

Similarity of Neural Network Representations Revisited

Simon Kornblith¹ Mohammad Norouzi¹ Honglak Lee^{1,2} Geoffrey Hinton¹

Abstract

Recent work has sought to understand the behavior of neural networks by comparing representations between layers and between different trained models. We examine methods for comparing neural network representations based on canonical correlation analysis (CCA). We show that CCA belongs to a family of statistics for measuring multivariate similarity, but that neither CCA nor any other statistic that is invariant to invertible linear transformation can measure meaningful similarities between representations of higher dimension than the number of data points. We introduce a similarity index that measures the relationship between representational similarity matrices and does not suffer from this limitation. This similarity index is equivalent to centered kernel alignment (CKA) and is also closely connected to CCA. Unlike CCA, CKA can reliably identify correspondences between representations in networks trained from different initializations.

1. Introduction

Across a wide range of machine learning tasks, deep neural networks enable learning powerful feature representations automatically from data. Despite impressive empirical advances of deep neural networks in solving various tasks, the problem of understanding and characterizing the neural network representations learned from data remains relatively under-explored. Previous work (e.g. Advani & Saxe (2017); Amari et al. (2018); Saxe et al. (2014)) has made progress in understanding the theoretical dynamics of the neural network training process. These studies are insightful, but fundamentally limited, because they ignore the complex interaction between the training dynamics and structured data. A window into the network's representation can provide more information about the interaction between machine learning algorithms and data than the value of the loss function alone.

¹Google Brain ²University of Michigan. Correspondence to: Simon Kornblith <skornblith@google.com>.

This paper investigates the problem of measuring similarities between deep neural network representations. An effective method for measuring representational similarity could help answer many interesting questions, including: (1) Do deep neural networks with the same architecture trained from different random initializations learn similar representations? (2) Can we establish correspondences between layers of different network architectures? (3) How similar are the representations learned using the same network architecture from different datasets?

We build upon previous studies investigating similarity between the representations of neural networks (Laakso & Cottrell, 2000; Li et al., 2015; Raghu et al., 2017; Morcos et al., 2018; Wang et al., 2018). We are also inspired by the extensive neuroscience literature that uses representational similarity analysis (Kriegeskorte et al., 2008a; Edelman, 1998) to compare representations across brain areas (Haxby et al., 2001; Freiwald & Tsao, 2010), individuals (Connolly et al., 2012), species (Kriegeskorte et al., 2008b), and behaviors (Elsayed et al., 2016), as well as between brains and neural networks (Yamins et al., 2014; Khaligh-Razavi & Kriegeskorte, 2014; Sussillo et al., 2015).

Our key contributions are summarized as follows:

- We discuss the invariance properties of similarity indexes and their implications for measuring similarity of neural network representations.
- We motivate and introduce *centered kernel alignment* (CKA) as a similarity index and analyze the relationship between CKA, linear regression, canonical correlation analysis (CCA), and related methods (Raghu et al., 2017; Morcos et al., 2018).
- We show that CKA is able to determine the correspondence between the hidden layers of neural networks trained from different random initializations and with different widths, scenarios where previously proposed similarity indexes fail.
- We verify that wider networks learn more similar representations, and show that the similarity of early layers saturates at fewer channels than later layers. We demonstrate that early layers, but not later layers, learn similar representations on different datasets.

Problem Statement

Let $X \in \mathbb{R}^{n \times p_1}$ denote a matrix of activations of p_1 neurons for n examples, and $Y \in \mathbb{R}^{n \times p_2}$ denote a matrix of activations of p_2 neurons for the same n examples. We assume that these matrices have been preprocessed to center the columns. Without loss of generality we assume that $p_1 \leq p_2$. We are concerned with the design and analysis of a scalar *similarity index* $s(X, Y)$ that can be used to compare representations within and across neural networks, in order to help visualize and understand the effect of different factors of variation in deep learning.

2. What Should Similarity Be Invariant To?

This section discusses the invariance properties of similarity indexes and their implications for measuring similarity of neural network representations. We argue that both intuitive notions of similarity and the dynamics of neural network training call for a similarity index that is invariant to orthogonal transformation and isotropic scaling, but not invertible linear transformation.

2.1. Invariance to Invertible Linear Transformation

A similarity index is invariant to invertible linear transformation if $s(X, Y) = s(XA, YB)$ for any full rank A and B . If activations X are followed by a fully-connected layer $f(X) = \sigma(XW + \beta)$, then transforming the activations by a full rank matrix A as $X' = XA$ and transforming the weights by the inverse A^{-1} as $W' = A^{-1}W$ preserves the output of $f(X)$. This transformation does not appear to change how the network operates, so intuitively, one might prefer a similarity index that is invariant to invertible linear transformation, as argued by [Raghu et al. \(2017\)](#).

However, a limitation of invariance to invertible linear transformation is that any invariant similarity index gives the same result for any representation of width greater than or equal to the dataset size, *i.e.* $p_2 \geq n$. We provide a simple proof in Appendix A.

Theorem 1. *Let X and Y be $n \times p$ matrices. Suppose s is invariant to invertible linear transformation in the first argument, *i.e.* $s(X, Z) = s(XA, Z)$ for arbitrary Z and any A with $\text{rank}(A) = p$. If $\text{rank}(X) = \text{rank}(Y) = n$, then $s(X, Z) = s(Y, Z)$.*

There is thus a practical problem with invariance to invertible linear transformation: Some neural networks, especially convolutional networks, have more neurons in some layers than there are examples the training dataset ([Springenberg et al., 2015](#); [Lee et al., 2018](#); [Zagoruyko & Komodakis, 2016](#)). It is somewhat unnatural that a similarity index could require more examples than were used for training.

A deeper issue is that neural network *training* is not invari-

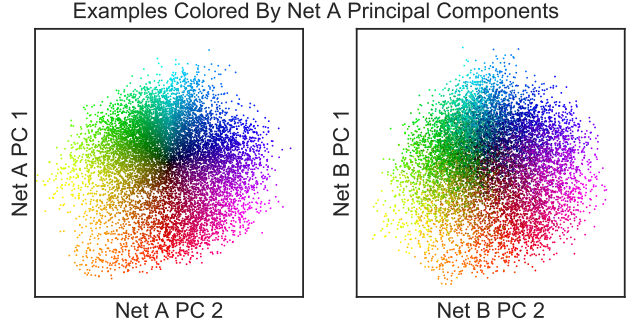


Figure 1. First principal components of representations of networks trained from different random initializations are similar. Each example from the CIFAR-10 test set is shown as a dot colored according to the value of the first two principal components of an intermediate layer of one network (left) and plotted on the first two principal components of the same layer of an architecturally identical network trained from a different initialization (right).

ant to arbitrary invertible linear transformation of inputs or activations. Even in the linear case, gradient descent converges first along the eigenvectors corresponding to the largest eigenvalues of the input covariance matrix ([LeCun et al., 1991](#)), and in cases of overparameterization or early stopping, the solution reached depends on the scale of the input. Similar results hold for gradient descent training of neural networks in the infinite width limit ([Jacot et al., 2018](#)). The sensitivity of neural networks training to linear transformation is further demonstrated by the popularity of batch normalization ([Ioffe & Szegedy, 2015](#)).

Invariance to invertible linear transformation implies that the scale of directions in activation space is irrelevant. Empirically, however, scale information is both consistent across networks and useful across tasks. Neural networks trained from different random initializations develop representations with similar large principal components, as shown in Figure 1. Consequently, Euclidean distances between examples, which depend primarily upon large principal components, are similar across networks. These distances are *meaningful*, as demonstrated by the success of perceptual loss and style transfer ([Gatys et al., 2016](#); [Johnson et al., 2016](#); [Dumoulin et al., 2017](#)). A similarity index that is invariant to invertible linear transformation ignores this aspect of the representation, and assigns the same score to networks that match only in large principal components or networks that match only in small principal components.

2.2. Invariance to Orthogonal Transformation

Rather than requiring invariance to any invertible linear transformation, one could require a weaker condition; invariance to orthogonal transformation, *i.e.* $s(X, Y) = s(XU, YV)$ for full-rank orthonormal matrices U and V such that $U^T U = I$ and $V^T V = I$.

Indexes invariant to orthogonal transformations do not share the limitations of indexes invariant to invertible linear transformation. When $p_2 > n$, indexes invariant to orthogonal transformation remain well-defined. Moreover, orthogonal transformations preserve scalar products and Euclidean distances between examples.

Invariance to orthogonal transformation seems desirable for neural networks trained by gradient descent. Invariance to orthogonal transformation implies invariance to permutation, which is needed to accommodate symmetries of neural networks (Chen et al., 1993; Orhan & Pitkow, 2018). In the linear case, orthogonal transformation of the input does not affect the dynamics of gradient descent training (LeCun et al., 1991), and for neural networks initialized with rotationally symmetric weight distributions, e.g. i.i.d. Gaussian weight initialization, training with fixed orthogonal transformations of activations yields the same distribution of training trajectories as untransformed activations, whereas an arbitrary linear transformation would not.

Given a similarity index $s(\cdot, \cdot)$ that is invariant to orthogonal transformation, one can construct a similarity index $s'(\cdot, \cdot)$ that is invariant to any invertible linear transformation by first orthonormalizing the columns of X and Y , and then applying $s(\cdot, \cdot)$. Given thin QR decompositions $X = Q_A R_A$ and $Y = Q_B R_B$ one can construct a similarity index $s'(X, Y) = s(Q_X, Q_Y)$, where $s'(\cdot, \cdot)$ is invariant to invertible linear transformation because orthonormal bases with the same span are related to each other by orthonormal transformation (see Appendix B).

2.3. Invariance to Isotropic Scaling

We expect similarity indexes to be invariant to isotropic scaling, i.e. $s(X, Y) = s(\alpha X, \beta Y)$ for any $\alpha, \beta \in \mathbb{R}^+$. That said, a similarity index that is invariant to both orthogonal transformation and non-isotropic scaling, i.e. rescaling of individual features, is invariant to any invertible linear transformation. This follows from the existence of the singular value decomposition of the transformation matrix. Generally, we are interested in similarity indexes that are invariant to isotropic but not necessarily non-isotropic scaling.

3. Comparing Similarity Structures

Our key insight is that instead of comparing multivariate features of an example in the two representations (e.g. via regression), one can first measure the similarity between every pair of *examples* in each representation separately, and then compare the similarity structures. In neuroscience, such matrices representing the similarities between examples are called representational similarity matrices (Kriegeskorte et al., 2008a). We show below that, if we use an inner product to measure similarity, the similarity between repre-

sentational similarity matrices reduces to another intuitive notion of pairwise feature similarity.

Dot Product-Based Similarity. A simple formula relates dot products between examples to dot products between features:

$$\langle \text{vec}(XX^T), \text{vec}(YY^T) \rangle = \text{tr}(XX^T YY^T) = \|Y^T X\|_F^2 \quad (1)$$

The elements of XX^T and YY^T are dot products between the representations of the i^{th} and j^{th} examples, and indicate the similarity between these examples according to the respective networks. The left-hand side of (1) thus measures the similarity between the inter-example similarity structures. The right-hand side yields the same result by measuring the similarity between *features* from X and Y , by summing the squared dot products between every pair.

Hilbert-Schmidt Independence Criterion. Equation 1 implies that, for centered X and Y :

$$\frac{1}{(n-1)^2} \text{tr}(XX^T YY^T) = \|\text{cov}(X^T, Y^T)\|_F^2 \quad (2)$$

The Hilbert-Schmidt Independence Criterion (Gretton et al., 2005) generalizes Equations 1 and 2 to inner products from reproducing kernel Hilbert spaces, where the squared Frobenius norm of the cross-covariance matrix becomes the squared Hilbert-Schmidt norm of the cross-covariance operator. Let $K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$ and let $L_{ij} = l(\mathbf{y}_i, \mathbf{y}_j)$ where k and l are two kernels. The empirical estimator of HSIC is:

$$\text{HSIC}(K, L) = \frac{1}{(n-1)^2} \text{tr}(KHLH) \quad (3)$$

where H is the centering matrix $H_n = I_n - \frac{1}{n} \mathbf{1}\mathbf{1}^T$. For linear kernels $k(\mathbf{x}, \mathbf{y}) = l(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \mathbf{y}$, HSIC yields (2).

Gretton et al. (2005) originally proposed HSIC as a test statistic for determining whether two sets of variables are independent. They prove that the empirical estimator converges to the population value at a rate of $1/\sqrt{n}$, and Song et al. (2007) provide an unbiased estimator. When k and l are universal kernels, $\text{HSIC} = 0$ implies independence, but HSIC is not an estimator of mutual information. HSIC is equivalent to maximum mean discrepancy between the joint distribution and the product of the marginal distributions, and HSIC with a specific kernel family is equivalent to distance covariance (Sejdinovic et al., 2013).

Centered Kernel Alignment. HSIC is not invariant to isotropic scaling, but it can be made invariant through normalization. This normalized index is known as centered kernel alignment (Cortes et al., 2012; Cristianini et al., 2002):

$$\text{CKA}(K, L) = \frac{\text{HSIC}(K, L)}{\sqrt{\text{HSIC}(K, K) \text{HSIC}(L, L)}} \quad (4)$$

Similarity Index	Formula	Invariant to		
		Invertible Linear Transform	Orthogonal Transform	Isotropic Scaling
Linear Reg. (R_{LR}^2)	$\ Q_Y^T X\ _F^2 / \ X\ _F^2$	Y only	✓	✓
CCA (R_{CCA}^2)	$\ Q_Y^T Q_X\ _F^2 / p_1$	✓	✓	✓
CCA ($\bar{\rho}_{CCA}$)	$\ Q_Y^T Q_X\ _* / p_1$	✓	✓	✓
SVCCA (R_{SVCCA}^2)	$\ (U_Y T_Y)^T U_X T_X\ _F^2 / \min(\ T_X\ _F^2, \ T_Y\ _F^2)$	If same subspace kept	✓	✓
SVCCA ($\bar{\rho}_{SVCCA}$)	$\ (U_Y T_Y)^T U_X T_X\ _* / \min(\ T_X\ _F^2, \ T_Y\ _F^2)$	If same subspace kept	✓	✓
PWCCA	$\sum_{i=1}^{p_1} \alpha_i \rho_i / \ \alpha\ _1, \alpha_i = \sum_j \langle \mathbf{h}_i, \mathbf{x}_j \rangle $	✗	✗	✓
Linear HSIC	$\ Y^T X\ _F^2 / (n-1)^2$	✗	✓	✗
Linear CKA	$\ Y^T X\ _F^2 / (\ X^T X\ _F \ Y^T Y\ _F)$	✗	✓	✓
RBF CKA	$\text{tr}(KHLH) / \sqrt{\text{tr}(KHKH)\text{tr}(LHLH)}$	✗	✓	✓*

Table 1. Summary of similarity methods investigated. Q_X and Q_Y are orthonormal bases for the columns of X and Y . U_X and U_Y are the left-singular vectors of X and Y sorted in descending order according to the corresponding singular vectors. $\|\cdot\|_*$ denotes the nuclear norm. T_X and T_Y are truncated identity matrices that select left-singular vectors such that the cumulative variance explained reaches some threshold. For RBF CKA, K and L are kernel matrices constructed by evaluating the RBF kernel between the examples as in Section 3, and H is the centering matrix $H_n = I_n - \frac{1}{n} \mathbf{1}\mathbf{1}^T$. See Appendix C for more detail about each technique.

*Invariance of RBF CKA to isotropic scaling depends on the procedure used to select the RBF kernel bandwidth parameter. In our experiments, we selected the bandwidth as a fraction of the median distance, which ensures that the similarity index is invariant to isotropic scaling.

For a linear kernel, CKA is equivalent to the RV coefficient (Robert & Escoufier, 1976) and to Tucker’s congruence coefficient (Tucker, 1951; Lorenzo-Seva & Ten Berge, 2006).

Kernel Selection. Below, we report results of CKA with a linear kernel and the RBF kernel $k(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\|\mathbf{x}_i - \mathbf{x}_j\|_2^2 / (2\sigma^2))$. For the RBF kernel, there are several possible strategies for selecting the bandwidth σ , which controls the extent to which similarity of small distances is emphasized over large distances. We set σ as a fraction of the median distance between examples. In practice, we find that RBF and linear kernels give similar results across most experiments, so we use linear CKA unless otherwise specified. Our framework extends to any valid kernel, including kernels equivalent to neural networks (Lee et al., 2018; Jacot et al., 2018; Garriga-Alonso et al., 2019; Novak et al., 2019).

4. Related Similarity Indexes

In this section, we briefly review linear regression, canonical correlation, and other related methods in the context of measuring similarity between neural network representations. We let Q_X and Q_Y represent any orthonormal bases for the columns of X and Y , i.e. $Q_X = X(X^T X)^{-1/2}$, $Q_Y = Y(Y^T Y)^{-1/2}$ or orthogonal transformations thereof. Table 1 summarizes the formulae and invariance properties of the indexes used in experiments. For a comprehensive general review of linear indexes for measuring multivariate similarity, see Ramsay et al. (1984).

Linear Regression. A simple way to relate neural network representations is via linear regression. One can fit

every feature in Y as a linear combination of features from X . A suitable summary statistic is the total fraction of variance explained by the fit:

$$R_{LR}^2 = 1 - \frac{\min_B \|Y - XB\|_F^2}{\|Y\|_F^2} = \frac{\|Q_Y^T X\|_F^2}{\|X\|_F^2} \quad (5)$$

We are unaware of any application of linear regression to measuring similarity of neural network representations, although Romero et al. (2015) used a least squares loss between activations of two networks to encourage thin and deep “student” networks to learn functions similar to wide and shallow “teacher” networks.

Canonical Correlation Analysis (CCA). Canonical correlation finds bases for two matrices such that, when the original matrices are projected onto these bases, the correlation is maximized. For $1 \leq i \leq p_1$, the i^{th} canonical correlation coefficient ρ_i is given by:

$$\begin{aligned} \rho_i = \max_{\mathbf{w}_X^i, \mathbf{w}_Y^i} \text{corr}(X\mathbf{w}_X^i, Y\mathbf{w}_Y^i) \\ \text{subject to } \forall_{j < i} X\mathbf{w}_X^i \perp X\mathbf{w}_X^j \\ \forall_{j < i} Y\mathbf{w}_Y^i \perp Y\mathbf{w}_Y^j \end{aligned} \quad (6)$$

The vectors $\mathbf{w}_X^i \in \mathbb{R}^{p_1}$ and $\mathbf{w}_Y^i \in \mathbb{R}^{p_2}$ that maximize ρ_i are the canonical weights, which transform the original data into canonical variables $X\mathbf{w}_X^i$ and $Y\mathbf{w}_Y^i$. The constraints in (6) enforce orthogonality of the canonical variables.

For the purpose of this work, we consider two summary

statistics of the goodness of fit of CCA:

$$R_{CCA}^2 = \frac{\sum_{i=1}^{p_1} \rho_i^2}{p_1} = \frac{\|Q_Y^T Q_X\|_F^2}{p_1} \quad (7)$$

$$\bar{\rho}_{CCA} = \frac{\sum_{i=1}^{p_1} \rho_i}{p_1} = \frac{\|Q_Y^T Q_X\|_*}{p_1} \quad (8)$$

where $\|\cdot\|_*$ denotes the nuclear norm. The mean squared CCA correlation R_{CCA}^2 is also known as *Yanai's GCD measure* (Ramsay et al., 1984), and several statistical packages report the sum of the squared canonical correlations $p_1 R_{CCA}^2 = \sum_{i=1}^{p_1} \rho_i^2$ under the name *Pillai's trace* (SAS Institute, 2015; StataCorp, 2015). The mean CCA correlation $\bar{\rho}_{CCA}$ was previously used to measure similarity between neural network representations in Raghu et al. (2017).

SVCCA. CCA is sensitive to perturbation when the condition number of X or Y is large (Golub & Zha, 1995). To improve robustness, *singular vector CCA* (SVCCA) performs CCA on truncated singular value decompositions of X and Y (Raghu et al., 2017; Mroueh et al., 2015; Kuss & Graepel, 2003). As formulated in Raghu et al. (2017), SVCCA keeps enough principal components of the input matrices to explain a fixed proportion of the variance, and drops remaining components. Thus, it is invariant to invertible linear transformation only if the retained subspace does not change.

Projection-Weighted CCA. Morcos et al. (2018) propose a different strategy to reduce the sensitivity of CCA to perturbation, which they term “projection-weighted canonical correlation” (PWCCA):

$$\rho_{PW} = \frac{\sum_{i=1}^c \alpha_i \rho_i}{\sum_{i=1}^c \alpha_i} \quad \alpha_i = \sum_j |\langle \mathbf{h}_i, \mathbf{x}_j \rangle| \quad (9)$$

where \mathbf{x}_j is the j^{th} column of X , and $\mathbf{h}_i = X \mathbf{w}_X^i$ is the vector of canonical variables formed by projecting X to the i^{th} canonical coordinate frame. As we show in Appendix C.3, PWCCA is closely related to linear regression, since:

$$R_{LR}^2 = \frac{\sum_{i=1}^c \alpha'_i \rho_i^2}{\sum_{i=1}^c \alpha'_i} \quad \alpha'_i = \sum_j \langle \mathbf{h}_i, \mathbf{x}_j \rangle^2 \quad (10)$$

Neuron Alignment Procedures. Other work has studied alignment between individual neurons, rather than alignment between subspaces. Li et al. (2015) examined correlation between the neurons in different neural networks, and attempt to find a bipartite match or semi-match that maximizes the sum of the correlations between the neurons, and then to measure the average correlations. Wang et al. (2018) proposed to search for subsets of neurons $\tilde{X} \subset X$ and $\tilde{Y} \subset Y$ such that, to within some tolerance, every neuron in \tilde{X} can be represented by a linear combination of neurons from \tilde{Y} and vice versa. They found that the maximum matching subsets are very small for intermediate layers.

Mutual Information. Among non-linear measures, one candidate is mutual information, which is invariant not only to invertible linear transformation, but to any invertible transformation. Li et al. (2015) previously used mutual information to measure neuronal alignment. In the context of comparing representations, we believe mutual information is not useful. Given any pair of representations produced by deterministic functions of the same input, mutual information between either and the input must be at least as large as mutual information between the representations. Moreover, in fully invertible neural networks (Dinh et al., 2017; Jacobsen et al., 2018), the mutual information between any two layers is equal to the entropy of the input.

5. Linear CKA versus CCA and Regression

Linear CKA is closely related to CCA and linear regression. If X and Y are centered, then Q_X and Q_Y are also centered, so:

$$R_{CCA}^2 = \text{CKA}(Q_X Q_X^T, Q_Y Q_Y^T) \sqrt{\frac{p_2}{p_1}} \quad (11)$$

When performing the linear regression fit of X with design matrix Y , $R_{LR}^2 = \|Q_Y^T X\|_F^2 / \|X\|_F^2$, so:

$$R_{LR}^2 = \text{CKA}(X X^T, Q_Y Q_Y^T) \frac{\sqrt{p_1} \|X^T X\|_F}{\|X\|_F^2} \quad (12)$$

When might we prefer linear CKA over CCA? One way to show the difference is to rewrite X and Y in terms of their singular value decompositions $X = U_X \Sigma_X V_X^T$, $Y = U_Y \Sigma_Y V_Y^T$. Let the i^{th} eigenvector of $X X^T$ (left-singular vector of X) be indexed as \mathbf{u}_X^i . Then R_{CCA}^2 is:

$$R_{CCA}^2 = \|U_Y^T U_X\|_F^2 / p_1 = \sum_{i=1}^{p_1} \sum_{j=1}^{p_2} \langle \mathbf{u}_X^i, \mathbf{u}_Y^j \rangle^2 / p_1 \quad (13)$$

Let the i^{th} eigenvalue of $X X^T$ (squared singular value of X) be indexed as λ_X^i . Linear CKA can be written as:

$$\begin{aligned} \text{CKA}(X X^T, Y Y^T) &= \frac{\|Y^T X\|_F^2}{\|X^T X\|_F \|Y^T Y\|_F} \\ &= \frac{\sum_{i=1}^{p_1} \sum_{j=1}^{p_2} \lambda_X^i \lambda_Y^j \langle \mathbf{u}_X^i, \mathbf{u}_Y^j \rangle^2}{\sqrt{\sum_{i=1}^{p_1} (\lambda_X^i)^2} \sqrt{\sum_{j=1}^{p_2} (\lambda_Y^j)^2}} \end{aligned} \quad (14)$$

Linear CKA thus resembles CCA weighted by the eigenvalues of the corresponding eigenvectors, *i.e.* the amount of variance in X or Y that each explains. SVCCA (Raghu et al., 2017) and projection-weighted CCA (Morcos et al., 2018) were also motivated by the idea that eigenvectors that correspond to small eigenvalues are less important, but

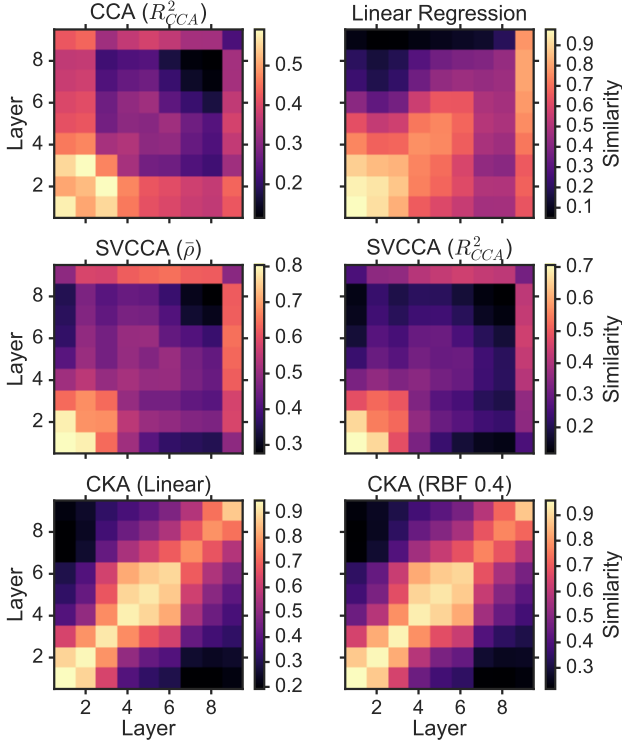


Figure 2. CKA reveals consistent relationships between layers of CNNs trained with different random initializations, whereas CCA, linear regression, and SVCCA do not. For linear regression, which is asymmetric, we plot R^2 for the fit of the layer on the x-axis with the layer on the y-axis. Results are averaged over 10 networks. See Table 2 for a numerical summary.

linear CKA incorporates this weighting symmetrically and can be computed without a matrix decomposition.

Comparison of (13) and (14) immediately suggests the possibility of alternative weightings of scalar products between eigenvectors. Indeed, as we show in Appendix D.1, the similarity index induced by “canonical ridge” regularized CCA (Vinod, 1976), when appropriately normalized, interpolates between R^2_{CCA} , linear regression, and linear CKA.

6. Results

6.1. A Sanity Check for Similarity Indexes

We propose a simple sanity check for similarity indexes: Given a pair of architecturally identical networks trained from different random initializations, for each layer in the first network, the most similar layer in the second network should be the architecturally corresponding layer. We train 10 networks and, for each layer of each network, we compute the accuracy with which we can find the corresponding layer in each of the other networks by maximum similarity. We then average the resulting accuracies. We compare CKA with CCA, SVCCA, PWCCA, and linear regression.

Index	Accuracy
CCA ($\bar{\rho}$)	1.4
CCA (R^2_{CCA})	10.6
SVCCA ($\bar{\rho}$)	9.9
SVCCA (R^2_{CCA})	15.1
PWCCA	11.1
Linear Reg.	45.4
Linear HSIC	22.2
CKA (Linear)	99.3
CKA (RBF 0.2)	80.6
CKA (RBF 0.4)	99.1
CKA (RBF 0.8)	99.3

Table 2. Accuracy of identifying corresponding layers based on maximum similarity for 10 architecturally identical 10-layer CNNs trained from different initializations, with logits layers excluded. For SVCCA, we used a truncation threshold of 0.99 as recommended in Raghu et al. (2017). For asymmetric indexes (PWCCA and linear regression) we symmetrized the similarity as $S + S^T$. CKA RBF kernel parameters reflect the fraction of the median Euclidean distance used as σ . Results not significantly different from the best result are bold-faced ($p < 0.05$, jackknife z-test).

We first investigate a simple VGG-like convolutional network based on All-CNN-C (Springenberg et al., 2015) (see Appendix E for architecture details). Figure 2 and Table 2 show that CKA passes our sanity check, but other methods perform substantially worse. For SVCCA, we experimented with a range of truncation thresholds, but no threshold revealed the layer structure (Appendix F.2); our results are consistent with those in Appendix E of Raghu et al. (2017).

We also investigate Transformer networks, where all layers are of equal width. In Appendix F.1, we show similarity between the 12 sublayers of the encoders of Transformer models (Vaswani et al., 2017) trained from different random initializations. All similarity indexes achieve non-trivial accuracy and thus pass the sanity check, although RBF CKA and R^2_{CCA} performed slightly better than other methods. However, we found that there are differences in feature scale between representations of feed-forward network and self-attention sublayers that CKA does not capture because it is invariant to non-isotropic scaling.

6.2. Using CKA to Understand Network Architectures

CKA can reveal pathology in neural networks representations. In Figure 3, we show CKA between layers of individual CNNs with different depths, where layers are repeated 2, 4, or 8 times. Doubling depth improved accuracy, but greater multipliers hurt accuracy. At 8x depth, CKA indicates that representations of more than half of the network are very similar to the last layer. We validated that these later layers do not refine the representation by training an ℓ^2 -regularized logistic regression classifier on each layer of the network. Classification accuracy in shallower architectures progressively improves with depth, but for the 8x deeper

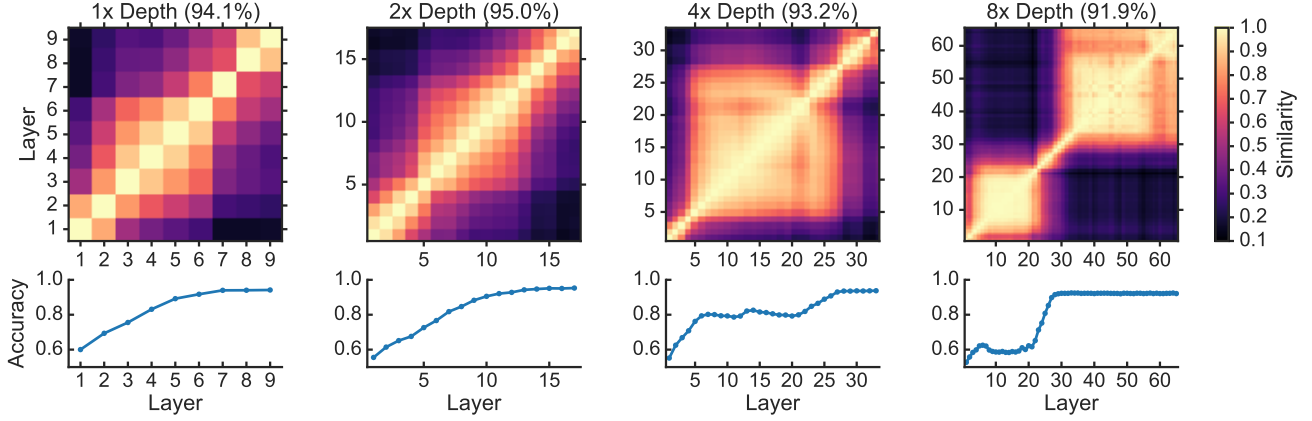


Figure 3. CKA reveals when depth becomes pathological. **Top**: Linear CKA between layers of individual networks of different depths on CIFAR-10. Titles show accuracy of each network. Later layers of the 8x depth network are similar to the last layer. **Bottom**: Accuracy of a logistic regression classifier trained on layers of the same networks is consistent with CKA.

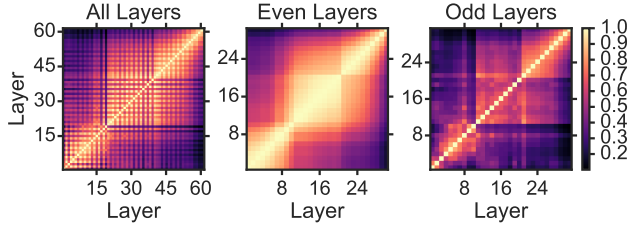


Figure 4. Linear CKA between layers of a ResNet-62 model. The grid pattern for ResNets in the left panel arises from the architecture. Right panels show similarity separately for even layer (post-residual) and odd layer (block interior) activations. Layers in the same block group (*i.e.* at the same feature map scale) are more similar than layers in different block groups.

network, accuracy plateaus less than halfway through the network. When applied to ResNets (He et al., 2016), CKA reveals no pathology (Figure 4). We instead observe a grid pattern that originates from the architecture: Post-residual activations are similar to other post-residual activations, but activations within blocks are not.

CKA is equally effective at revealing relationships between layers of different architectures. Figure 5 shows the relationship between different layers of networks with and without residual connections. CKA indicates that, as networks are made deeper, the new layers are effectively inserted in between the old layers. Other similarity indexes fail to reveal meaningful relationships between different architectures, as we show in Appendix F.5.

In Figure 6, we show CKA between networks with different layer widths. Like Morcos et al. (2018), we find that increasing layer width leads to more similar representations between networks. As width increases, CKA approaches 1; CKA of earlier layers saturates faster than later layers. Networks are generally more similar to other networks of the same width than they are to the widest network we trained.

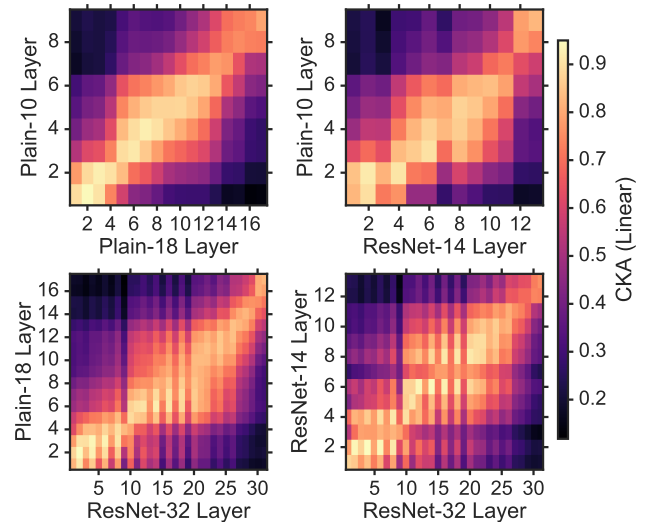


Figure 5. Linear CKA between layers of networks with different architectures.

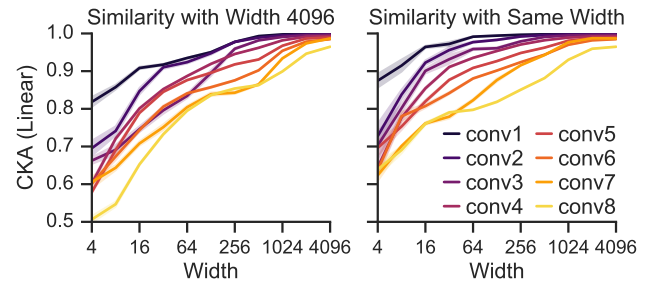


Figure 6. Layers become more similar to each other and to wide networks as width increases, but similarity of earlier layers saturates first. **Left**: Similarity of networks with the widest network we trained. **Middle**: Similarity of networks with other networks of the same width trained from random initialization. All CKA values are computed between 10 networks; shaded regions reflect jackknife standard error.

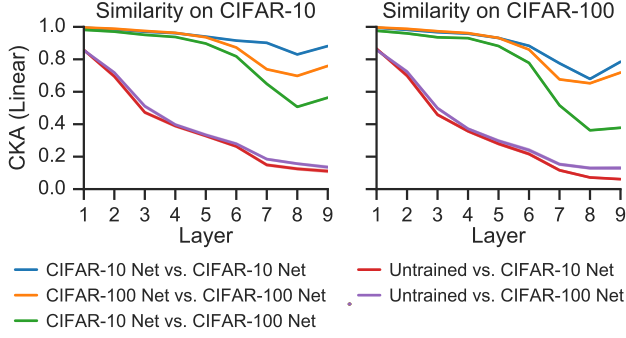


Figure 7. CKA shows that models trained on different datasets (CIFAR-10 and CIFAR-100) develop similar representations, and these representations differ from untrained models. The left panel shows similarity between the same layer of different models on the CIFAR-10 test set, while the right panel shows similarity computed on CIFAR-100 test set. CKA is averaged over 10 models of each type (45 pairs).

6.3. Similar Representations Across Datasets

CKA can also be used to compare networks trained on different datasets. In Figure 7, we show that models trained on CIFAR-10 and CIFAR-100 develop similar representations in their early layers. These representations require training; similarity with untrained networks is much lower. We further explore similarity between layers of untrained networks in Appendix F.3.

6.4. Analysis of the Shared Subspace

Equation 14 suggests a way to further elucidating what CKA is measuring, based on the action of one representational similarity matrix (RSM) YY^T applied to the eigenvectors \mathbf{u}_X^i of the other RSM XX^T . By definition, $XX^T\mathbf{u}_X^i$ points in the same direction as \mathbf{u}_X^i , and its norm $\|XX^T\mathbf{u}_X^i\|_2$ is the corresponding eigenvalue. The degree of scaling and rotation by YY^T thus indicates how similar the action of YY^T is to XX^T , for each eigenvector of XX^T . For visualization purposes, this approach is somewhat less useful than the CKA summary statistic, since it does not collapse the similarity to a single number, but it provides a more complete picture of what CKA measures. Figure 8 shows that, for large eigenvectors, XX^T and YY^T have similar actions, but the rank of the subspace where this holds is substantially lower than the dimensionality of the activations. In the penultimate (global average pooling) layer, the dimensionality of the shared subspace is approximately 10, which is the number of classes in the CIFAR-10 dataset.

7. Conclusion and Future Work

Measuring similarity between the representations learned by neural networks is an ill-defined problem, since it is not entirely clear what aspects of the representation a similarity

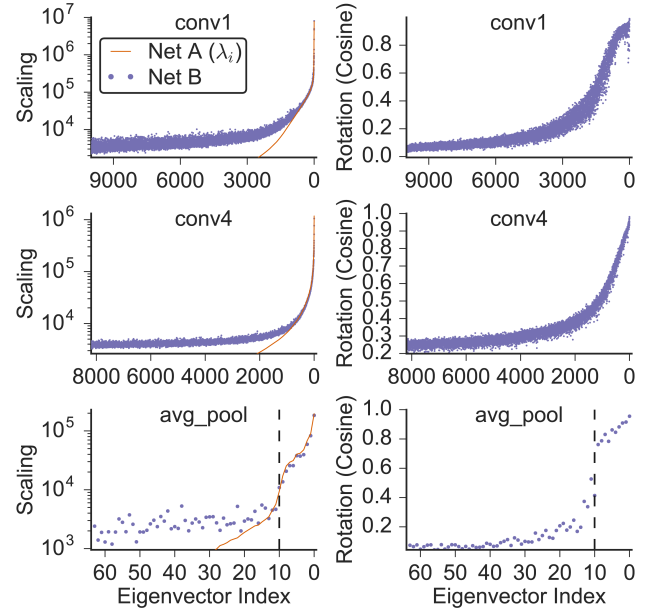


Figure 8. The shared subspace of two Tiny-10 networks trained from random initialization is spanned primarily by the eigenvectors corresponding to the largest eigenvalues. Each row represents a different network layer. Note that the average pooling layer has only 64 units. **Left:** Scaling of the eigenvectors \mathbf{u}_X^i of the RSM XX^T from network A by RSMs of networks A and B. Orange lines show $\|XX^T\mathbf{u}_X^i\|_2$, i.e. the eigenvalues. Purple dots show $\|YY^T\mathbf{u}_X^i\|_2$, the scaling of the eigenvectors of the RSM of network A by the RSM of network B. **Right:** Cosine of the rotation by the RSM of network B, $(\mathbf{u}_X^i)^T YY^T \mathbf{u}_X^i / \|YY^T \mathbf{u}_X^i\|_2$.

index should focus on. Previous work has suggested that there is little similarity between intermediate layers of neural networks trained from different random initializations (Raghu et al., 2017; Wang et al., 2018). We propose CKA as a method for comparing representations of neural networks, and show that it consistently identifies correspondences between layers, not only in the same network trained from different initializations, but across entirely different architectures, whereas other methods do not. We also provide a unified framework for understanding the space of similarity indexes, as well as an empirical framework for evaluation.

We show that CKA captures intuitive notions of similarity, i.e. that neural networks trained from different initializations should be similar to each other. However, it remains an open question whether there exist kernels beyond the linear and RBF kernels that would be better for analyzing neural network representations. Moreover, there are other potential choices of weighting in Equation 14 that may be more appropriate in certain settings. We leave these questions as future work. Nevertheless, CKA seems to be much better than previous methods at finding correspondences between the learned representations in hidden layers of neural networks.

Acknowledgements

We thank Gamaleldin Elsayed, Jaehoon Lee, Paul-Henri Mignot, Maithra Raghu, Samuel L. Smith, and Alex Williams for comments on the manuscript, Rishabh Agarwal for ideas, and Aliza Elkin for support.

References

- Advani, M. S. and Saxe, A. M. High-dimensional dynamics of generalization error in neural networks. *arXiv preprint arXiv:1710.03667*, 2017.
- Amari, S.-i., Ozeki, T., Karakida, R., Yoshida, Y., and Okada, M. Dynamics of learning in mlp: Natural gradient and singularity revisited. *Neural Computation*, 30(1): 1–33, 2018.
- Björck, Å. and Golub, G. H. Numerical methods for computing angles between linear subspaces. *Mathematics of Computation*, 27(123):579–594, 1973.
- Chen, A. M., Lu, H.-m., and Hecht-Nielsen, R. On the geometry of feedforward neural network error surfaces. *Neural Computation*, 5(6):910–927, 1993.
- Connolly, A. C., Guntupalli, J. S., Gors, J., Hanke, M., Halchenko, Y. O., Wu, Y.-C., Abdi, H., and Haxby, J. V. The representation of biological classes in the human brain. *Journal of Neuroscience*, 32(8):2608–2618, 2012.
- Cortes, C., Mohri, M., and Rostamizadeh, A. Algorithms for learning kernels based on centered alignment. *Journal of Machine Learning Research*, 13(Mar):795–828, 2012.
- Cristianini, N., Shawe-Taylor, J., Elisseeff, A., and Kandola, J. S. On kernel-target alignment. In *Advances in Neural Information Processing Systems*, pp. 367–373, 2002.
- Dinh, L., Sohl-Dickstein, J., and Bengio, S. Density estimation using real nvp. In *International Conference on Learning Representations*, 2017.
- Dumoulin, V., Shlens, J., and Kudlur, M. A learned representation for artistic style. *International Conference on Learning Representations*, 2, 2017.
- Edelman, S. Representation is representation of similarities. *Behavioral and Brain Sciences*, 21(4):449–467, 1998.
- Elsayed, G. F., Lara, A. H., Kaufman, M. T., Churchland, M. M., and Cunningham, J. P. Reorganization between preparatory and movement population responses in motor cortex. *Nature Communications*, 7:13239, 2016.
- Freiwald, W. A. and Tsao, D. Y. Functional compartmentalization and viewpoint generalization within the macaque face-processing system. *Science*, 330(6005):845–851, 2010.
- Garriga-Alonso, A., Rasmussen, C. E., and Aitchison, L. Deep convolutional networks as shallow gaussian processes. In *International Conference on Learning Representations*, 2019.
- Gatys, L. A., Ecker, A. S., and Bethge, M. Image style transfer using convolutional neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2414–2423, 2016.
- Golub, G. H. and Zha, H. The canonical correlations of matrix pairs and their numerical computation. In *Linear Algebra for Signal Processing*, pp. 27–49. Springer, 1995.
- Gretton, A., Bousquet, O., Smola, A., and Schölkopf, B. Measuring statistical dependence with hilbert-schmidt norms. In *International Conference on Algorithmic Learning Theory*, pp. 63–77. Springer, 2005.
- Haxby, J. V., Gobbini, M. I., Furey, M. L., Ishai, A., Schouten, J. L., and Pietrini, P. Distributed and overlapping representations of faces and objects in ventral temporal cortex. *Science*, 293(5539):2425–2430, 2001.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 770–778, 2016.
- Ioffe, S. and Szegedy, C. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International Conference on Machine Learning*, pp. 448–456, 2015.
- Jacobsen, J.-H., Smeulders, A. W., and Oyallon, E. i-RevNet: Deep invertible networks. In *International Conference on Learning Representations*, 2018.
- Jacot, A., Gabriel, F., and Hongler, C. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in Neural Information Processing Systems*, pp. 8571–8580, 2018.
- Johnson, J., Alahi, A., and Fei-Fei, L. Perceptual losses for real-time style transfer and super-resolution. In *European Conference on Computer Vision*, pp. 694–711. Springer, 2016.
- 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.
- Kriegeskorte, N., Mur, M., and Bandettini, P. A. Representational similarity analysis-connecting the branches of systems neuroscience. *Frontiers in Systems Neuroscience*, 2:4, 2008a.

- Kriegeskorte, N., Mur, M., Ruff, D. A., Kiani, R., Bodurka, J., Esteky, H., Tanaka, K., and Bandettini, P. A. Matching categorical object representations in inferior temporal cortex of man and monkey. *Neuron*, 60(6):1126–1141, 2008b.
- Kuss, M. and Graepel, T. The geometry of kernel canonical correlation analysis. Technical report, Max Planck Institute for Biological Cybernetics, 2003.
- Laakso, A. and Cottrell, G. Content and cluster analysis: assessing representational similarity in neural systems. *Philosophical Psychology*, 13(1):47–76, 2000.
- LeCun, Y., Kanter, I., and Solla, S. A. Second order properties of error surfaces: Learning time and generalization. In *Advances in Neural Information Processing Systems*, pp. 918–924, 1991.
- Lee, J., Sohl-dickstein, J., Pennington, J., Novak, R., Schoenholz, S., and Bahri, Y. Deep neural networks as gaussian processes. In *International Conference on Learning Representations*, 2018.
- Li, Y., Yosinski, J., Clune, J., Lipson, H., and Hopcroft, J. Convergent learning: Do different neural networks learn the same representations? In Storcheus, D., Ros-tamizadeh, A., and Kumar, S. (eds.), *Proceedings of the 1st International Workshop on Feature Extraction: Modern Questions and Challenges at NIPS 2015*, volume 44 of *Proceedings of Machine Learning Research*, pp. 196–212, Montreal, Canada, 11 Dec 2015. PMLR.
- Lorenzo-Seva, U. and Ten Berge, J. M. Tucker’s congruence coefficient as a meaningful index of factor similarity. *Methodology*, 2(2):57–64, 2006.
- Morcos, A., Raghu, M., and Bengio, S. Insights on representational similarity in neural networks with canonical correlation. *Advances in Neural Information Processing Systems 31*, pp. 5732–5741, 2018.
- Mroueh, Y., Marcheret, E., and Goel, V. Asymmetrically weighted CCA and hierarchical kernel sentence embedding for multimodal retrieval. *arXiv preprint arXiv:1511.06267*, 2015.
- Novak, R., Xiao, L., Bahri, Y., Lee, J., Yang, G., Abolafia, D. A., Pennington, J., and Sohl-dickstein, J. Bayesian deep convolutional networks with many channels are gaussian processes. In *International Conference on Learning Representations*, 2019.
- Orhan, E. and Pitkow, X. Skip connections eliminate singularities. In *International Conference on Learning Representations*, 2018.
- Press, W. H. Canonical correlation clarified by singular value decomposition, 2011. URL <http://numerical.recipes/whp/notes/CanonCorrBySVD.pdf>.
- Raghu, M., Gilmer, J., Yosinski, J., and Sohl-Dickstein, J. Svcca: Singular vector canonical correlation analysis for deep learning dynamics and interpretability. In Guyon, I., Luxburg, U. V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems 30*, pp. 6076–6085. Curran Associates, Inc., 2017.
- Ramsay, J., ten Berge, J., and Styán, G. Matrix correlation. *Psychometrika*, 49(3):403–423, 1984.
- Robert, P. and Escoufier, Y. A unifying tool for linear multivariate statistical methods: the RV-coefficient. *Applied Statistics*, pp. 257–265, 1976.
- Romero, A., Ballas, N., Kahou, S. E., Chassang, A., Gatta, C., and Bengio, Y. Fitnets: Hints for thin deep nets. In *International Conference on Learning Representations*, 2015.
- SAS Institute. *Introduction to Regression Procedures*. 2015. URL <https://support.sas.com/documentation/onlinedoc/stat/141/introreg.pdf>.
- Saxe, A. M., McClelland, J. L., and Ganguli, S. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. In *International Conference on Learning Representations*, 2014.
- Sejdinovic, D., Sriperumbudur, B., Gretton, A., and Fukumizu, K. Equivalence of distance-based and rkhs-based statistics in hypothesis testing. *The Annals of Statistics*, pp. 2263–2291, 2013.
- Smith, S. L., Turban, D. H., Hamblin, S., and Hammerla, N. Y. Offline bilingual word vectors, orthogonal transformations and the inverted softmax. In *International Conference on Learning Representations*, 2017.
- Song, L., Smola, A., Gretton, A., Borgwardt, K. M., and Bedo, J. Supervised feature selection via dependence estimation. In *Proceedings of the 24th international conference on Machine learning*, pp. 823–830. ACM, 2007.
- Springenberg, J. T., Dosovitskiy, A., Brox, T., and Ried-miller, M. Striving for simplicity: The all convolutional net. In *International Conference on Learning Representations Workshop*, 2015.
- StataCorp. *Stata Multivariate Statistics Reference Manual*. 2015. URL <https://www.stata.com/manuals14/mv.pdf>.

- Sussillo, D., Churchland, M. M., Kaufman, M. T., and Shenoy, K. V. A neural network that finds a naturalistic solution for the production of muscle activity. *Nature Neuroscience*, 18(7):1025, 2015.
- Tucker, L. R. A method for synthesis of factor analysis studies. Technical report, Educational Testing Service, Princeton, NJ, 1951.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. Attention is all you need. In *Advances in Neural Information Processing Systems*, pp. 5998–6008, 2017.
- Vinod, H. D. Canonical ridge and econometrics of joint production. *Journal of Econometrics*, 4(2):147–166, 1976.
- Wang, L., Hu, L., Gu, J., Wu, Y., Hu, Z., He, K., and Hopcroft, J. E. Towards understanding learning representations: To what extent do different neural networks learn the same representation. In *Advances in Neural Information Processing Systems*, pp. 9607–9616, 2018.
- 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.
- Zagoruyko, S. and Komodakis, N. Wide residual networks. In Richard C. Wilson, E. R. H. and Smith, W. A. P. (eds.), *Proceedings of the British Machine Vision Conference (BMVC)*, pp. 87.1–87.12. BMVA Press, September 2016. ISBN 1-901725-59-6. doi: 10.5244/C.30.87.

A. Proof of Theorem 1

Theorem. Let X and Y be $n \times p$ matrices. Suppose s is invariant to invertible linear transformation in the first argument, i.e. $s(X, Z) = s(XA, Z)$ for arbitrary Z and any A with $\text{rank}(A) = p$. If $\text{rank}(X) = \text{rank}(Y) = n$, then $s(X, Z) = s(Y, Z)$.

Proof. Let

$$X' = \begin{bmatrix} X \\ K_X \end{bmatrix} \quad Y' = \begin{bmatrix} Y \\ K_Y \end{bmatrix}$$

where K_X is a basis for the null space of the rows of X and K_Y is a basis for the null space of the rows of Y . Then let $A = X'^{-1}Y'$.

$$\begin{bmatrix} X \\ K_X \end{bmatrix} A = \begin{bmatrix} Y \\ K_Y \end{bmatrix} \implies XA = Y$$

Because X' and Y' have rank p by construction, A also has rank p . Thus, $s(X, Z) = s(XA, Z) = s(Y, Z)$. \square

B. Orthogonalization and Invariance to Invertible Linear Transformation

Here we show that any similarity index that is invariant to orthogonal transformation can be made invariant to invertible linear transformation by orthogonalizing the columns of the input.

Proposition 1. Let X be an $n \times p$ matrix of full column rank and let A be an invertible $p \times p$ matrix. Let $X = Q_X R_X$ and $XA = Q_{XA} R_{XA}$, where $Q_X^T Q_X = Q_{XA}^T Q_{XA} = I$ and R_X and R_{XA} are invertible. If $s(\cdot, \cdot)$ is invariant to orthogonal transformation, then $s(Q_X, Y) = s(Q_{XA}, Y)$.

Proof. Let $B = R_X A R_{XA}^{-1}$. Then $Q_X B = Q_{XA}$, and B is an orthogonal transformation:

$$B^T B = B^T Q_X^T Q_X B = Q_{XA}^T Q_{XA} = I$$

Thus $s(Q_X, Y) = s(Q_X B, Y) = s(Q_{XA}, Y)$. \square

C. CCA and Linear Regression

C.1. Linear Regression

Consider the linear regression fit of the columns of an $n \times m$ matrix C with an $n \times p$ matrix A :

$$\hat{B} = \arg \min_B \|C - AB\|_F^2 = (A^T A)^{-1} A^T C$$

Let $A = QR$, the thin QR decomposition of A . Then the fitted values are given by:

$$\begin{aligned} \hat{C} &= A \hat{B} \\ &= A(A^T A)^{-1} A^T C \\ &= QR(R^T Q^T QR)^{-1} R^T Q^T C \\ &= QRR^{-1}(R^T)^{-1} R^T Q^T C \\ &= QQ^T C \end{aligned}$$

The residuals $E = C - \hat{C}$ are orthogonal to the fitted values, i.e.

$$\begin{aligned} E^T \hat{C} &= (C - QQ^T C)^T QQ^T C \\ &= C^T QQ^T C - C^T QQ^T C = 0 \end{aligned}$$

Thus:

$$\begin{aligned}
 \|E\|_F^2 &= \text{tr}(E^T E) \\
 &= \text{tr}(E^T C - E^T \hat{C}) \\
 &= \text{tr}((C - \hat{C})^T C) \\
 &= \text{tr}(C^T C) - \text{tr}(C^T Q Q^T C) \\
 &= \|C\|_F^2 - \|Q^T C\|_F^2
 \end{aligned} \tag{15}$$

Assuming that C was centered by subtracting its column means prior to the linear regression fit, the total fraction of variance explained by the fit is:

$$R^2 = 1 - \frac{\|E\|_F^2}{\|C\|_F^2} = 1 - \frac{\|C\|_F^2 - \|Q^T C\|_F^2}{\|C\|_F^2} = \frac{\|Q^T C\|_F^2}{\|C\|_F^2} \tag{16}$$

Although we have assumed that Q is obtained from QR decomposition, any orthonormal basis with the same span will suffice, because orthogonal transformations do not change the Frobenius norm.

C.2. CCA

Let X be an $n \times p_1$ matrix and Y be an $n \times p_2$ matrix, and let $p = \min(p_1, p_2)$. Given the thin QR decompositions of X and Y , $X = Q_X R_X$, $Y = Q_Y R_Y$ such that $Q_X^T Q_X = I$, $Q_Y^T Q_Y = I$, the canonical correlations ρ_i are the singular values of $A = Q_X^T Q_Y$ (Björck & Golub, 1973; Press, 2011) and thus the square roots of the eigenvalues of $A^T A$. The squared canonical correlations ρ_i^2 are the eigenvalues of $A^T A = Q_Y^T Q_X Q_X^T Q_Y$. Their sum is $\sum_{i=1}^p \rho_i^2 = \text{tr}(A^T A) = \|Q_Y^T Q_X\|_F^2$.

Now consider the linear regression fit of the columns of Q_X with Y . Assume that Q_X has zero mean. Substituting Q_Y for Q and Q_X for C in Equation 16, and noting that $\|Q_X\|_F^2 = p_1$:

$$R^2 = \frac{\|Q_Y^T Q_X\|_F^2}{p_1} = \frac{\sum_{i=1}^p \rho_i^2}{p_1} \tag{17}$$

C.3. Projection-Weighted CCA

Morcos et al. (2018) proposed to compute projection-weighted canonical correlation as:

$$\bar{\rho}_{\text{PW}} = \frac{\sum_{i=1}^c \alpha_i \rho_i}{\sum_{i=1}^c \alpha_i} \quad \alpha_i = \sum_j |\langle \mathbf{h}_i, \mathbf{x}_j \rangle|$$

where the \mathbf{x}_j are the columns of X , and the \mathbf{h}_i are the canonical variables formed by projecting X to the canonical coordinate frame. Below, we show that if we modify $\bar{\rho}_{\text{PW}}$ by squaring the dot products and ρ_i , we recover linear regression. Specifically:

$$R_{\text{MPW}}^2 = \frac{\sum_{i=1}^c \alpha'_i \rho_i^2}{\sum_{i=1}^c \alpha'_i} = R_{\text{LR}}^2 \quad \alpha'_i = \sum_j \langle \mathbf{h}_i, \mathbf{x}_j \rangle^2$$

Our derivation begins by forming the SVD $Q_X^T Q_Y = U \Sigma V^T$. Σ is a diagonal matrix of the canonical correlations ρ_i , and the matrix of canonical variables $H = Q_X U$. Then R_{MPW}^2 is:

$$\begin{aligned}
 R_{\text{MPW}}^2 &= \frac{\|X^T H \Sigma\|_F^2}{\|X^T H\|_F^2} \\
 &= \frac{\text{tr}((X^T H \Sigma)^T (X^T H \Sigma))}{\text{tr}((X^T H)^T (X^T H))} \\
 &= \frac{\text{tr}(\Sigma H^T X X^T H \Sigma)}{\text{tr}(H^T X X^T H)} \\
 &= \frac{\text{tr}(X^T H \Sigma^2 H^T X)}{\text{tr}(X^T H H^T X)} \\
 &= \frac{\text{tr}(R_X^T Q_X^T H \Sigma^2 H^T Q_X R_X)}{\text{tr}(R_X^T Q_X^T Q_X U U^T Q_X^T Q_X R_X)}
 \end{aligned} \tag{18}$$

Noting that $Q_X^T H = U$ and $U\Sigma = Q_X^T Q_Y V$:

$$\begin{aligned}
 R_{\text{MPW}}^2 &= \frac{\text{tr}(R_X^T U \Sigma^2 U^T R_X)}{\text{tr}(R_X^T Q_X^T Q_X R_X)} \\
 &= \frac{\text{tr}(R_X^T Q_X^T Q_Y V \Sigma U^T R_X)}{\text{tr}(X^T X)} \\
 &= \frac{\text{tr}(X^T Q_Y Q_Y^T Q_X R_X)}{\text{tr}(X^T X)} \\
 &= \frac{\text{tr}(X^T Q_Y Q_Y^T X)}{\text{tr}(X^T X)} \\
 &= \frac{\|Q_Y^T X\|_F^2}{\|X\|_F^2}
 \end{aligned}$$

Substituting Q_Y for Q and X for C in Equation 16:

$$R_{\text{LR}}^2 = \frac{\|Q_Y^T X\|_F^2}{\|X\|_F^2} = R_{\text{MPW}}^2$$

D. Notes on Other Methods

D.1. Canonical Ridge

Beyond CCA, we could also consider the “canonical ridge” regularized CCA objective (Vinod, 1976):

$$\begin{aligned}
 \sigma_i &= \max_{\mathbf{w}_X^i, \mathbf{w}_Y^i} \frac{(X \mathbf{w}_X^i)^T (Y \mathbf{w}_Y^i)}{\sqrt{\|X \mathbf{w}_X^i\|^2 + \kappa_X \|\mathbf{w}_X^i\|_2^2} \sqrt{\|Y \mathbf{w}_Y^i\|^2 + \kappa_Y \|\mathbf{w}_Y^i\|_2^2}} \\
 \text{subject to } &\forall_{j < i} (\mathbf{w}_X^i)^T (X^T X + \kappa_X I) \mathbf{w}_X^j = 0 \\
 &\forall_{j < i} (\mathbf{w}_Y^i)^T (Y^T Y + \kappa_Y I) \mathbf{w}_Y^j = 0
 \end{aligned} \tag{19}$$

Given the singular value decompositions $X = U_X \Sigma_X V_X^T$ and $Y = U_Y \Sigma_Y V_Y^T$, one can form “partially orthogonalized” bases $\tilde{Q}_X = U_X \Sigma_X (\Sigma_X^2 + \kappa_X I)^{-1/2}$ and $\tilde{Q}_Y = U_Y \Sigma_Y (\Sigma_Y^2 + \kappa_Y I)^{-1/2}$. Given the singular value decomposition of their product $\tilde{U} \tilde{\Sigma} \tilde{V}^T = \tilde{Q}_X^T \tilde{Q}_Y$, the canonical weights are given by $W_X = V_X (\Sigma_X^2 + \kappa_X I)^{-1/2} \tilde{U}$ and $W_Y = V_Y (\Sigma_Y^2 + \kappa_Y I)^{-1/2} \tilde{V}$, as previously shown by Mroueh et al. (2015). As in the unregularized case (Equation 13), there is a convenient expression for the sum of the squared singular values $\sum \tilde{\sigma}_i^2$ in terms of the eigenvalues and eigenvectors of XX^T and YY^T . Let the i^{th} left-singular vector of X (eigenvector of XX^T) be indexed as \mathbf{u}_X^i and let the i^{th} eigenvalue of XX^T (squared singular value of X) be indexed as λ_X^i , and similarly let the left-singular vectors of YY^T be indexed as \mathbf{u}_Y^i and the eigenvalues as λ_Y^i . Then:

$$\sum_{i=1}^{p_1} \tilde{\sigma}_i^2 = \|\tilde{Q}_Y^T \tilde{Q}_X\|_F^2 \tag{20}$$

$$= \|(\Sigma_Y^2 + \kappa_Y I)^{-1/2} \Sigma_Y U_Y^T U_X \Sigma_X (\Sigma_X^2 + \kappa_X I)^{-1/2}\|_F^2 \tag{21}$$

$$= \sum_{i=1}^{p_1} \sum_{j=1}^{p_2} \frac{\lambda_X^i \lambda_Y^j}{(\lambda_X^i + \kappa_X)(\lambda_Y^j + \kappa_Y)} \langle \mathbf{u}_X^i, \mathbf{u}_Y^j \rangle^2 \tag{22}$$

Unlike in the unregularized case, the singular values σ_i do not measure the correlation between the canonical variables. Instead, they become arbitrarily small as κ_X or κ_Y increase. Thus, we need to normalize the statistic to remove the dependency on the regularization parameters.

Applying von Neumann's trace inequality yields a bound:

$$\sum_{i=1}^{p_1} \tilde{\sigma}_i^2 = \text{tr}(\tilde{Q}_Y \tilde{Q}_Y^T \tilde{Q}_X \tilde{Q}_X^T) \quad (23)$$

$$= \text{tr}((U_Y \Sigma_Y^2 (\Sigma_Y^2 + \kappa_Y I)^{-1} U_Y^T)(U_X \Sigma_X^2 (\Sigma_X^2 + \kappa_X I)^{-1} U_X^T)) \quad (24)$$

$$\leq \sum_{i=1}^{p_1} \frac{\lambda_X^i \lambda_Y^i}{(\lambda_X^i + \kappa_X)(\lambda_Y^i + \kappa_Y)} \quad (25)$$

Applying the Cauchy-Schwarz inequality to (25) yields the alternative bounds:

$$\sum_{i=1}^{p_1} \tilde{\sigma}_i^2 \leq \sqrt{\sum_{i=1}^{p_1} \left(\frac{\lambda_X^i}{\lambda_X^i + \kappa_X} \right)^2} \sqrt{\sum_{i=1}^{p_1} \left(\frac{\lambda_Y^i}{\lambda_Y^i + \kappa_Y} \right)^2} \quad (26)$$

$$\leq \sqrt{\sum_{i=1}^{p_1} \left(\frac{\lambda_X^i}{\lambda_X^i + \kappa_X} \right)^2} \sqrt{\sum_{i=1}^{p_2} \left(\frac{\lambda_Y^i}{\lambda_Y^i + \kappa_Y} \right)^2} \quad (27)$$

A normalized form of (22) could be produced by dividing by any of (25), (26), or (27).

If $\kappa_X = \kappa_Y = 0$, then (25) and (26) are equal to p_1 . In this case, (22) is simply the sum of the squared canonical correlations, so normalizing by either of these bounds recovers R_{CCA}^2 .

If $\kappa_Y = 0$, then as $\kappa_X \rightarrow \infty$, normalizing by the bound from (25) recovers R^2 :

$$\lim_{\kappa_X \rightarrow \infty} \frac{\sum_{i=1}^{p_1} \sum_{j=1}^{p_2} \frac{\lambda_X^i \lambda_Y^j}{(\lambda_X^i + \kappa_X)(\lambda_Y^j + 0)} \langle \mathbf{u}_X^i, \mathbf{u}_Y^j \rangle^2}{\sum_{i=1}^{p_1} \frac{\lambda_X^i \lambda_Y^i}{(\lambda_X^i + \kappa_X)(\lambda_Y^i + 0)}} \quad (28)$$

$$= \lim_{\kappa_X \rightarrow \infty} \frac{\sum_{i=1}^{p_1} \sum_{j=1}^{p_2} \frac{\lambda_X^i}{\left(\frac{\lambda_X^i}{\kappa_X} + 1 \right)} \langle \mathbf{u}_X^i, \mathbf{u}_Y^j \rangle^2}{\sum_{i=1}^{p_1} \frac{\lambda_X^i}{\left(\frac{\lambda_X^i}{\kappa_X} + 1 \right)}} \quad (29)$$

$$= \frac{\sum_{i=1}^{p_1} \sum_{j=1}^{p_2} \lambda_X^i \langle \mathbf{u}_X^i, \mathbf{u}_Y^j \rangle^2}{\sum_{i=1}^{p_1} \lambda_X^i} \quad (30)$$

$$= \frac{\|U_Y^T U_X \Sigma_X\|_F^2}{\|X\|_F^2} = \frac{\|Q_Y^T X\|_F^2}{\|X\|_F^2} = R_{\text{LR}}^2 \quad (31)$$

The bound from (27) differs from the bounds in (25) and (26) because it is multiplicatively separable in X and Y . Normalizing by this bound leads to CKA($\tilde{Q}_X \tilde{Q}_X^T, \tilde{Q}_Y \tilde{Q}_Y^T$):

$$\frac{\sum_{i=1}^{p_1} \sum_{j=1}^{p_2} \frac{\lambda_X^i \lambda_Y^j}{(\lambda_X^i + \kappa_X)(\lambda_Y^j + \kappa_Y)} \langle \mathbf{u}_X^i, \mathbf{u}_Y^j \rangle^2}{\sqrt{\sum_{i=1}^{p_1} \left(\frac{\lambda_X^i}{\lambda_X^i + \kappa_X} \right)^2} \sqrt{\sum_{i=1}^{p_2} \left(\frac{\lambda_Y^i}{\lambda_Y^i + \kappa_Y} \right)^2}} \quad (32)$$

$$= \frac{\|\tilde{Q}_Y^T \tilde{Q}_X\|_F^2}{\|\tilde{Q}_X^T \tilde{Q}_X\|_F \|\tilde{Q}_Y^T \tilde{Q}_Y\|_F} = \text{CKA}(\tilde{Q}_X \tilde{Q}_X^T, \tilde{Q}_Y \tilde{Q}_Y^T) \quad (33)$$

Moreover, setting $\kappa_X = \kappa_Y = \kappa$ and taking the limit as $\kappa \rightarrow \infty$, the normalization from (27) leads to $\text{CKA}(XX^T, YY^T)$:

$$\lim_{\kappa \rightarrow \infty} \frac{\sum_{i=1}^{p_1} \sum_{j=1}^{p_2} \frac{\lambda_X^i \lambda_Y^j}{(\lambda_X^i + \kappa)(\lambda_Y^j + \kappa)} \langle \mathbf{u}_X^i, \mathbf{u}_Y^j \rangle^2}{\sqrt{\sum_{i=1}^{p_1} \left(\frac{\lambda_X^i}{\lambda_X^i + \kappa} \right)^2} \sqrt{\sum_{j=1}^{p_2} \left(\frac{\lambda_Y^j}{\lambda_Y^j + \kappa} \right)^2}} \quad (34)$$

$$= \lim_{\kappa \rightarrow \infty} \frac{\sum_{i=1}^{p_1} \sum_{j=1}^{p_2} \frac{\lambda_X^i \lambda_Y^j}{\left(\frac{\lambda_X^i}{\kappa} + 1 \right) \left(\frac{\lambda_Y^j}{\kappa} + 1 \right)} \langle \mathbf{u}_X^i, \mathbf{u}_Y^j \rangle^2}{\sqrt{\sum_{i=1}^{p_1} \left(\frac{\lambda_X^i}{\frac{\lambda_X^i}{\kappa} + 1} \right)^2} \sqrt{\sum_{j=1}^{p_2} \left(\frac{\lambda_Y^j}{\frac{\lambda_Y^j}{\kappa} + 1} \right)^2}} \quad (35)$$

$$= \frac{\sum_{i=1}^{p_1} \sum_{j=1}^{p_2} \lambda_X^i \lambda_Y^j \langle \mathbf{u}_X^i, \mathbf{u}_Y^j \rangle^2}{\sqrt{\sum_{i=1}^{p_1} (\lambda_X^i)^2} \sqrt{\sum_{j=1}^{p_2} (\lambda_Y^j)^2}} \quad (36)$$

$$= \text{CKA}(XX^T, YY^T)$$

Overall, the hyperparameters of the canonical ridge objective make it less useful for exploratory analysis. These hyperparameters could be selected by cross-validation, but this is computationally expensive, and the resulting estimator would be biased by sample size. Moreover, our goal is not to map representations of networks to a common space, but to measure the similarity between networks. Appropriately chosen regularization will improve out-of-sample performance of the mapping, but it makes the meaning of “similarity” more ambiguous.

D.2. The Orthogonal Procrustes Problem

The orthogonal Procrustes problem consists of finding an orthogonal rotation in feature space that produces the smallest error:

$$\hat{Q} = \arg \min_Q \|Y - XQ\|_F^2 \text{ subject to } Q^T Q = I \quad (37)$$

The objective can be written as:

$$\begin{aligned} \|Y - XQ\|_F^2 &= \text{tr}((Y - XQ)^T(Y - XQ)) \\ &= \text{tr}(Y^T Y) - \text{tr}(Y^T XQ) - \text{tr}(Q^T X^T Y) + \text{tr}(Q^T X^T XQ) \\ &= \|Y\|_F^2 + \|X\|_F^2 - 2\text{tr}(Y^T XQ) \end{aligned} \quad (38)$$

Thus, an equivalent objective is:

$$\hat{Q} = \arg \max_Q \text{tr}(Y^T XQ) \text{ subject to } Q^T Q = I \quad (39)$$

The solution is $\hat{Q} = UV^T$ where $U\Sigma V^T = X^T Y$, the singular value decomposition. At the maximum of (39):

$$\text{tr}(Y^T X\hat{Q}) = \text{tr}(V\Sigma U^T UV^T) = \text{tr}(\Sigma) = \|X^T Y\|_* = \|Y^T X\|_* \quad (40)$$

which is similar to what we call “dot product-based similarity” (Equation 1), but with the squared Frobenius norm of $Y^T X$ (the sum of the squared singular values) replaced by the nuclear norm (the sum of the singular values). The Frobenius norm of $Y^T X$ can be obtained as the solution to a similar optimization problem:

$$\|Y^T X\|_F = \max_W \text{tr}(Y^T XW) \text{ subject to } \text{tr}(W^T W) = 1 \quad (41)$$

In the context of neural networks, [Smith et al. \(2017\)](#) previously proposed using the solution to the orthogonal Procrustes problem to align word embeddings from different languages, and demonstrated that it outperformed CCA.

E. Architecture Details

All non-ResNet architectures are based on All-CNN-C (Springenberg et al., 2015), but none are architecturally identical. The Plain-10 model is very similar, but we place the final linear layer after the average pooling layer and use batch normalization because these are common choices in modern architectures. We use these models because they train in minutes on modern hardware.

Tiny-10
3×3 conv. 16-BN-ReLu $\times 2$
3×3 conv. 32 stride 2-BN-ReLu
3×3 conv. 32-BN-ReLu $\times 2$
3×3 conv. 64 stride 2-BN-ReLu
3×3 conv. 64 valid padding-BN-ReLu
1×1 conv. 64-BN-ReLu
Global average pooling
Logits

Table E.1. The Tiny-10 architecture, used in Figures 2, 8, F.3, . The average Tiny-10 model achieved 89.4% accuracy.

Plain- $(8n + 2)$
3×3 conv. 96-BN-ReLu $\times (3n - 1)$
3×3 conv. 96 stride 2-BN-ReLu
3×3 conv. 192-BN-ReLu $\times (3n - 1)$
3×3 conv. 192 stride 2-BN-ReLu
3×3 conv. 192 BN-ReLu $\times (n - 1)$
3×3 conv. 192 valid padding-BN-ReLu
1×1 conv. 192-BN-ReLu $\times n$
Global average pooling
Logits

Table E.2. The Plain- $(8n + 2)$ architecture, used in Figures 3, 5, 7, F.4, F.5, F.6, and F.7. Mean accuracies: Plain-10, 93.9%; Plain-18: 94.8%; Plain-34: 93.7%; Plain-66: 91.3%

Width- n
3×3 conv. n -BN-ReLu $\times 2$
3×3 conv. n stride 2-BN-ReLu
3×3 conv. n -BN-ReLu $\times 2$
3×3 conv. n stride 2-BN-ReLu
3×3 conv. n valid padding-BN-ReLu
1×1 conv. n -BN-ReLu
Global average pooling
Logits

Table E.3. The architectures used for width experiments in Figure 6.

F. Additional Experiments

F.1. Sanity Check for Transformer Encoders

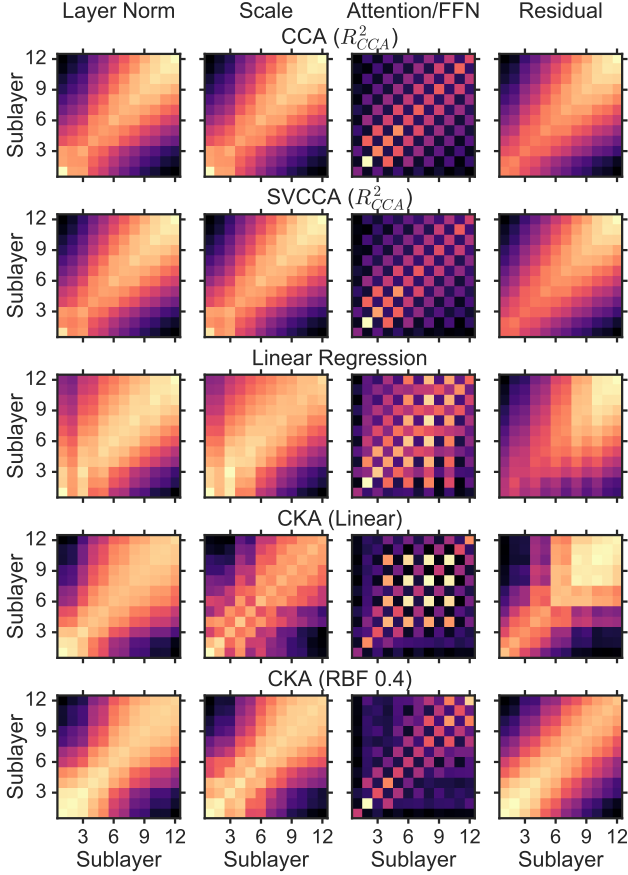


Figure F.1. All similarity indices broadly reflect the structure of Transformer encoders. Similarity indexes are computed between the 12 sublayers of Transformer encoders, for each of the 4 possible places in each sublayer that representations may be taken (see Figure F.2), averaged across 10 models trained from different random initializations.

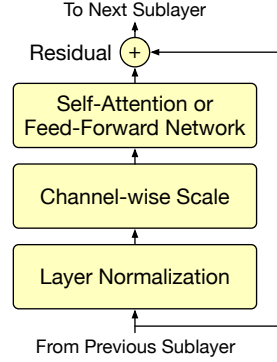


Figure F.2. Architecture of a single sublayer of the Transformer encoder used for our experiments. The full encoder includes 12 sublayers, alternating between self-attention and feed-forward network sublayers.

Index	Layer Norm	Scale	Attn/FFN	Residual
CCA ($\bar{\rho}$)	85.3	85.3	94.9	90.9
CCA (R^2_{CCA})	87.8	87.8	95.3	95.2
SVCCA ($\bar{\rho}$)	78.2	83.0	89.5	75.9
SVCCA (R^2_{CCA})	85.4	86.9	90.8	84.7
PWCCA	88.5	88.9	96.1	87.0
Linear Reg.	78.1	83.7	76.0	36.9
CKA (Linear)	78.6	95.6	86.0	73.6
CKA (RBF 0.2)	76.5	73.1	70.5	76.2
CKA (RBF 0.4)	92.3	96.5	89.1	98.1
CKA (RBF 0.8)	80.8	95.8	93.6	90.0

Table F.1. Accuracy of identifying corresponding sublayers based maximum similarity, for 10 architecturally identical 12-sublayer Transformer encoders at the 4 locations in each sublayer after which the representation may be taken (see Figure F.2). Results not significantly different from the best result are bold-faced ($p < 0.05$, jackknife z-test).

When applied to Transformer encoders, all similarity indexes we investigated passed the sanity check described in Section 6.1. In Figure F.1, we show similarity between the 12 sublayers of the encoders of 10 Transformer models (45 pairs) (Vaswani et al., 2017) trained from different random initializations to perform English to German translation. Each Transformer sublayer contains four operations, shown in Figure F.2, and results vary based which operation the representation is taken after. Table F.1 shows the accuracy with which we can identify corresponding layers between network pairs by maximal similarity.

The Transformer architecture alternates between self-attention and feed-forward network sublayers. The checkerboard pattern in similarity plots for the Attention/FFN layer in Figure F.1 indicates that representations of feed-forward network sublayers are more similar to other feed-forward network sublayers than to self-attention sublayers, and similarly, representations of self-attention sublayers are more similar to other self-attention sublayers than to feed-forward network layers. CKA also reveals a checkerboard pattern for activations after the channel-wise scale operation (before the self-attention/feed-forward network operation) that other methods do not. Because CCA is invariant to non-isotropic scaling, CCA similarities before and after channel-wise scaling are identical. Thus, CCA cannot capture this structure, even though it is consistent across different networks.

F.2. SVCCA at Alternative Thresholds

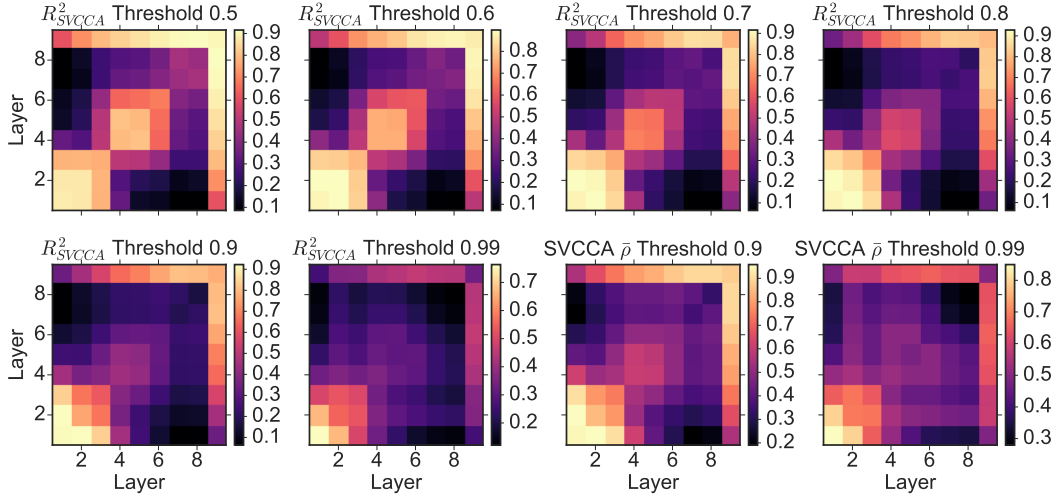


Figure F.3. Same as Figure 2 row 2, but for more SVCCA thresholds than the 0.99 threshold suggested by Raghu et al. (2017). No threshold reveals the structure of the network.

F.3. CKA at Initialization

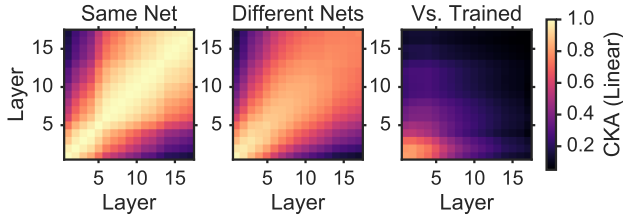


Figure F.4. Similarity of the Plain-18 network at initialization. **Left:** Similarity between layers of the same network. **Middle:** Similarity between untrained networks with different initializations. **Right:** Similarity between untrained and trained networks.

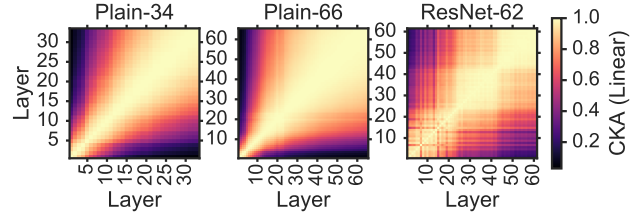


Figure F.5. Similarity between layers at initialization for deeper architectures.

F.4. Additional CKA Results

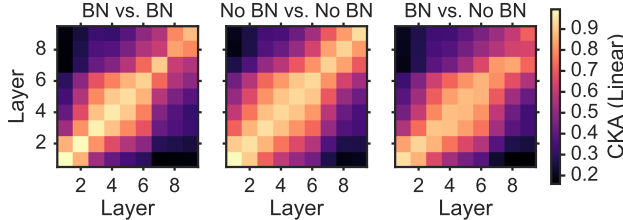


Figure F.6. Networks with and without batch normalization trained from different random initializations learn similar representations according to CKA. The largest difference between networks is at the last convolutional layer. Optimal hyperparameters were separately selected for the batch normalized network (93.9% average accuracy) and the network without batch normalization (91.5% average accuracy).

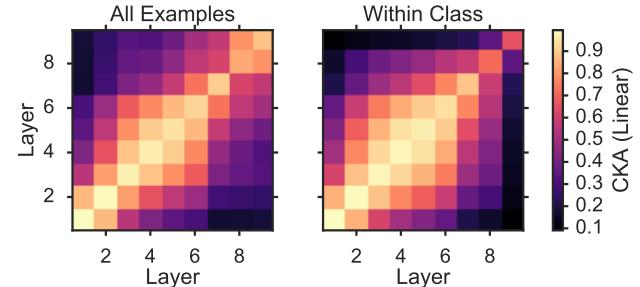


Figure F.7. Within-class CKA is similar to CKA based on all examples. To measure within-class CKA, we computed CKA separately for examples belonging to each CIFAR-10 class based on representations from Plain-10 networks, and averaged the resulting CKA values across classes.

F.5. Similarity Between Different Architectures with Other Indexes

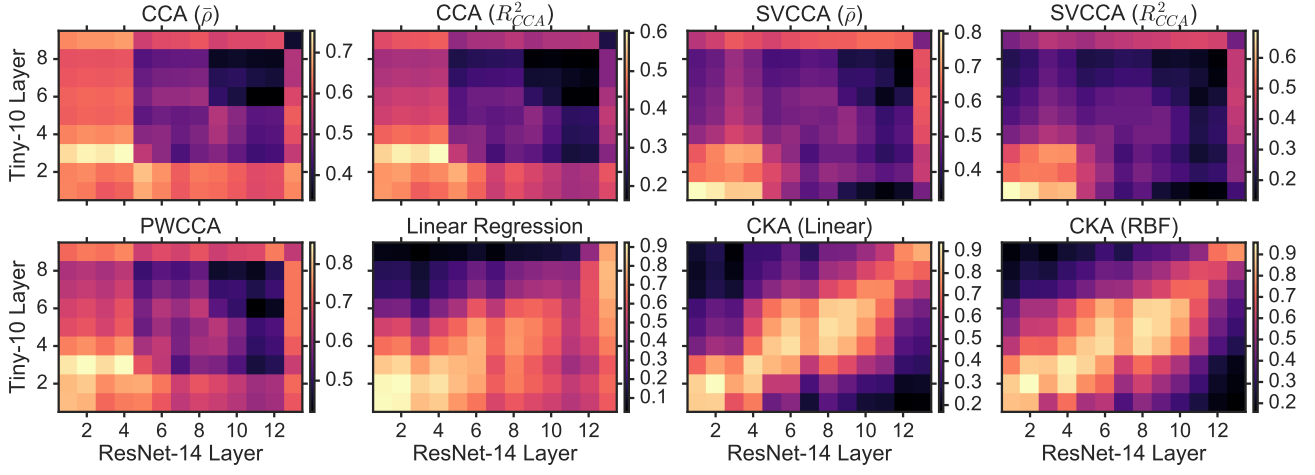


Figure F.8. Similarity between layers of different architectures (Tiny-10 and ResNet-14) for all methods investigated. Only CKA reveals meaningful correspondence. SVCCA results resemble Figure 7 of Raghu et al. (2017). In order to achieve better performance for CCA-based techniques, which are sensitive to the number of examples used to compute similarity, all plots show similarity on the CIFAR-10 training set.