and \tilde{N}_x neurons. We consider sampling only in the neuron/feature axis. The sample data $\mathbf{X} \in \mathbb{R}^{P \times N_x}$ is obtained by applying an $\tilde{N}_x \times N_x$ random projection matrix \mathbf{R}_x on $\tilde{\mathbf{X}}$

$$\mathbf{X} = \tilde{\mathbf{X}}\mathbf{R}_x, \quad (\mathbf{R}_x)_{ij} \sim \mathcal{N}\left(0, \frac{1}{N_x}\right). \quad (\text{S1})$$

The population and sample Gram matrices and their corresponding eigen-components are denoted as

$$\begin{aligned} \tilde{\Sigma}_x &= \tilde{\mathbf{X}}\tilde{\mathbf{X}}^\top = \sum_{i=1}^P \tilde{\lambda}_i |\tilde{u}_i\rangle\langle\tilde{u}_i|, \\ \Sigma_x &= \mathbf{X}\mathbf{X}^\top = \sum_{i=1}^P \lambda_i |u_i\rangle\langle u_i|. \end{aligned} \quad (\text{S2})$$

In Random Matrix Theory (RMT), it is often convenient to consider matrices of the form $\mathbf{M} = \sqrt{\mathbf{C}}\mathbf{W}\sqrt{\mathbf{C}}$, where $\mathbf{W} = \mathbf{R}\mathbf{R}^\top$ is a random Wishart matrix and \mathbf{C} is a deterministic square matrix. We first put Σ_x into this form to simplify our calculations (Knowles & Yin, 2017). The sample Gram matrix can be written in terms of the SVD components of $\tilde{\mathbf{X}} = \mathbf{U}\tilde{\Lambda}_x^{1/2}\mathbf{V}^\top$

$$\Sigma_x = \tilde{\mathbf{X}}\mathbf{R}_x\mathbf{R}_x^\top\tilde{\mathbf{X}}^\top = \mathbf{U}\tilde{\Lambda}_x^{1/2}\left(\mathbf{V}^\top\mathbf{R}_x\mathbf{R}_x^\top\mathbf{V}\right)\tilde{\Lambda}_x^{1/2}\mathbf{U}^\top, \quad (\text{S3})$$

where $\tilde{\Lambda}_x \in \mathbb{R}^{P \times \tilde{N}_x}$ is a diagonal matrix, and $\mathbf{U} \in \mathbb{R}^{P \times P}$ and $\mathbf{V} \in \mathbb{R}^{\tilde{N}_x \times \tilde{N}_x}$ orthogonal matrices. Since deterministic orthogonal transformations of Wishart matrices are again Wishart matrices, we get:

$$\Sigma_x = \mathbf{U}\tilde{\Lambda}_x^{1/2}\mathbf{W}_x\tilde{\Lambda}_x^{1/2}\mathbf{U}^\top, \quad (\text{S4})$$

where $\mathbf{W}_x = \mathbf{V}^\top\mathbf{R}_x\mathbf{R}_x^\top\mathbf{V}$ is a random Wishart matrix with aspect ratio $\phi_x = \tilde{N}_x/N_x$. We divide our discussion into two cases:

- When $P \geq \tilde{N}_x$, the eigenvalue matrix can be completed to a $P \times P$ -matrix by zero padding and replacing \mathbf{W}_x with a Wishart matrix with $q_x = P/N_x$. Using the orthogonality of \mathbf{U} , this allows us to express Σ_x as

$$\Sigma_x = (\mathbf{U}\tilde{\Lambda}_x^{1/2}\mathbf{U}^\top)(\mathbf{U}\mathbf{W}_x\mathbf{U}^\top)(\mathbf{U}\tilde{\Lambda}_x^{1/2}\mathbf{U}^\top) = \sqrt{\tilde{\Sigma}_x}\mathbf{W}_x\sqrt{\tilde{\Sigma}_x}, \quad (\text{S5})$$

where \mathbf{W}_x is a Wishart matrix with aspect ratio $q_x = P/N_x$.

- When $P < N_x$, the eigenvalue matrix and the Wishart matrix can be written as

$$\tilde{\Lambda}_x = \begin{pmatrix} \tilde{\Lambda}'_x & \mathbf{0} \end{pmatrix}, \quad \mathbf{W}_x = \begin{pmatrix} \mathbf{R}_1 \\ \mathbf{R}_2 \end{pmatrix} \begin{pmatrix} \mathbf{R}_1^\top & \mathbf{R}_2^\top \end{pmatrix}, \quad (\text{S6})$$

where the $P \times P$ matrix $\tilde{\Lambda}'_x$ is the non-zero part of $\tilde{\Lambda}_x$ and $\mathbf{R}_1 \in \mathbb{R}^{P \times \tilde{N}_x}$, $\mathbf{R}_2 \in \mathbb{R}^{(\tilde{N}_x - P) \times \tilde{N}_x}$ are two projection matrices. Plugging these back in, we arrive at the same form as the previous case.

In both cases, the statistics of Σ_x does not depend explicitly on \tilde{N}_x .

A.2. Eigenvalue statistics of sample Gram matrices

One of the main objectives of RMT is to understand the eigenvalue distribution of random matrices in terms of deterministic quantities (Potters & Bouchaud, 2020). Here, we review some classical results on the eigenvalue statistics of random

matrices of the form $\Sigma = \sqrt{\tilde{\Sigma}} \mathbf{W} \sqrt{\tilde{\Sigma}}$ where \mathbf{W} is a $P \times N$ Wishart matrix with ratio $q = \frac{P}{N}$. Here, Σ and $\tilde{\Sigma}$ are the sample and population Gram matrices, and they have the following eigendecompositions

$$\Sigma = \sum_{i=1}^P \lambda_i |u_i\rangle\langle u_i|, \quad \tilde{\Sigma} = \sum_{i=1}^P \tilde{\lambda}_i |\tilde{u}_i\rangle\langle \tilde{u}_i|. \quad (\text{S7})$$

We denote their (discrete-)eigenvalue distribution by $\rho(\lambda)$ and $\tilde{\rho}(\tilde{\lambda})$:

$$\rho(\lambda) = \frac{1}{P} \sum_{i=1}^P \delta(\lambda - \lambda_i), \quad \tilde{\rho}(\tilde{\lambda}) = \frac{1}{P} \sum_{i=1}^P \delta(\tilde{\lambda} - \tilde{\lambda}_i). \quad (\text{S8})$$

We define the resolvent of the random matrix \mathbf{X} and its trace as

$$\mathbf{G}(z) = (z - \Sigma)^{-1} = \sum_{i=1}^P \frac{|u_i\rangle\langle u_i|}{z - \lambda_i}. \quad (\text{S9})$$

The Stieltjes transform of the empirical spectral distribution is defined as

$$\mathbf{g}^P(z) := \int \frac{\rho(\lambda)}{z - \lambda} d\lambda = \frac{1}{P} \text{Tr} \mathbf{G}(z). \quad (\text{S10})$$

In large P limit, this quantity is self-averaging and there is a deterministic equivalent $\mathbf{g}(z) \sim \mathbf{g}^P(z)$ given by the self-consistent equation

$$\mathbf{g}(z) = \int \frac{\tilde{\rho}(\tilde{\lambda})}{z - \tilde{\lambda}(1 - q + qz\mathbf{g}(z))} d\tilde{\lambda}, \quad (\text{S11})$$

which only depends on the deterministic eigenvalues $\tilde{\rho}_x(\tilde{\lambda})$ and the ratio $q = P/N$ (Potters & Bouchaud, 2020). In practical applications, $\tilde{\rho}_x(\tilde{\lambda})$ is often replaced by with the uniform measure over the population eigenvalues $\{\tilde{\lambda}_i\}$ as defined in Eq. (S8). This remarkable result was first obtained in (Marchenko & Pastur, 1967) for white Wishart matrices (for which $\tilde{\Sigma} = \mathbf{I}$).

Due to the equivalence $\mathbf{g}(z) \sim \mathbf{g}^P(z)$ in large P limit, these two integrals are equivalent

$$\int \frac{\rho(\lambda)}{z - \lambda} d\lambda \xrightarrow{P \rightarrow \infty} \int \frac{\tilde{\rho}(\tilde{\lambda})}{z - \tilde{\lambda}(1 - q + qz\mathbf{g}(z))} d\tilde{\lambda}, \quad (\text{S12})$$

from which one can obtain the density of the limiting spectral density using the inversion formula (Bun et al., 2017)

$$\rho(\lambda) = \frac{1}{\pi} \lim_{\eta \rightarrow 0^+} \text{Im} \mathbf{g}(\lambda - i\eta). \quad (\text{S13})$$

The Stieltjes transform also connects to the effective regularization in ridge regression (Bordelon et al., 2020; Jacot et al., 2020; Canatar et al., 2021; Atanasov et al., 2024). We define a new function $\kappa(z)$ as

$$\kappa(z) := -\frac{z}{1 - q + qz\mathbf{g}(z)}, \quad \mathbf{g}(z) = z^{-1} - q^{-1}(z^{-1} + \kappa(z)^{-1}) \quad (\text{S14})$$

and express Eq. (S11) in terms of this quantity:

$$\mathbf{g}(z) = \frac{\kappa(z)}{z} \int \frac{\tilde{\rho}(\tilde{\lambda})}{\tilde{\lambda} + \kappa(z)} d\tilde{\lambda} = z^{-1} - q^{-1}(z^{-1} + \kappa(z)^{-1}). \quad (\text{S15})$$

Then, we obtain a new self-consistent equation for κ

$$\kappa(z) = -z + \kappa(z) \int \frac{q\tilde{\lambda}}{\tilde{\lambda} + \kappa(z)} \tilde{\rho}(\tilde{\lambda}) d\tilde{\lambda}, \quad (\text{S16})$$

which is also known as the Silverstein equation (Silverstein & Choi, 1995). Expressing this in terms of the discrete population eigenvalues, and evaluating it at $z = -\lambda$, we get

$$\kappa = \lambda + \kappa \frac{1}{N} \sum_{i=1}^P \frac{\tilde{\lambda}_i}{\tilde{\lambda}_i + \kappa}, \quad (\text{S17})$$

which is exactly the equation for the renormalized ridge parameter in (Canatar et al., 2021; Atanasov et al., 2024) with the scaling $\tilde{\lambda}_i \rightarrow N\tilde{\lambda}_i$.

A.3. Eigenvector statistics of sample Gram matrices and the self-overlap matrix

This result from Eq. (S11) can also be generalized to the resolvent matrix itself (Knowles & Yin, 2017; Bun et al., 2017), which becomes diagonal in the population eigenbasis:

$$\mathbf{G}(z) = \sum_{i=1}^P \frac{|u_i\rangle\langle u_i|}{z - \lambda_i} \sim \sum_{i=1}^P \frac{|\tilde{u}_i\rangle\langle \tilde{u}_i|}{z - \tilde{\lambda}_i(1 - q + qz\mathbf{g}(z))}, \quad (\text{S18})$$

where the integral over eigenvalues is replaced by the discrete measure over population eigenvalues. This allows us to study the eigenvector statistics by analyzing the quantity

$$\langle \tilde{u}_j | \mathbf{G}(z) | \tilde{u}_j \rangle = \sum_{i=1}^P \frac{\langle u_i | \tilde{u}_j \rangle^2}{z - \lambda_i} \sim \frac{1}{z - \tilde{\lambda}_j(1 - q + qz\mathbf{g}(z))}. \quad (\text{S19})$$

In large P limit, the sum over empirical eigenvalues become an integral:

$$\langle \tilde{u}_j | \mathbf{G}(z) | \tilde{u}_j \rangle \xrightarrow{P \rightarrow \infty} \int \frac{Q(\lambda, \tilde{\lambda}_j)}{z - \lambda} \rho(\lambda) d\lambda, \quad (\text{S20})$$

where we defined $Q(\lambda_i, \tilde{\lambda}_j) := P \langle u_i | \tilde{u}_j \rangle^2$ is the overlap between the i^{th} sample eigenvector and the j^{th} population eigenvector. Now, we can obtain $Q(\lambda_i, \tilde{\lambda}_j)$ using the following inversion formula

$$Q(\lambda_i, \tilde{\lambda}_j) = \frac{1}{\pi \rho(\lambda_i)} \lim_{\eta \rightarrow 0^+} \text{Im} \langle \tilde{u}_j | \mathbf{G}(\lambda_i - i\eta) | \tilde{u}_j \rangle. \quad (\text{S21})$$

Using the equivalence in Eq. (S19) and evaluating this expression explicitly:

$$Q(\lambda_i, \tilde{\lambda}_j) = \frac{q\lambda_i\tilde{\lambda}_j}{\left[\tilde{\lambda}_j(1 - q) - \lambda_i + q\lambda_i\tilde{\lambda}_j\mathfrak{h}(\lambda_i)\right]^2 + \left[q\lambda_i\tilde{\lambda}_j\pi\rho(\lambda_j)\right]^2}, \quad (\text{S22})$$

we get an explicit formula for eigenvector overlaps (Ledoit & Péché, 2011; Bun et al., 2017), where $\rho(\lambda_i)$ is given by Eq. (S13) and $\mathfrak{h}(z)$ is its Hilbert transform:

$$\mathfrak{h}(z) = \text{p.v.} \int \frac{\rho(\lambda)}{z - \lambda} d\lambda. \quad (\text{S23})$$

and can be obtained from the Stieltjes transform via

$$\lim_{\eta \rightarrow 0^+} \mathbf{g}(z - i\eta) = \mathfrak{h}(z) + i\pi\rho(z). \quad (\text{S24})$$

A.4. Overlap formula for two Gram matrices

Here, we provide a short review of the work by Bun et al. (2018) which derives an overlap formula between eigenvectors from random matrices. We consider observations from two representations $\mathbf{X} \in \mathbb{R}^{P \times N_x}$ and $\mathbf{Y} \in \mathbb{R}^{P \times N_y}$ in response to a common set of inputs of size P . Their sample Gram matrices have decompositions:

$$\Sigma_x = \mathbf{X}\mathbf{X}^\top = \sum_{i=1}^P \lambda_i |u_i\rangle\langle u_i|, \quad \Sigma_y = \mathbf{Y}\mathbf{Y}^\top = \sum_{a=1}^P \mu_a |w_a\rangle\langle w_a|. \quad (\text{S25})$$

We assume that \mathbf{X} and \mathbf{Y} are observations sampled from the underlying population features $\tilde{\mathbf{X}} \in \mathbb{R}^{P \times \tilde{N}_x}$ and $\tilde{\mathbf{Y}} \in \mathbb{R}^{P \times \tilde{N}_y}$ through independent random projections. The corresponding population Gram matrices are decomposed as:

$$\tilde{\Sigma}_x = \tilde{\mathbf{X}}\tilde{\mathbf{X}}^\top = \sum_{i=1}^P \tilde{\lambda}_i |\tilde{u}_i\rangle\langle\tilde{u}_i|, \quad \tilde{\Sigma}_y = \tilde{\mathbf{Y}}\tilde{\mathbf{Y}}^\top = \sum_{a=1}^P \tilde{\mu}_a |\tilde{w}_a\rangle\langle\tilde{w}_a|. \quad (\text{S26})$$

We consider two sample data matrices $\mathbf{X} \in \mathbb{R}^{P \times N_x}$ and $\mathbf{Y} \in \mathbb{R}^{P \times N_y}$. In Sec. A.1, we showed that the sample Gram matrices can be expressed in terms of the population ones as:

$$\begin{aligned} \Sigma_x &= \sqrt{\tilde{\Sigma}_x} \mathbf{W}_x \sqrt{\tilde{\Sigma}_x}, \\ \Sigma_y &= \sqrt{\tilde{\Sigma}_y} \mathbf{W}_y \sqrt{\tilde{\Sigma}_y}, \end{aligned} \quad (\text{S27})$$

where the Wishart matrices \mathbf{W}_x and \mathbf{W}_y have aspect ratios $q_x = P/N_x$ and $q_y = P/N_y$, respectively. Resolvents of the sample Gram matrices are

$$\mathbf{G}_x(z) \equiv (z - \Sigma_x)^{-1} = \sum_{i=1}^P \frac{|u_i\rangle\langle u_i|}{z - \lambda_i}, \quad \mathbf{G}_y(z') \equiv (z' - \Sigma_y)^{-1} = \sum_{a=1}^P \frac{|w_a\rangle\langle w_a|}{z' - \mu_a}. \quad (\text{S28})$$

We want to compute

$$\psi_P(z, z') = \mathbb{E} \left[\frac{1}{P} \text{Tr} [\mathbf{G}_x(z) \mathbf{G}_y(z')] \right] = \mathbb{E} \left[\frac{1}{P^2} \sum_{i,a=1}^P \frac{P \langle u_i | w_a \rangle^2}{(z - \lambda_i)(z' - \mu_a)} \right], \quad (\text{S29})$$

where the expectation is over random realizations of sample Gram matrices (Bun et al., 2018). In the limit $P \rightarrow \infty$, as empirical eigenvalues become continuous, this object approaches a deterministic function

$$\psi_P(z, z') \sim \psi(z, z') = \int \frac{\rho_x(\lambda) \rho_y(\mu)}{(z - \lambda)(z' - \mu)} M(\lambda, \mu) d\lambda d\mu, \quad M(\lambda_i, \mu_a) \sim \mathbb{E} \left[P \langle u_i | w_a \rangle^2 \right] \quad (\text{S30})$$

Here, $\rho_x(\lambda)$, $\rho_y(\mu)$ are the eigenvalue densities of Σ_x , Σ_y given by Eq. (S13). The function $M(\lambda_i, \mu_a) \sim \mathbb{E} \left[P \langle u_i | w_a \rangle^2 \right]$ denotes the expected overlap between two eigenvectors associated with eigenvalues λ_i and μ_a , and it is the central object for our analysis since it directly appears in CCA and CKA. This quantity can be obtained by computing $\psi(\lambda_i - i\eta, \mu_a + i\eta')$, collecting the term proportional to $\eta\eta'$ and taking the limit $\eta, \eta' \rightarrow 0$ (Bun et al., 2018):

$$\begin{aligned} \psi(\lambda_i - i\eta, \mu_a + i\eta') &= \int \frac{(\lambda_i - \lambda + i\eta) \rho_x(\lambda)}{(\lambda_i - \lambda)^2 + \eta^2} \frac{(\mu_a - \mu - i\eta') \rho_y(\mu)}{(\mu_a - \mu)^2 + \eta'^2} M(\lambda, \mu) d\lambda d\mu \\ &= \int \frac{\eta \rho_x(\lambda)}{(\lambda_i - \lambda)^2 + \eta^2} \frac{\eta' \rho_y(\mu)}{(\mu_a - \mu)^2 + \eta'^2} M(\lambda, \mu) d\lambda d\mu + (\dots) \\ &\underset{\eta, \eta' \rightarrow 0}{=} \pi^2 \rho_x(\lambda_i) \rho_y(\mu_a) M(\lambda_i, \mu_a) + (\dots) \end{aligned} \quad (\text{S31})$$

To simplify, we will assume that the population eigenvectors form a complete set of basis:

$$\mathbf{I} = \sum_{i=1}^P |\tilde{u}_i\rangle\langle\tilde{u}_i| = \sum_{a=1}^P |\tilde{w}_a\rangle\langle\tilde{w}_a|. \quad (\text{S32})$$

Then each resolvent in Eq. (S29) can be expressed in these bases:

$$\begin{aligned} \mathbf{G}_x(z) &= \sum_{i,j} |\tilde{u}_i\rangle\langle\tilde{u}_j| \Phi_{ij}^x(z), \quad \Phi_{ij}^x(z) := \langle\tilde{u}_i | \mathbf{G}_x(z) | \tilde{u}_j \rangle, \\ \mathbf{G}_y(z) &= \sum_{a,b} |\tilde{w}_a\rangle\langle\tilde{w}_b| \Phi_{ab}^y(z'), \quad \Phi_{ab}^y(z') := \langle\tilde{w}_a | \mathbf{G}_y(z) | \tilde{w}_b \rangle, \end{aligned} \quad (\text{S33})$$

where Φ_{ij}^x and Φ_{ab}^y are the matrix elements of resolvents $\mathbf{G}_x(z)$ and $\mathbf{G}_y(z')$ in their respective deterministic bases. Then, Eq. (S29) simplifies to

$$\psi_P(z, z') = \mathbb{E} \left[\frac{1}{P} \sum_{i,j,a,b} \Phi_{ij}^x(z) \tilde{C}_{ja} \Phi_{ab}^y(z') \tilde{C}_{bi}^\top \right] = \frac{1}{P} \sum_{i,j,a,b} \mathbb{E}[\Phi_{ij}^x(z)] \tilde{C}_{ja} \mathbb{E}[\Phi_{ab}^y(z')] \tilde{C}_{bi}^\top, \quad (\text{S34})$$

where we defined the deterministic overlap matrix elements $\tilde{C}_{ia} := \langle \tilde{u}_i | \tilde{w}_a \rangle$. In the second equality, we assumed that two resolvents are independent, reducing the problem to computing the expected resolvent of a single Gram matrix.

As discussed around Eq. (S19), the resolvent \mathbf{G}_x has a limiting value for $P \rightarrow \infty$ that is diagonal in the corresponding deterministic basis (Bun et al., 2017), and its matrix elements are given by:

$$\Phi_{ij}^x(z) = \frac{\delta_{ij}}{z - \tilde{\lambda}_i(1 - q_x + q_x z \mathfrak{g}_x(z))} + \mathcal{O}(P^{-1/2}), \quad (\text{S35})$$

where $\mathfrak{g}_x(z)$ satisfies the self-consistency condition in Eq. (S11).

In order to compute the overlap $M(\lambda_i, \mu_a)$, we use Eq. (S31) and collect the term proportional to $\eta\eta'$. Thanks to Eq. (S34) and Eq. (S35), this term simplifies to:

$$\pi^2 \rho_x(\lambda_i) \rho_y(\mu_a) M(\lambda_i, \mu_a) = \frac{1}{P} \sum_{j,b} \left(\lim_{\eta \rightarrow 0} \text{Im} \Phi_{jj}^x(\lambda_i - i\eta) \right) \tilde{C}_{jb}^2 \left(\lim_{\eta' \rightarrow 0} \text{Im} \Phi_{bb}^y(\mu_a - i\eta') \right). \quad (\text{S36})$$

Defining

$$Q_x(\lambda_i, \tilde{\lambda}_j) := \frac{1}{\pi \rho_x(\lambda_i)} \lim_{\eta \rightarrow 0} \text{Im} \Phi_{jj}^x(\lambda_i - i\eta), \quad Q_y(\mu_a, \tilde{\mu}_b) := \frac{1}{\pi \rho_y(\mu_a)} \lim_{\eta' \rightarrow 0} \text{Im} \Phi_{bb}^y(\mu_a - i\eta') \quad (\text{S37})$$

we get an equation for M as

$$M(\lambda_i, \mu_a) = \frac{1}{P} \sum_{j,b} Q_x(\lambda_i, \tilde{\lambda}_j) \tilde{C}_{jb}^2 Q_y(\mu_a, \tilde{\mu}_b). \quad (\text{S38})$$

Here, Q_x and Q_y were already calculated in Eq. (S22). Identifying the following quantities

$$\begin{aligned} Q_{ij}^x &\equiv \mathbb{E} \langle u_i | \tilde{u}_j \rangle^2 = \frac{1}{P} Q_x(\lambda_i, \tilde{\lambda}_j), & Q_{ab}^y &\equiv \mathbb{E} \langle w_a | \tilde{w}_b \rangle^2 = \frac{1}{P} Q_y(\mu_a, \tilde{\mu}_b), \\ M_{ia} &\equiv \mathbb{E} \langle u_i | w_a \rangle^2 = \frac{1}{P} M(\lambda_i, \mu_a), & \tilde{M}_{ia} &:= \langle \tilde{u}_i | \tilde{w}_a \rangle^2 = \tilde{C}_{ia}^2, \end{aligned} \quad (\text{S39})$$

we get our main result (Bun et al., 2018):

$$\begin{aligned} \mathbf{M} &= \mathbf{Q}^x \tilde{\mathbf{M}} \mathbf{Q}^{y^\top}, \\ Q_{ij}^x &= \frac{1}{P} \frac{q_x \lambda_i \tilde{\lambda}_j}{\left[\tilde{\lambda}_j(1 - q_x) - \lambda_i + q_x \lambda_i \tilde{\lambda}_j \mathfrak{h}_x(\lambda_i) \right]^2 + \left[q_x \lambda_i \tilde{\lambda}_j \pi \rho_x(\lambda_i) \right]^2}, \\ Q_{ab}^y &= \frac{1}{P} \frac{q_y \mu_a \tilde{\mu}_b}{\left[\tilde{\mu}_b(1 - q_y) - \mu_a + q_y \mu_a \tilde{\mu}_b \mathfrak{h}_y(\mu_a) \right]^2 + \left[q_y \mu_a \tilde{\mu}_b \pi \rho_y(\mu_a) \right]^2}. \end{aligned} \quad (\text{S40})$$

A.5. Statistics of sample eigenvalues and its concentration properties

As we discussed in the main text, the practical usage of Eq. (S22) requires computing the expectation value of individual sample eigenvalues. Eq.S40, treating i^{th} biggest sample eigenvalue as deterministic and plugging $\eta = 1/\sqrt{P}$

For sufficient conditions, we can show that the sample resolvent $\mathfrak{g}(z)$ self-averages. In this case, sample eigenvalue density $\rho(\lambda)$ converges in law. Here, we show that for practical use of Eq.S21, Eq.S40, we can treat i^{th} biggest eigenvalue as effectively deterministic in its most probable position.

Specifically, we demonstrate that for a large number of eigenvalues P , the most probable i -th largest eigenvalue λ_i satisfies

$$\int_{\lambda_i}^{\infty} \rho(\lambda) d\lambda = \frac{i}{P}, \quad (\text{S41})$$

and that the fluctuations around this most probable is $O(1/\sqrt{P})$.

Consider a set of P eigenvalues $\{\lambda_1, \lambda_2, \dots, \lambda_P\}$ drawn independently from the probability density $\rho(\lambda)$. We order these eigenvalues in descending order:

$$\lambda_{(1)} \geq \lambda_{(2)} \geq \dots \geq \lambda_{(P)},$$

where $\lambda_{(i)}$ denotes the i^{th} largest eigenvalue. To find the most probable value λ_i for the i^{th} largest eigenvalue, we focus on the probability that *exactly* i eigenvalues exceed a threshold $\bar{\lambda}$. If we define

$$F(\bar{\lambda}) = \int_{\bar{\lambda}}^{\infty} \rho(\lambda') d\lambda',$$

then the probability that exactly i out of P samples exceed $\bar{\lambda}$ is given by the binomial expression

$$F(\bar{\lambda}, P, i) = \binom{P}{i} [F(\bar{\lambda})]^i [1 - F(\bar{\lambda})]^{P-i}.$$

We determine the threshold $\bar{\lambda}_i$ that maximizes $F(\bar{\lambda}, P, i)$ by setting its derivative (with respect to $\bar{\lambda}$) to zero. From this calculation, one obtains the simple condition

$$F(\bar{\lambda}_i) = \frac{i}{P}.$$

Equivalently, since $F(\lambda) = \int_{\lambda}^{\infty} \rho(\lambda') d\lambda'$, the most probable i^{th} largest eigenvalue λ_i satisfies

$$\int_{\lambda_i}^{\infty} \rho(\lambda) d\lambda = \frac{i}{P}.$$

Now we calculate approximations for fluctuation around this most probable position. Let's analyze $F(\bar{\lambda}, P, i)$ near $\bar{\lambda}_i$. Write $\bar{\lambda} = \lambda_i + \delta\lambda$ and expand $F(\bar{\lambda})$ in a Taylor series about λ_i :

$$F(\bar{\lambda}) = F(\lambda_i + \delta\lambda) \approx F(\lambda_i) + \left. \frac{dF}{d\lambda} \right|_{\lambda_i} \delta\lambda + \frac{1}{2} \left. \frac{d^2F}{d\lambda^2} \right|_{\lambda_i} (\delta\lambda)^2 + \dots$$

Since $F(\lambda_i) = \frac{i}{P}$ and λ_i is determined by maximizing $F(\bar{\lambda}, P, i)$, the first derivative of F at λ_i vanishes:

$$\left. \frac{dF}{d\lambda} \right|_{\lambda_i} = 0,$$

thus

$$F(\bar{\lambda}) \approx \frac{i}{P} + \frac{1}{2} F''(\lambda_i) (\delta\lambda)^2.$$

(We expect $F''(\lambda_i) < 0$ since $F(\lambda)$ decreases with λ .)

Substituting this expansion back into $\binom{P}{i} [F(\bar{\lambda})]^i [1 - F(\bar{\lambda})]^{P-i}$, we find that the dominant dependence on $\delta\lambda$ appears in a Gaussian-like factor

$$\exp\left(-\frac{1}{2} |F''(\lambda_i)| P (\delta\lambda)^2\right).$$

This indicates that $\bar{\lambda}$ (the threshold that yields exactly i exceedances) is peaked sharply around λ_i with a variance

$$\sigma_i^2 = \frac{1}{-F''(\lambda_i) P}.$$

In summary, most probable i -th largest eigenvalue $\lambda_{(i)}$ is determined by

$$\int_{\lambda_i}^{\infty} \rho(\lambda) d\lambda = \frac{i}{P},$$

with fluctuation $O(1/\sqrt{P})$.

A.6. Statistics of sample eigenvalues and its concentration properties

Note that unlike eigenvalue density converges in law, eigenvector statistics Eq. (S21) is noisy even when $P \rightarrow \infty$ (Potters & Bouchaud, 2020). In this case, we define Q matrix as expectation over different trials as in Eq.S39. Equivalently, this could be obtained by averaging over small eigenvalue interval, which could be done by plugging small $\eta = 1/\sqrt{P}$ to extract pole. Note that this $1/\sqrt{P}$ is also obtained by analyzing fluctuation around most probable i -th biggest eigenvalue as in above. This is essentially averaging over cauchy distribution centered in λ with width η . Thus for practical usage of Eq.S40, we simply plug this most likely i -th eigenvalue (Bun et al., 2017), with $\eta = 1/\sqrt{P}$.

B. Relation to regression based similarity measures

Regression Score is not an representational similarity measure but commonly used for scoring model closeness to brain (Schrimpf et al., 2018; Canatar et al., 2024). Here, we discuss how our theoretical analysis for the overlap matrix \mathbf{M} also can be applied to the regression setting. Regression score measures how well a model's activations \mathbf{X} predict neural responses \mathbf{Y} via a linear probe. Concretely, one performs ridge regression on a training subset $(\mathbf{X}_{1:p}, \mathbf{Y}_{1:p})$ of size $p < P$, obtaining:

$$\hat{\mathbf{X}}(p) = \mathbf{Y} \hat{\beta}(p). \quad (\text{S42})$$

$$\hat{\beta}(p) = \arg \min_{\beta} \|\mathbf{Y}_{1:p}\beta - \mathbf{X}_{1:p}\|_F^2 + \alpha_{\text{reg}} \|\beta\|_F^2, \quad (\text{S43})$$

Then regression score gives the neural prediction error,

$$E_g(p) = \frac{\|\hat{\mathbf{X}}(p) - \mathbf{X}\|_F^2}{\|\mathbf{X}\|_F^2}, \quad (\text{S44})$$

Note that this error can be decomposed to each error mode, where $E_g(p) = \sum_i \widetilde{W}_i(p)$ where $\widetilde{W}_i(p) := \frac{\kappa^2}{1-\gamma} \frac{W_i}{(p\lambda_i + \kappa)^2}$.

The quantity W_i denotes the projection of target labels on the i^{th} -model eigenvalue end hence can be expressed in terms of eigencomponents, $W_i = \sum_j \frac{\lambda_j}{\sum_k \lambda_k} M_{ij}$. However, calculating W_i assumes that there is access to population level eigenvalues and poses a problem with limited data. In future work, we would like to test if our analyses help improve the reliability of regression based similarity methods.

C. Theory of Power Law Spectrum

Here, we consider the case where population spectrum obeys a power law:

$$\tilde{\lambda}_k = \left(\frac{k}{P}\right)^{-s}, \quad k = 1, \dots, P, \quad s > 1 \quad (\text{S45})$$

where we normalized eigenvalue indices explicitly by P . For large P , the population density becomes:

$$\tilde{\rho}(\tilde{\lambda}) = \frac{1}{P} \sum_{k=1}^P \delta(\tilde{\lambda} - \tilde{\lambda}_k) \sim \frac{1}{P} \int_1^P \delta(\tilde{\lambda} - \tilde{\lambda}_k) dk, \quad (\text{S46})$$

We change the variables to $\mu := \tilde{\lambda}_k$ for which we get:

$$d\mu = -sP^s k^{-s-1} dk = -\frac{s}{P} \mu^{1+1/s} dk. \quad (\text{S47})$$

In the limit $P \rightarrow \infty$, the density becomes

$$\tilde{\rho}(\tilde{\lambda}) = \frac{1}{s} \int_1^\infty \mu^{-1-1/s} \delta(\tilde{\lambda} - \mu) d\mu = \gamma \tilde{\lambda}^{-1-\gamma}, \quad \tilde{\lambda} \in [1, \infty], \quad \gamma = s^{-1}, \quad (\text{S48})$$

where we defined $\gamma \in [0, 1]$ for notational convenience. Note that, in this definition, the expectation value of $\tilde{\lambda}$ diverges.

Next, we solve the self-consistent equation for the Stieltjes transform Eq. (S11) which reads:

$$\mathbf{g}(z) = \int \frac{\tilde{\rho}(\tilde{\lambda})}{z - \tilde{\lambda}(1 - q + qz\mathbf{g}(z))} d\tilde{\lambda} = \frac{\kappa}{z} \int_1^\infty \frac{\gamma \tilde{\lambda}^{-1-\gamma}}{\kappa - \tilde{\lambda}} d\tilde{\lambda}, \quad \kappa := \frac{z}{z - \tilde{\lambda}(1 - q + qz\mathbf{g}(z))} \quad (\text{S49})$$

where we defined κ to simplify notation. This integral has an analytical solution expressed in terms of hypergeometric functions (Bahri et al., 2024):

$$\int_1^\infty \frac{\tilde{\lambda}^{-1-\gamma}}{\kappa - \tilde{\lambda}} d\tilde{\lambda} = -\frac{1}{1 + \gamma} {}_2F_1(1, 1 + \gamma, 2 + \gamma, \kappa). \quad (\text{S50})$$

In order to solve the self-consistent equation analytically, we need to expand the ${}_2F_1$:

$$\int_1^\infty \frac{\tilde{\lambda}^{-1-\gamma}}{\kappa - \tilde{\lambda}} d\tilde{\lambda} = -\pi \kappa^{-\gamma-1} (\cot(\pi\gamma) - i) - \kappa^{-1} \sum_{n=0}^\infty \frac{1}{n - \gamma} \kappa^{-n} \quad (\text{S51})$$

which can be truncated. Here, we keep all terms and arrange the self-consistent equation in the form

$$\frac{\mathbf{g}'(z)}{\gamma} = -\pi \kappa^{-\gamma} (\cot(\pi\gamma) - i) - \sum_{n=0}^\infty \frac{1}{n - \gamma} \kappa^{-n}, \quad \mathbf{g}'(z) := z\mathbf{g}(z) \quad (\text{S52})$$

where we defined $\mathbf{g}'(z)$ in terms of which κ becomes

$$\kappa = \frac{z}{z - \tilde{\lambda}(1 + q(\mathbf{g}'(z) - 1))}. \quad (\text{S53})$$

We expand the r.h.s. of Eq. (S52) in terms of $\mathbf{g}'(z)$:

$$\begin{aligned} \frac{\mathbf{g}'(z)}{\gamma} &= \frac{1}{\gamma} - \pi (\cot(\pi\gamma) - i) \left(\frac{z}{1 - q} \right)^{-\gamma} - \sum_{n=1}^\infty \frac{1}{(n - \gamma)} \left(\frac{z}{1 - q} \right)^{-n} \\ &\quad - \mathbf{g}'(z) \frac{q}{1 - q} \left(\pi (\cot(\pi\gamma) - i) \gamma \left(\frac{z}{1 - q} \right)^{-\gamma} + \sum_{n=1}^\infty \frac{n}{(n - \gamma)} \left(\frac{z}{1 - q} \right)^{-n} \right) + \mathcal{O}(\mathbf{g}'(z)^2). \end{aligned} \quad (\text{S54})$$

In order to obtain an analytical solution, we must assume $\mathbf{g}'(z) \gg 1$ and keep only the linear term above. Then, the solution to self-consistent equation becomes:

$$\mathbf{g}(z) = z^{-1} \frac{\frac{1}{\gamma} - \pi (\cot(\pi\gamma) - i) \left(\frac{z}{1 - q} \right)^{-\gamma} - \sum_{n=1}^\infty \frac{1}{(n - \gamma)} \left(\frac{z}{1 - q} \right)^{-n}}{\frac{1}{\gamma} + \frac{q}{1 - q} \left(\pi (\cot(\pi\gamma) - i) \gamma \left(\frac{z}{1 - q} \right)^{-\gamma} + \sum_{n=1}^\infty \frac{n}{(n - \gamma)} \left(\frac{z}{1 - q} \right)^{-n} \right)}. \quad (\text{S55})$$

Next, we compute the sample eigenvalue density $\rho(\lambda)$ and its Hilbert transform $\mathfrak{h}(\lambda)$ by computing

$$\lim_{\eta \rightarrow 0^+} \mathbf{g}(\lambda - i\eta) = \mathfrak{h}(\lambda) + i\pi\rho(\lambda). \quad (\text{S56})$$

This is an extremely tedious calculation which we perform using Mathematica¹. Furthermore, we expand the results in q and, assuming $q \ll 1$, keep only the linear term. In this regime, the leading order behavior of $\rho(\lambda)$ and $\mathfrak{h}(\lambda)$ looks like:

$$\begin{aligned} \rho(\lambda) &= \gamma \lambda^{-1-\gamma} \left(1 - q\gamma \left(2\pi\gamma \cot(\pi\gamma) \lambda^{-\gamma} + \sum_{n=1}^\infty \frac{n + \gamma}{n - \gamma} \lambda^{-n} \right) \right) + \mathcal{O}(q^2) \\ \mathfrak{h}(\lambda) &= \lambda^{-1} \left(1 - \lambda^{-\gamma} \pi\gamma \cot(\pi\gamma) - \lambda^{-1} \frac{\gamma}{1 - \gamma} \right) \\ &\quad + \pi\gamma^2 q \left(\pi\gamma \lambda^{-2\gamma-1} (\cot^2(\pi\gamma) - 1) + \lambda^{-\gamma-2} \frac{(\gamma + 1) \cot(\pi\gamma)}{1 - \gamma} \right) + \mathcal{O}(q^2, \lambda^3). \end{aligned} \quad (\text{S57})$$

¹We will provide the Mathematica file.

Here, we did not include higher order terms for $\mathfrak{h}(\lambda)$ to avoid clutter.

Finally, we use the formula for estimating sample eigenvalues Eq. (S41) for which we obtain an explicit formula:

$$\mathfrak{F}(\lambda, q; \gamma) := \int_{\lambda}^{\infty} \rho(\lambda) d\lambda = \lambda^{-\gamma} \left(1 - q\gamma^2 \left(\lambda^{-\gamma} \pi \cot \pi\gamma + \sum_{n=1}^{\infty} \frac{1}{n-\gamma} \lambda^{-n} \right) \right). \quad (\text{S58})$$

Here, semi-column separates sample related arguments that we have access empirically (λ_i, q) and the only population related quantity γ . Hence, using the following relation (Ledoit & Wolf, 2016; Bun et al., 2017)

$$\mathfrak{F}(\lambda_i, q; \gamma) = \frac{i}{P} \quad (\text{S59})$$

we can either predict the shape of empirical eigenvalues given the decay rate of population spectrum (forward), or infer the population decay rate given the empirical observations of eigenvalues (backward). Finally, we numerically test our theory and obtain perfect agreement with empirical data in Fig.S1.

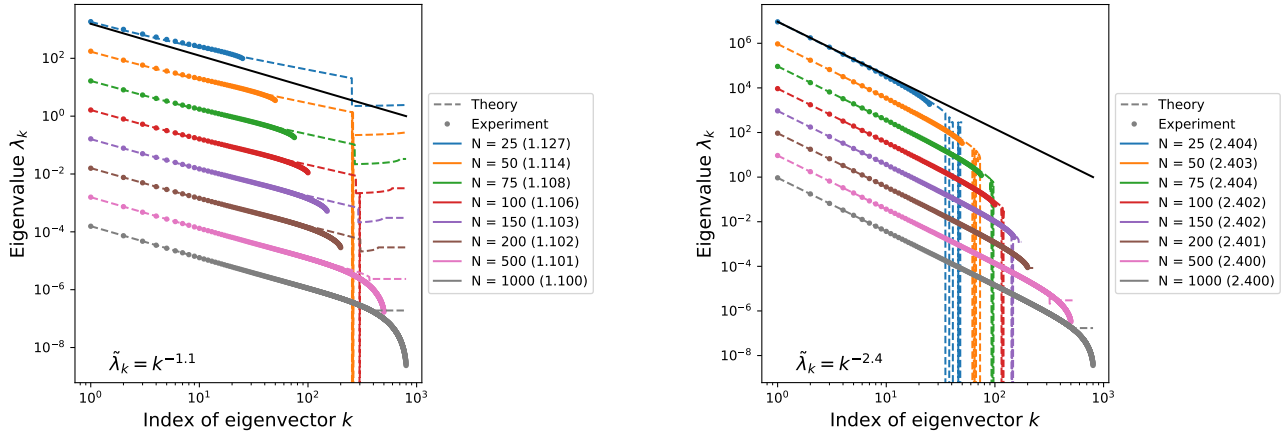


Figure S1: For a population spectrum with $\tilde{\lambda}_k = k^{-1.1}$ (Left) and $\tilde{\lambda}_k = k^{-2.4}$ (Right), we show the spectra of the empirical eigenvalues for different N . Black solid line indicates the true eigenvalue decay. The numbers in parentheses in the legend indicate the inferred true decay rate from a population of N . In the regime $s < 2$ ($\gamma > 0.5$), the empirical eigenvalues are always overestimated (Left), and in the regime $s \geq 2$ ($\gamma < 0.5$) they are always underestimated (Right).

D. Experimental Details

Code for all experiments are provided with supplemental material.

D.1. Synthetic Data

For the synthetic experiments, we generate a population activation matrix in $\mathbb{R}^{P \times \tilde{N}}$ whose Gram matrix follows a chosen spectral distribution (e.g., a power-law). We then form the sample activation matrix by projecting onto a random subset (or random linear subspace) of size N , yielding $\mathbb{R}^{P \times N}$. This procedure enables us to directly control the underlying population eigenvalues and eigenvectors, facilitating clean comparisons between sample-level and population-level similarity measures.

D.2. Brain Data

We employ a set of publicly available neural recordings from primate visual cortex (e.g., V2) and compare these against the representations of various vision models, similarly to the methodology in (Canatar et al., 2024). In total, we evaluate 32 models spanning supervised, self-supervised, and adversarially trained architectures, including well-known families such as ResNet, DenseNet, MobileNet, EfficientNet, and Vision Transformers. We extract intermediate-layer activations for each

model on the same set of visual stimuli used in the neural recordings, applying the standard preprocessing routines (e.g., image resizing, ImageNet normalization).

Within each model, we select one or more representative layers (e.g., post-ReLU or transformer blocks). We then compute Gram matrices from those activations, matching the dimensionality of the neural dataset. In scenarios where the dataset contains more neurons than we wish to analyze, we project the data into a lower-dimensional subspace of size N . Finally, we compute representational similarity (e.g., CKA or (SV)CCA) between these model-derived Gram matrices and the neural Gram matrices, both in their raw (sample) forms and using our denoising procedure for backward inference.

E. Another denoising method: truncated inverse

We utilize a truncated Singular Value Decomposition (SVD) to obtain a regularized estimate of $\tilde{\mathbf{M}}$:

$$\tilde{\mathbf{M}} = \mathbf{V}\Sigma_{\text{trunc}}^{-1}\mathbf{U}^{\top}\mathbf{M}, \tag{S60}$$

where $\mathbf{Q}^{(x)} = \mathbf{U}\Sigma\mathbf{V}^{\top}$ is the SVD of $\mathbf{Q}^{(x)}$, and $\Sigma_{\text{trunc}}^{-1}$ is the truncated inverse of the singular values, defined as:

$$\left(\Sigma_{\text{trunc}}^{-1}\right)_{ii} = \begin{cases} \frac{1}{\sigma_i} & \text{if } i \leq \tau, \\ 0 & \text{otherwise,} \end{cases} \tag{S61}$$