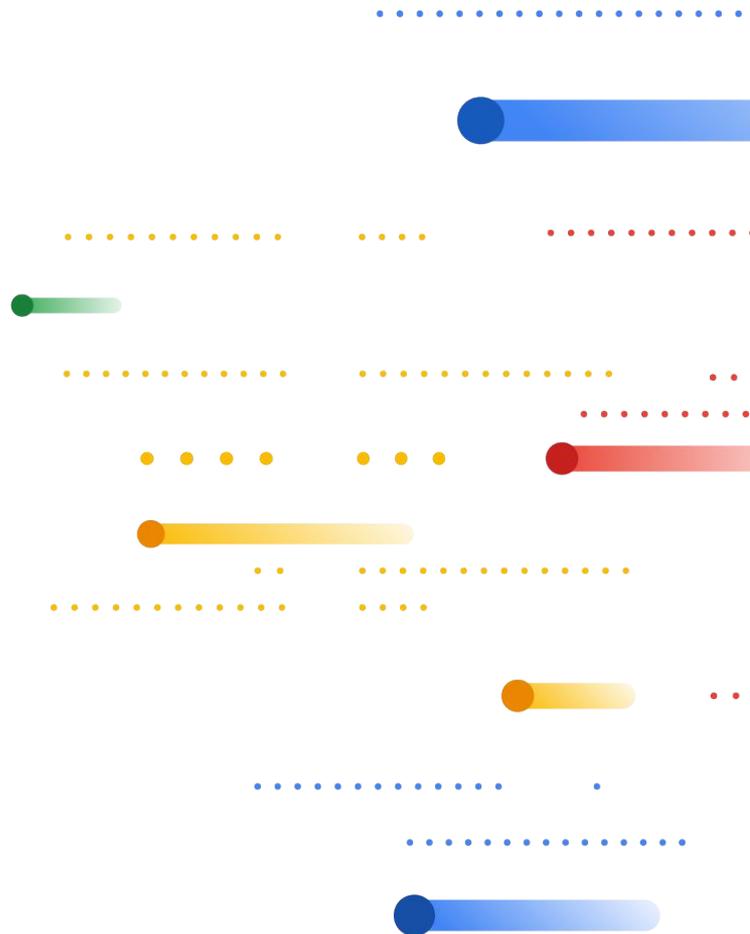


Similarity of Neural Network Representations Revisited

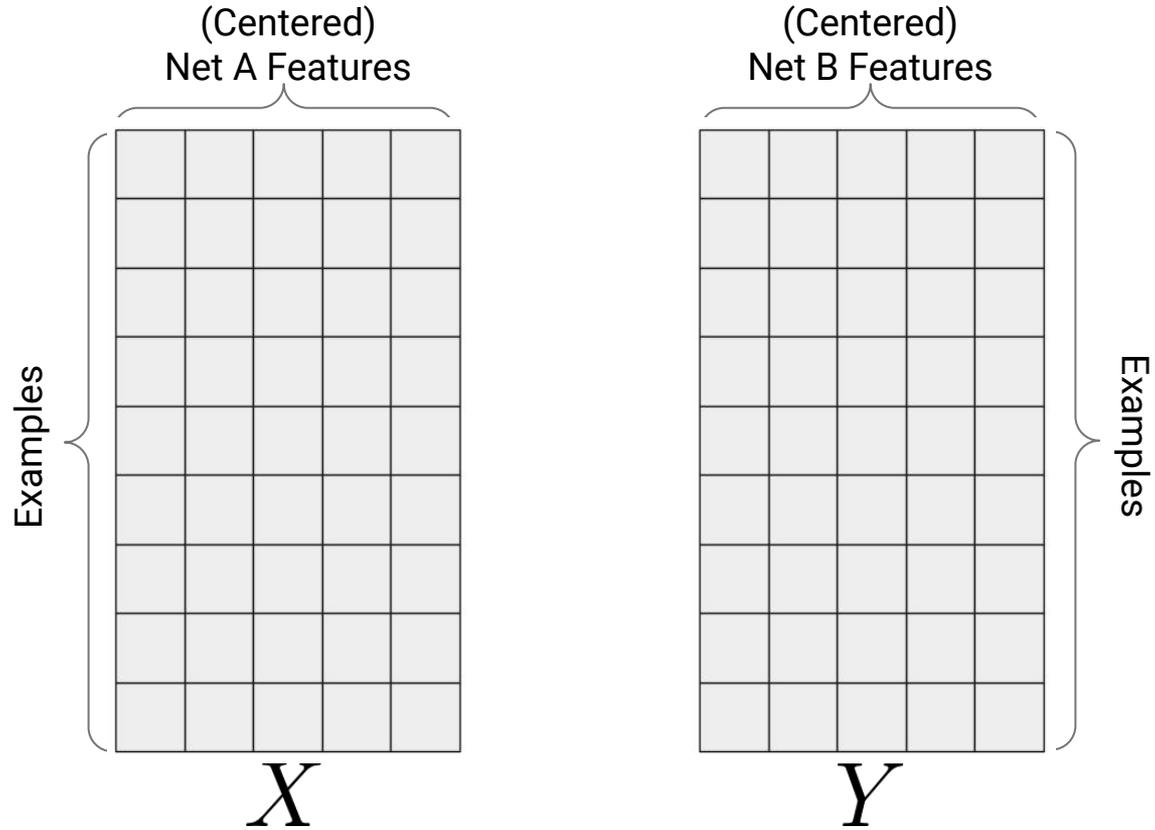
Simon Kornblith, Mohammad Norouzi, Honglak Lee, Geoffrey Hinton



Motivation

- We need tools to understand *trained* neural networks
 - Neural network training involves interactions between an algorithm and structured data
 - We don't know the structure of the data
- One way to understand trained neural networks is by comparing their representations

What is a Representation?



Comparing Features = Comparing Examples

$$\|X^T Y\|_F^2 = \langle \text{vec}(X X^T), \text{vec}(Y Y^T) \rangle$$

Sum of squared dot products
(similarities) between features

Dot product between reshaped inter-example
similarity matrices

Comparing Features = Comparing Examples

$$\|X^T Y\|_F^2 = \langle \text{vec}(X X^T), \text{vec}(Y Y^T) \rangle$$

$$\frac{\|X^T Y\|_F^2}{\|X^T X\|_F \|Y^T Y\|_F} = \frac{\langle \text{vec}(X X^T), \text{vec}(Y Y^T) \rangle}{\|X X^T\|_F \|Y Y^T\|_F}$$

Comparing Features = Comparing Examples

$$\|X^T Y\|_F^2 = \langle \text{vec}(X X^T), \text{vec}(Y Y^T) \rangle$$

$$\frac{\|X^T Y\|_F^2}{\|X^T X\|_F \|Y^T Y\|_F} = \frac{\langle \text{vec}(X X^T), \text{vec}(Y Y^T) \rangle}{\|X X^T\|_F \|Y Y^T\|_F}$$

Centered kernel alignment (CKA) (Cortes et al., 2012)

RV-coefficient (Robert & Escoufier, 1976)

Tucker's congruence coefficient (Tucker, 1951)

The Kernel Trick

$$\frac{\langle \text{vec}(XX^T), \text{vec}(YY^T) \rangle}{\|XX^T\|_F \|YY^T\|_F} \rightarrow \frac{\langle \text{vec}(\tilde{K}), \text{vec}(\tilde{L}) \rangle}{\|\tilde{K}\|_F \|\tilde{L}\|_F}$$

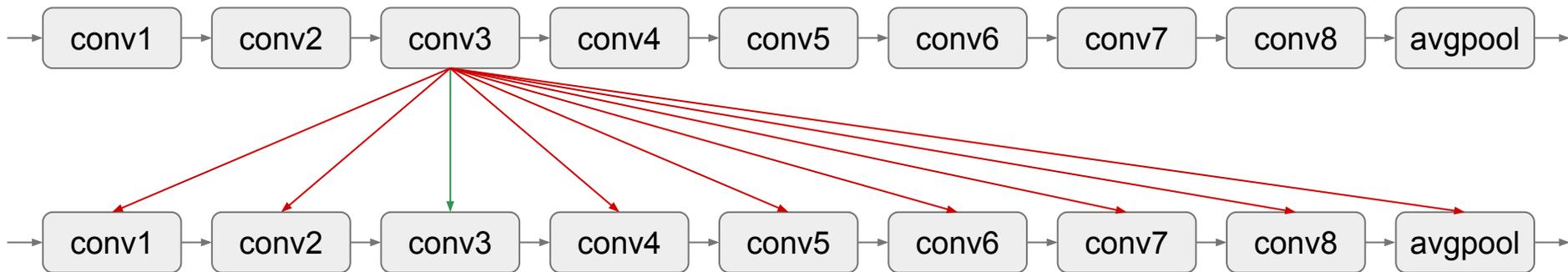
$$K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j) \quad \tilde{K} = HKH$$

$$L_{ij} = l(\mathbf{y}_i, \mathbf{y}_j) \quad \tilde{L} = HLH$$

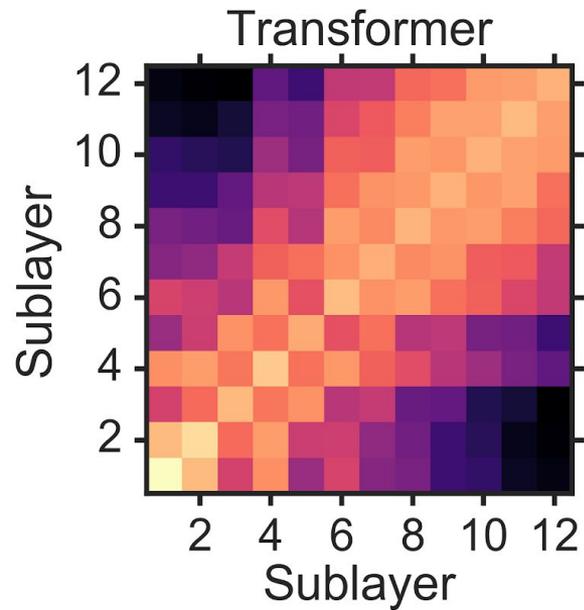
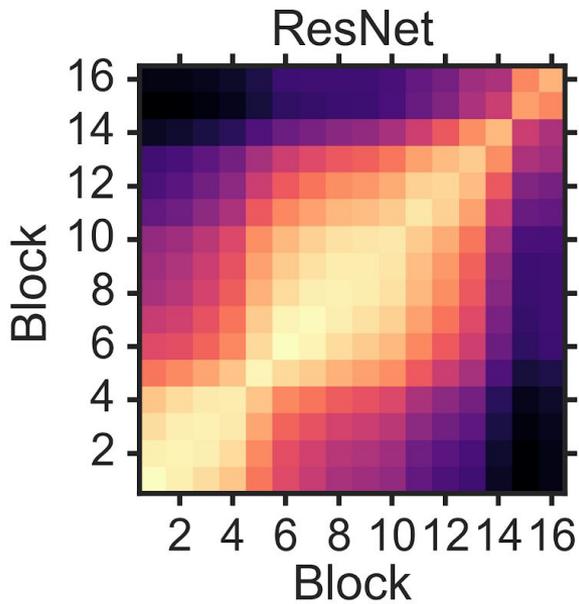
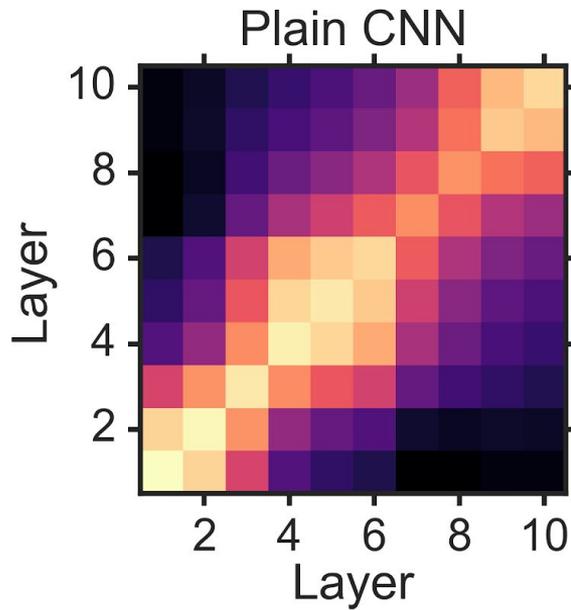
H is the centering matrix

A Sanity Check for Similarity

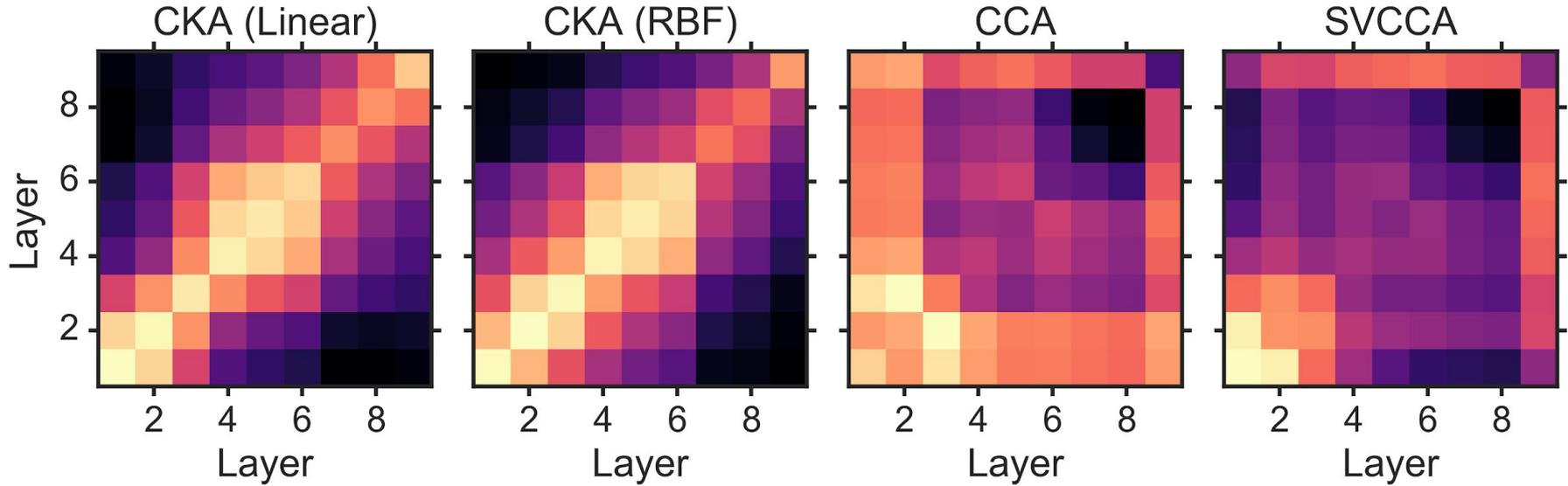
Given two architecturally identical networks A and B trained from different random initializations, a layer from net A should be most similar to the architecturally corresponding layer in net B



A Sanity Check for Similarity



A Sanity Check for Similarity



CKA Reveals Network Pathology

1x Depth (94.1% on CIFAR-10)



CKA Reveals Network Pathology

1x Depth (94.1%)



2x Depth (95.0%)



CKA Reveals Network Pathology

1x Depth (94.1%)



2x Depth (95.0%)



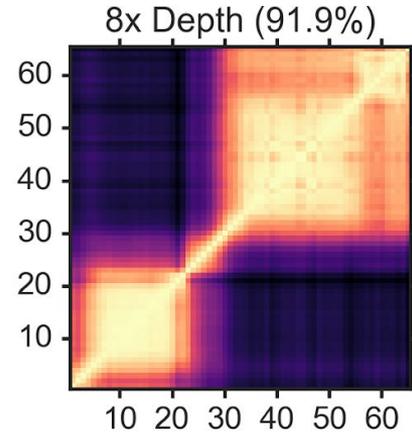
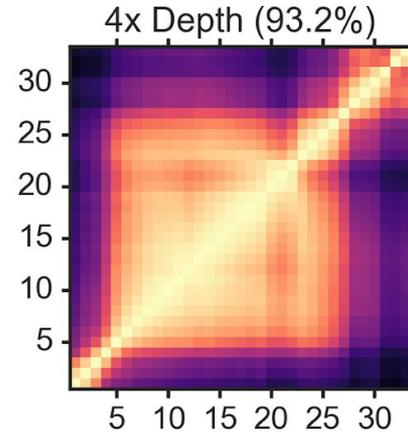
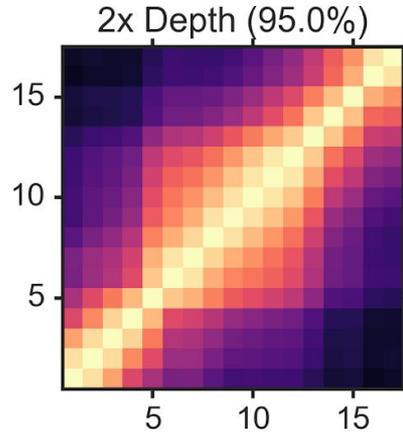
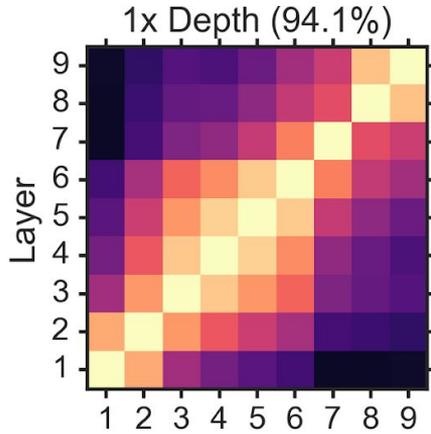
4x Depth (93.2%)



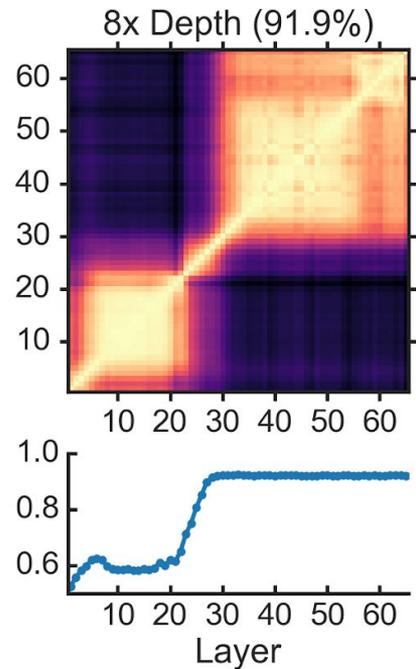
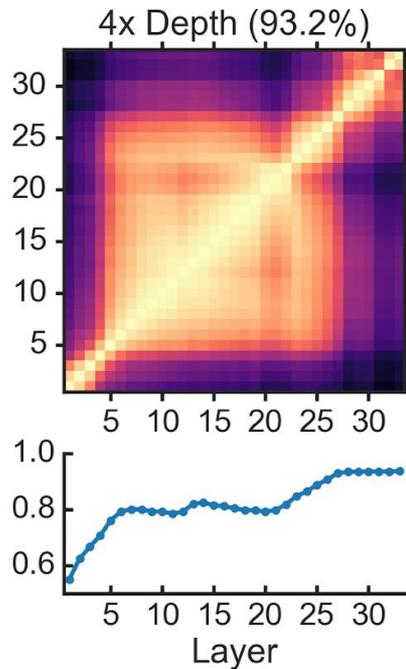
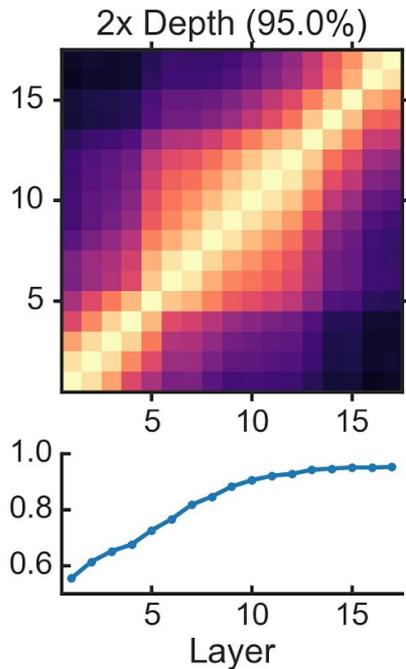
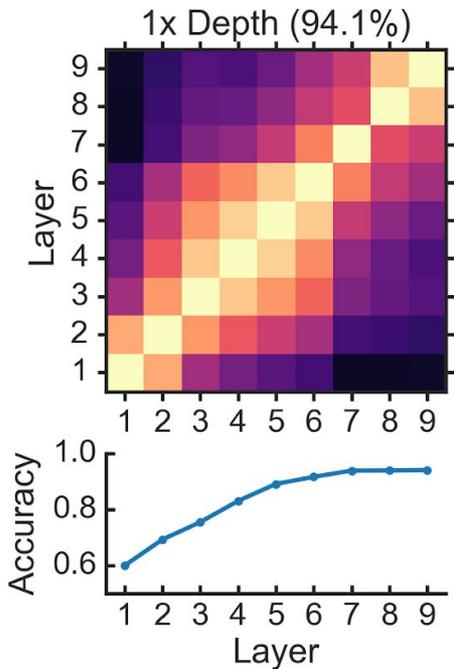
8x Depth (91.9%)



CKA Reveals Network Pathology



CKA Reveals Network Pathology



Thank You!

[cka-similarity.github.io](https://github.com/cka-similarity)