

## 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  $\rho_i$  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  $\rho_i^2$  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  $\rho_i$ , 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  $Q_X^T Q_Y = U \Sigma V^T$ .  $\Sigma$  is a diagonal matrix of the canonical correlations  $\rho_i$ , and the matrix of canonical variables  $H = Q_X U$ . Then  $R_{\text{MPW}}^2$  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):

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

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  $U \Sigma \tilde{V}^T = \tilde{Q}_X^T \tilde{Q}_Y$ , 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  $i^{\text{th}}$  left-singular vector of  $X$  (eigenvector of  $XX^T$ ) be indexed as  $\mathbf{u}_X^i$  and let the  $i^{\text{th}}$  eigenvalue of  $XX^T$  (squared singular value of  $X$ ) be indexed as  $\lambda_X^i$ , and similarly let the left-singular vectors of  $YY^T$  be indexed as  $\mathbf{u}_Y^i$  and the eigenvalues as  $\lambda_Y^i$ . 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  $\sigma_i$  do not measure the correlation between the canonical variables. Instead, they become arbitrarily small as  $\kappa_X$  or  $\kappa_Y$  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  $R_{\text{CCA}}^2$ .

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

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

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

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

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

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( $\tilde{Q}_X \tilde{Q}_X^T, \tilde{Q}_Y \tilde{Q}_Y^T$ ):

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

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

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

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

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

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

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

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
$3 \times 3$ conv. 16-BN-ReLu $\times 2$
$3 \times 3$ conv. 32 stride 2-BN-ReLu
$3 \times 3$ conv. 32-BN-ReLu $\times 2$
$3 \times 3$ conv. 64 stride 2-BN-ReLu
$3 \times 3$ conv. 64 valid padding-BN-ReLu
$1 \times 1$ 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- $(8n + 2)$
$3 \times 3$ conv. 96-BN-ReLu $\times (3n - 1)$
$3 \times 3$ conv. 96 stride 2-BN-ReLu
$3 \times 3$ conv. 192-BN-ReLu $\times (3n - 1)$
$3 \times 3$ conv. 192 stride 2-BN-ReLu
$3 \times 3$ conv. 192 BN-ReLu $\times (n - 1)$
$3 \times 3$ conv. 192 valid padding-BN-ReLu
$1 \times 1$ 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- $n$
$3 \times 3$ conv. $n$ -BN-ReLu $\times 2$
$3 \times 3$ conv. $n$ stride 2-BN-ReLu
$3 \times 3$ conv. $n$ -BN-ReLu $\times 2$
$3 \times 3$ conv. $n$ stride 2-BN-ReLu
$3 \times 3$ conv. $n$ valid padding-BN-ReLu
$1 \times 1$ 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

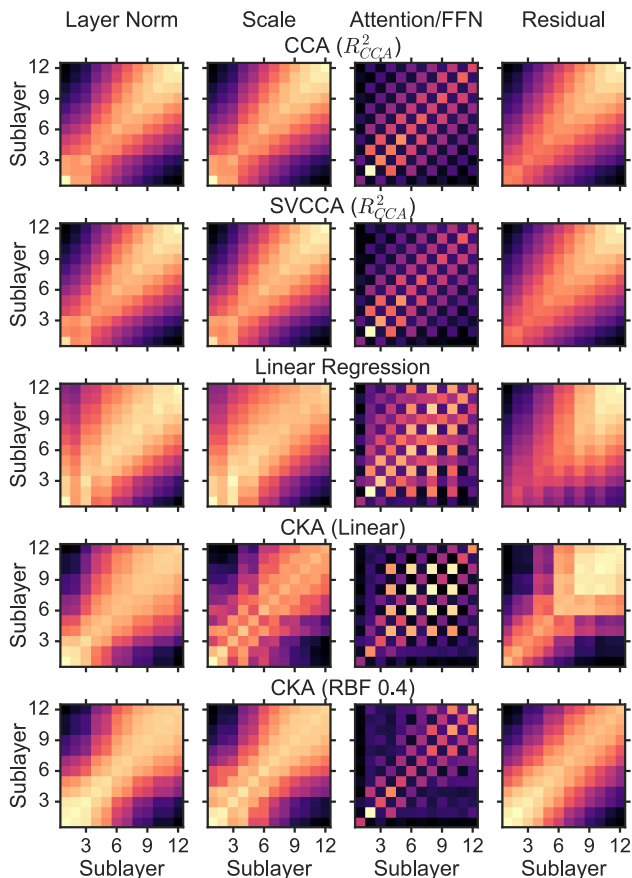


Figure F.1. All similarity indices broadly reflect the structure of Transformer encoders. Similarity indexes are computed between the 12 sublayers of Transformer encoders, for each of the 4 possible places in each sublayer that representations may be taken (see Figure F.2), averaged across 10 models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

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.

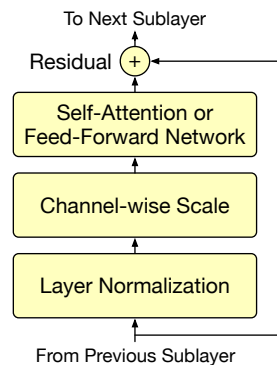


Figure F.2. Architecture of a single sublayer of the Transformer encoder used for our experiments. The full encoder includes 12 sublayers, alternating between self-attention and feed-forward network sublayers.

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

**F.2. SVCCA at Alternative Thresholds**

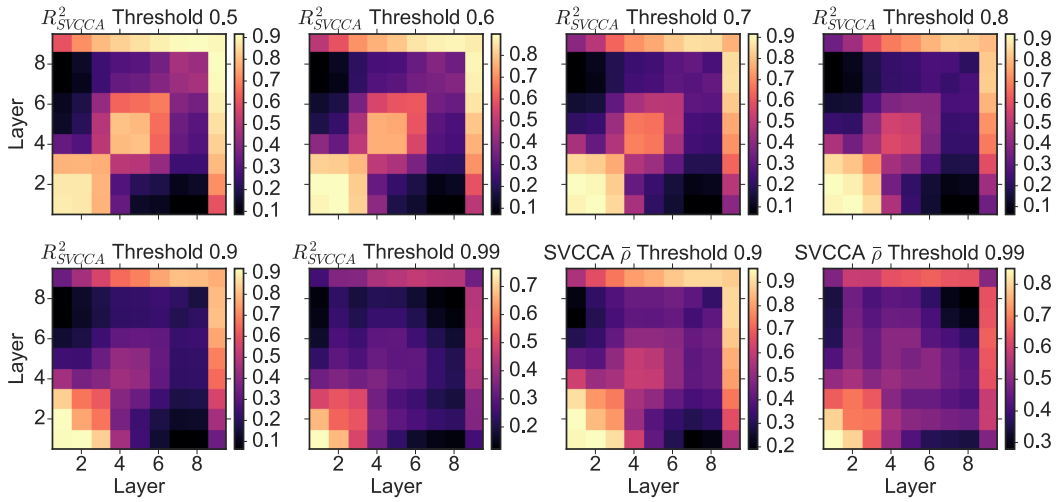


Figure F.3. Same as Figure 2 row 2, but for more SVCCA thresholds than the 0.99 threshold suggested by Raghu et al. (2017). No threshold reveals the structure of the network.

**F.3. CKA at Initialization**

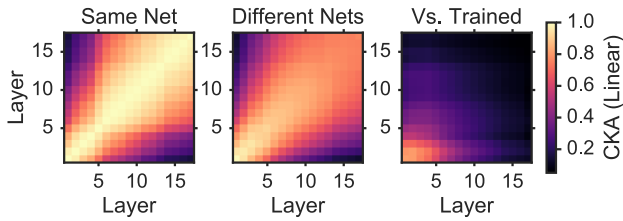


Figure F.4. Similarity of the Plain-18 network at initialization. **Left:** Similarity between layers of the same network. **Middle:** Similarity between untrained networks with different initializations. **Right:** Similarity between untrained and trained networks.

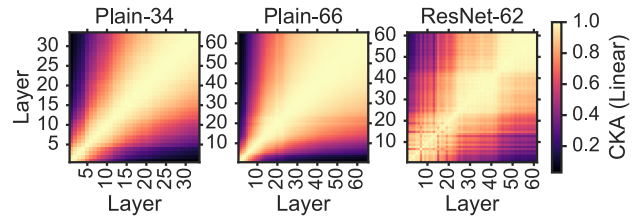


Figure F.5. Similarity between layers at initialization for deeper architectures.

**F.4. Additional CKA Results**

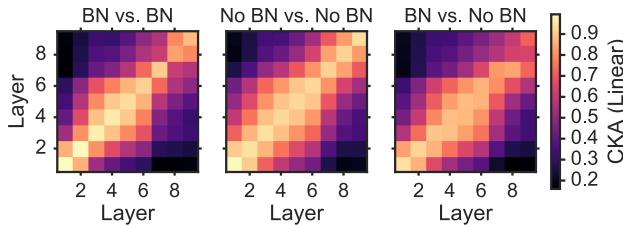


Figure F.6. Networks with and without batch normalization trained from different random initializations learn similar representations according to CKA. The largest difference between networks is at the last convolutional layer. Optimal hyperparameters were separately selected for the batch normalized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

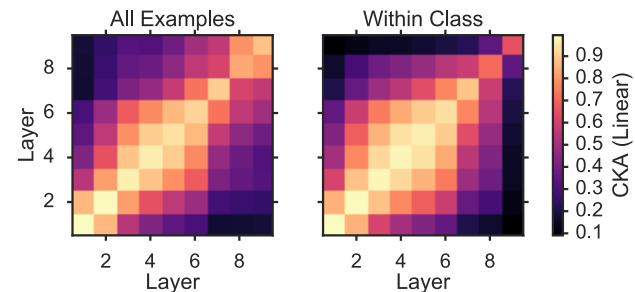


Figure F.7. Within-class CKA is similar to CKA based on all examples. To measure within-class CKA, we computed CKA separately for examples belonging to each CIFAR-10 class based on representations from Plain-10 networks, and averaged the resulting CKA values across classes.



E.5. Similarity Between Different Architectures with Other Indexes

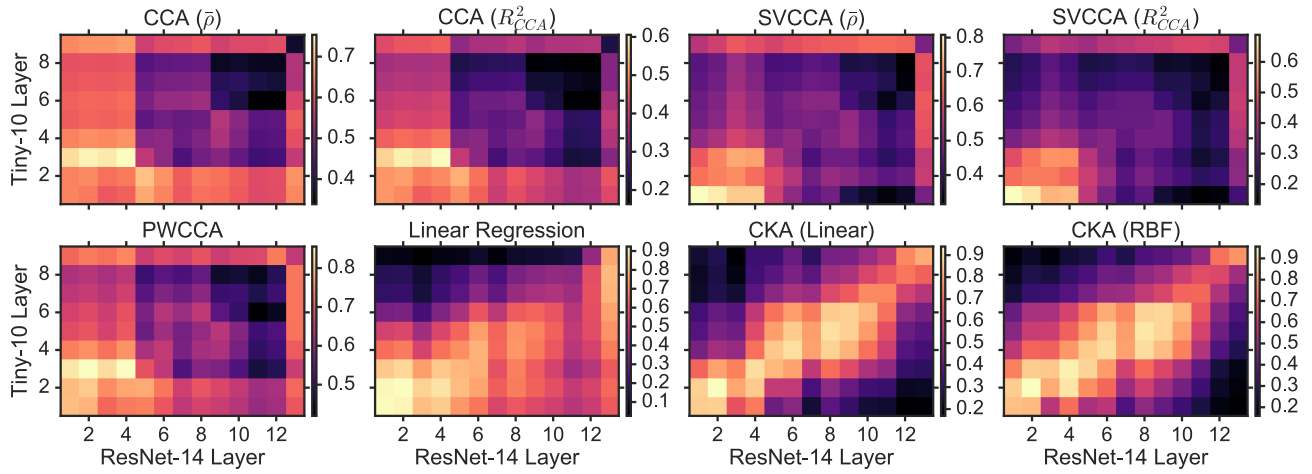


Figure F.8. Similarity between layers of different architectures (Tiny-10 and ResNet-14) for all methods investigated. Only CKA reveals meaningful correspondence. SVCCA results resemble Figure 7 of Raghu et al. (2017). In order to achieve better performance for CCA-based techniques, which are sensitive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