

# From **CNN** to **ViT**: Evaluating the Impact of Attention Block On the Learned Internal Representation

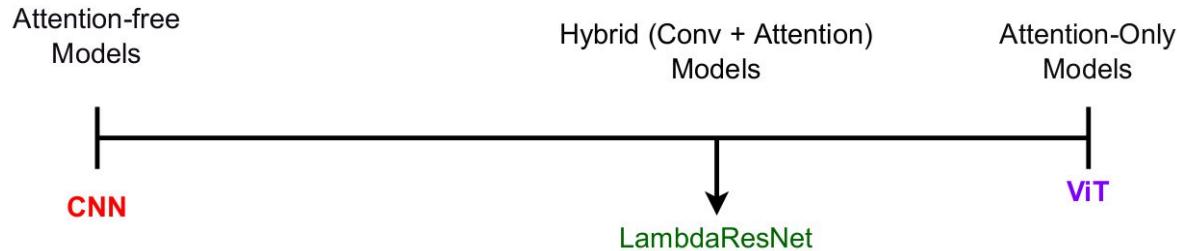
Animesh Basak Chowdhury (abc586)

Nandan Kumar Jha (nj2049)

# Executive summary

**Problem Statement (Goal):** Self-attention operator exhibits quadratic complexity.  
How much Self-attention do we really need?

**Solution/Approach:** Investigate the entire spectrum of **CNN**(attention-free) model to **Lambda Networks** (conv + self-attention) to **ViT**(pure-attention) w/ HSIC metric

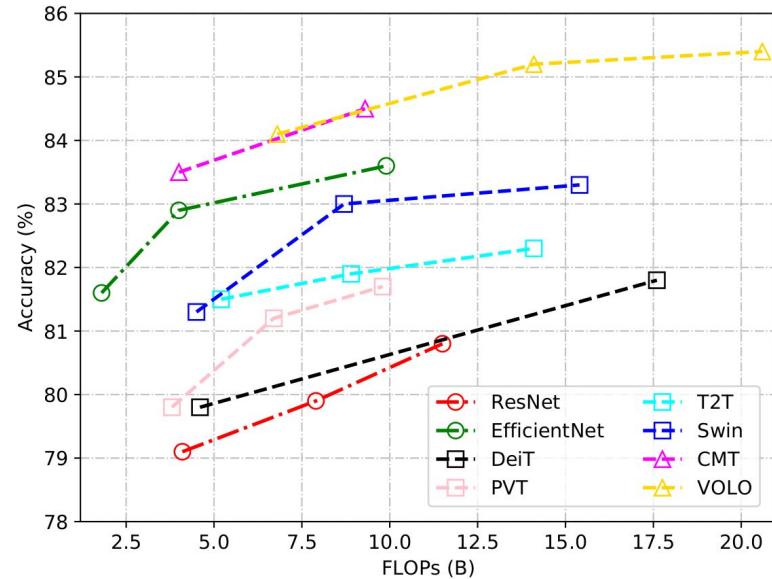
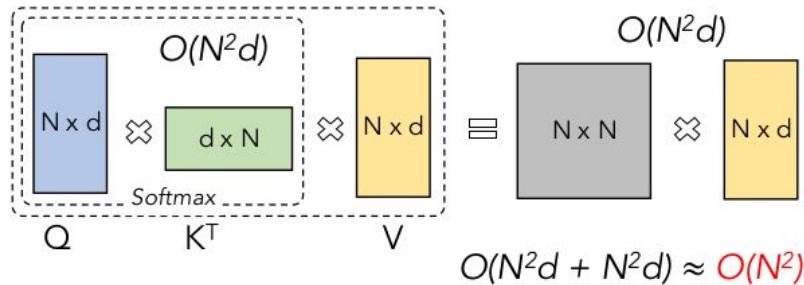


**Value/Benefit:** Visualizing the (learned) internal representation is a *self-explanatory, effective, and robust* way for evaluating the impact of different architectures and operators

# Problem Statement: CNN vs Transformer

CNNs are **inferior** to pure-attention based transformer models

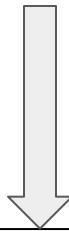
**Self-attention: Quadratic** computational and memory complexity



[Han, TPAMI'22]

# Problem Statement: How good are the hybrid models?

**1Q:** Do we really need pure attention based models?

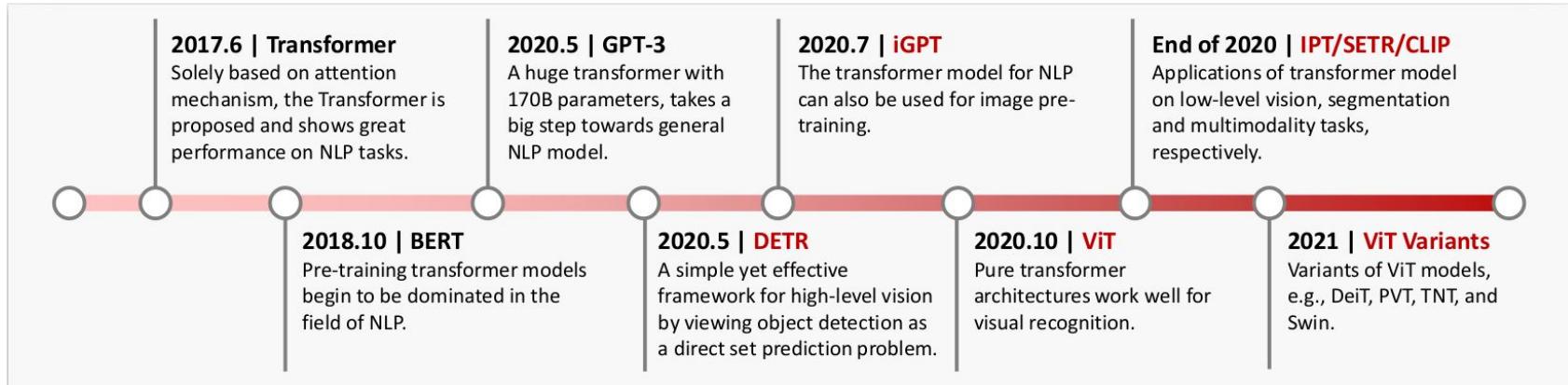


**2Q:** How good are the Hybrid models (CNN w/ self-attention in fewer layers)?

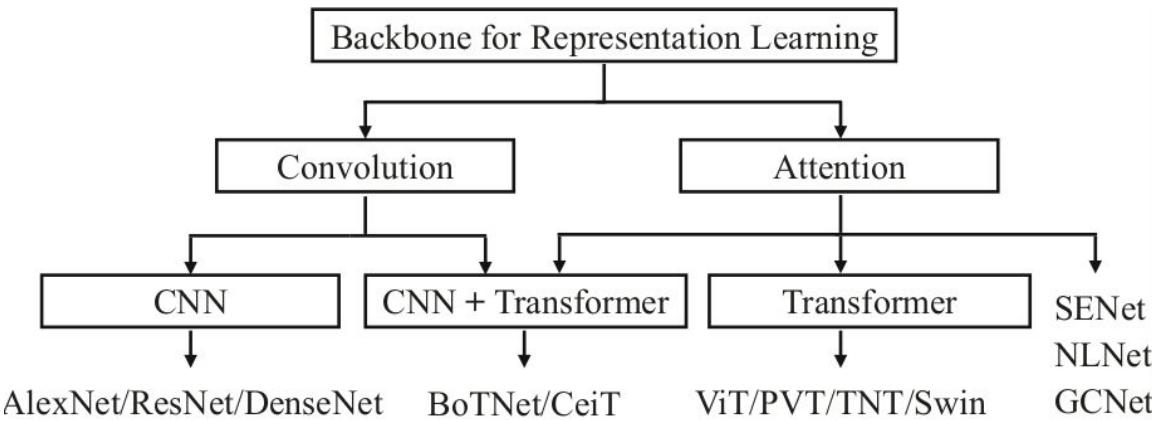


**3Q:** How self-attention changes the (learned)internal representation?

# Background: Model Complexity (CNN vs ViT)

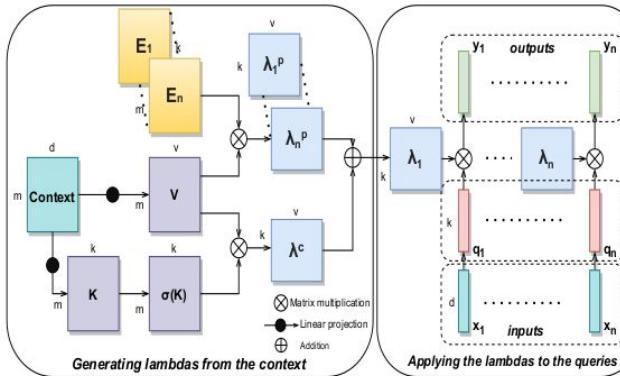


[A Survey on Vision Transformer: [Han, TPAMI'22](#)]



# Background: Hybrid Network w/ Linear Self-attention [Bello, ICLR'21]

Architecture	Params (M)	Throughput	top-1
C → C → C → C	25.6	7240 ex/s	76.9
L → C → C → C	25.5	1880 ex/s	77.3
L → L → C → C	25.0	1280 ex/s	77.2
L → L → L → C	21.7	1160 ex/s	77.8
L → L → L → L	15.0	1160 ex/s	78.4
C → L → L → L	15.1	2200 ex/s	78.3
C → C → L → L	15.4	4980 ex/s	78.3
C → C → C → L	18.8	7160 ex/s	77.3



Name	Description
$ k ,  v $	query, value depth
$\mathbf{X} \in \mathbb{R}^{n \times d}$	inputs
$\mathbf{C} \in \mathbb{R}^{m \times d}$	context
$\mathbf{Q} = \mathbf{X}\mathbf{W}_Q \in \mathbb{R}^{n \times k}$	queries
$\mathbf{K} = \mathbf{C}\mathbf{W}_K \in \mathbb{R}^{m \times k}$	keys
$\mathbf{V} = \mathbf{C}\mathbf{W}_V \in \mathbb{R}^{m \times v}$	values
$\sigma(\mathbf{K}) = \text{softmax}(\mathbf{K}, \text{axis}=m)$	normalized keys
$\mathbf{E}_n \in \mathbb{R}^{m \times k}$	relative position embeddings
$\mathbf{\lambda}^c = \bar{\mathbf{K}}^T \mathbf{V} \in \mathbb{R}^{ k  \times  v }$	content lambda
$\mathbf{\lambda}^p = \mathbf{E}_n^T \mathbf{V} \in \mathbb{R}^{ k  \times  v }$	position lambdas
$\mathbf{\lambda}_n = \mathbf{\lambda}^c + \mathbf{\lambda}^p \in \mathbb{R}^{ k  \times  v }$	lambdas

Lambda convolution are used in **selective** layers only!

$$\boldsymbol{\lambda}_n = \sum_m (\bar{\mathbf{k}}_m + \mathbf{e}_{nm}) \mathbf{v}_m^T = \underbrace{\bar{\mathbf{K}}^T \mathbf{V}}_{\text{content lambda}} + \underbrace{\mathbf{E}_n^T \mathbf{V}}_{\text{position lambda}} \in \mathbb{R}^{|k| \times |v|}$$

Model	Block Configuration	Lambda layers in c4
LambdaResNet-50	[ 3-4-6-3 ]	3
LambdaResNet-101	[ 3-4-23-3 ]	6, 12, 18
LambdaResNet-152	[ 3-8-36-3 ]	5, 10, 15, 20, 25, 30
LambdaResNet-200	[ 3-24-36-3 ]	5, 10, 15, 20, 25, 30
LambdaResNet-270	[ 4-29-53-4 ]	8, 16, 24, 32, 40, 48
LambdaResNet-350	[ 4-36-72-4 ]	10, 20, 30, 40, 50, 60
LambdaResNet-420	[ 4-44-87-4 ]	10, 20, 30, 40, 50, 60, 70, 80

Most of the benefits come from the placing lambda-layers in the **last stages** of the networks.

# Technical Challenges

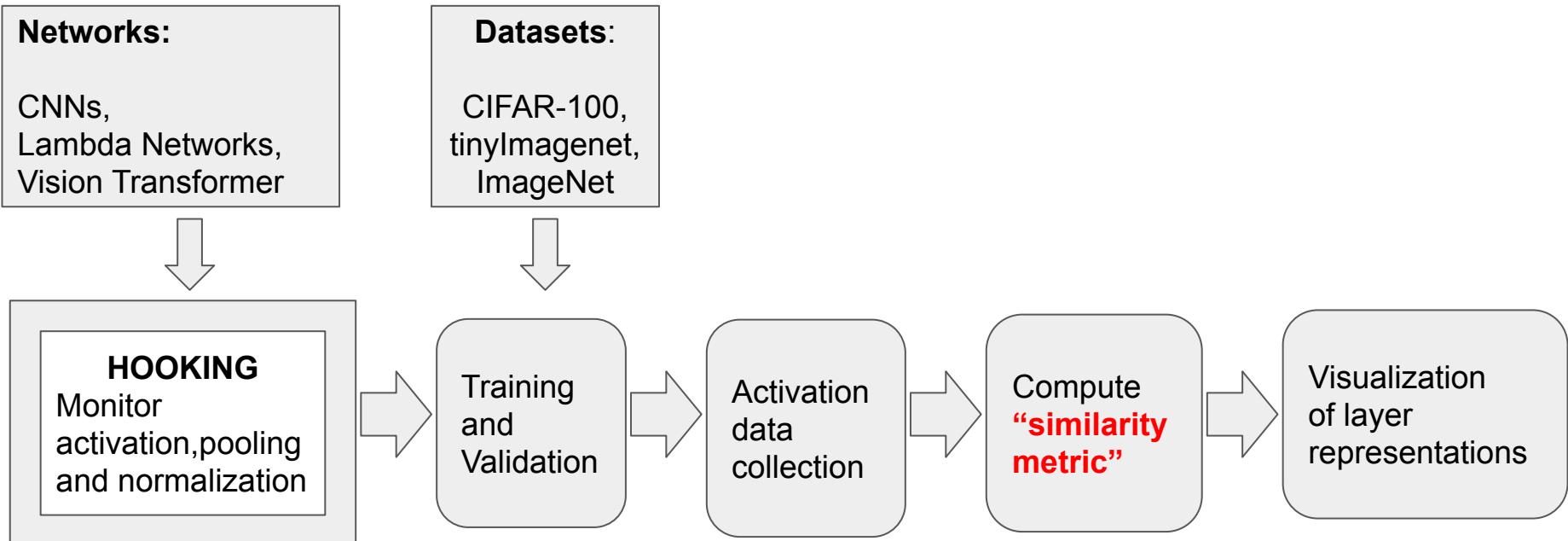
**Challenge1:** What is the right way to investigate the impact of conv vs self-attention operator?

**Hint:** A single metric (e.g., Top-1 val accuracy) *is not enough* to describe the impact of different architecture and operators on the learning of networks.

**Challenge2:** How to visualize the fundamental change on the learned representation of a **over-parameterized** network with **high-dimensional** data manifold?

**Hint:** A distance-metric/function that can capture the similarity-dissimilarity in **non-euclidean** space/geometry.

# Representation similarity in vision models: A visualization approach



# Similarity metric

What does it **measure**?

Similarity between two learned representation in high-dimensional space.

**Hilbert-Schmidt Independence Criterion(HSIC)**

$$\langle \text{vec}(XX^T), \text{vec}(YY^T) \rangle$$

$$= \text{tr}(XX^TYY^T) = \\ \|\text{cov}(X^T, Y^T)\|_F^2$$

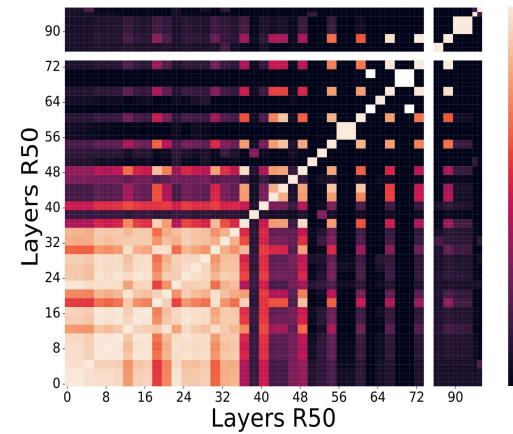
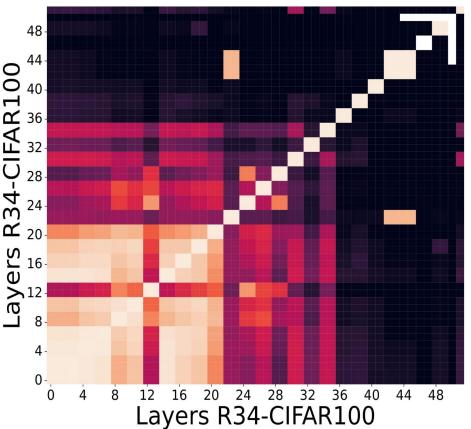
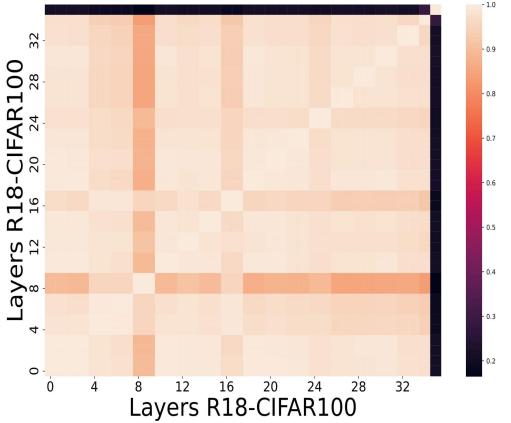
Similarity Index	Formula	Invertible Linear Transform	Orthogonal Transform	Isotropic Scaling
Linear Reg. ( $R_{LR}^2$ )	$\ Q_Y^T X\ _F^2 / \ X\ _F^2$	$Y$ only	✓	✓
CCA ( $R_{CCA}^2$ )	$\ Q_Y^T Q_X\ _F^2 / p_1$	✓	✓	✓
CCA ( $\bar{\rho}_{CCA}$ )	$\ Q_Y^T Q_X\ _* / p_1$	✓	✓	✓
SVCCA ( $R_{SVCCA}^2$ )	$\ (U_Y T_Y)^T U_X T_X\ _F^2 / \min(\ T_X\ _F^2, \ T_Y\ _F^2)$	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 h_i, 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	$\text{tr}(KHLH) / \sqrt{\text{tr}(KHKH)\text{tr}(LHLH)}$	✗	✓	✓*

**Centered Kernel Alignment (CKA)**

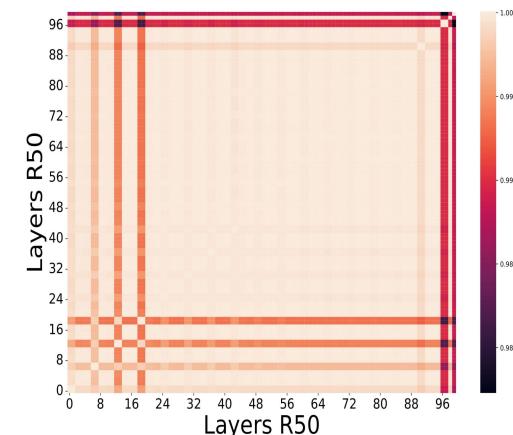
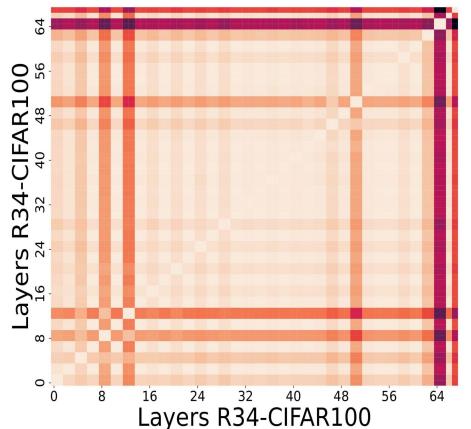
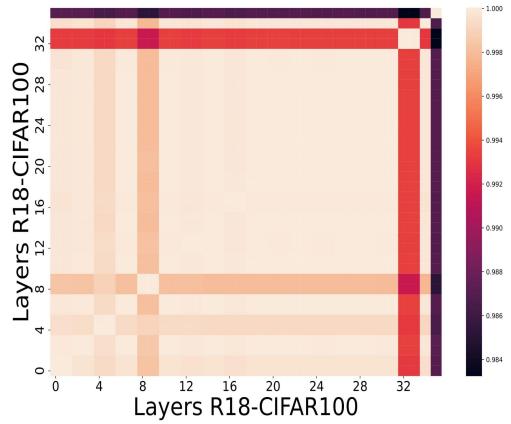
$$\text{CKA}(X, Y) = \text{HSIC}(X, Y) / \{\text{HSIC}(X, X)^* \text{HSIC}(Y, Y)\}^{1/2}$$

# Visualizing CNNs

Random weights

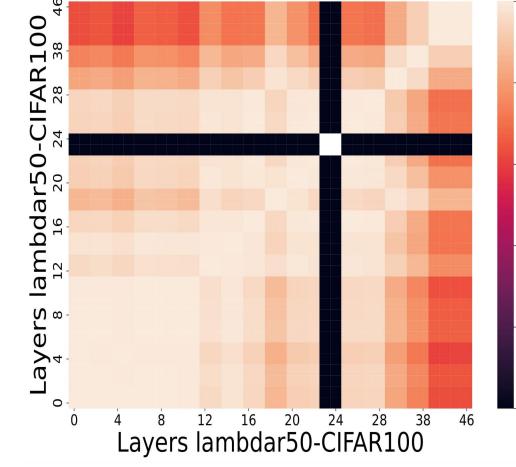
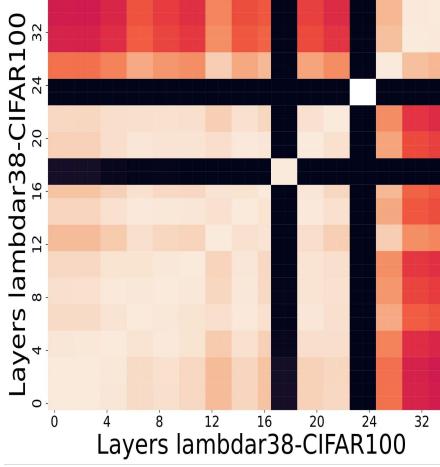
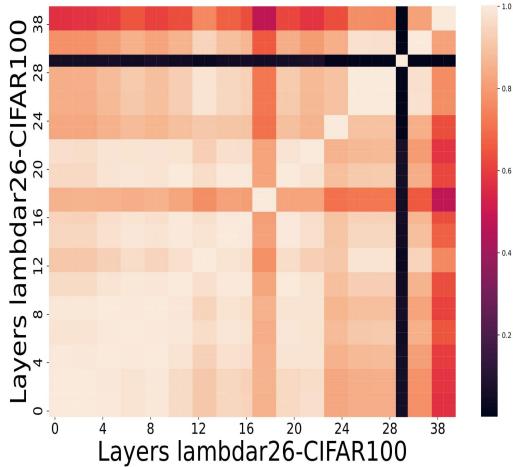


Pre-trained weights

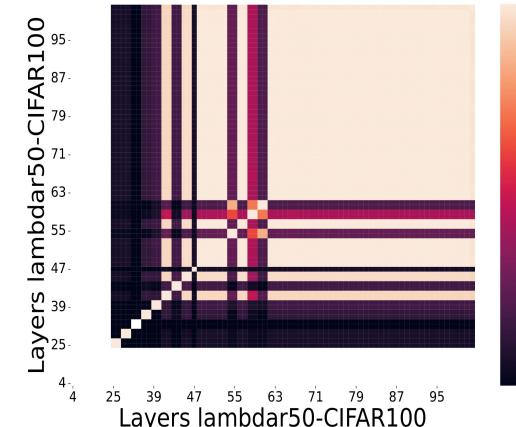
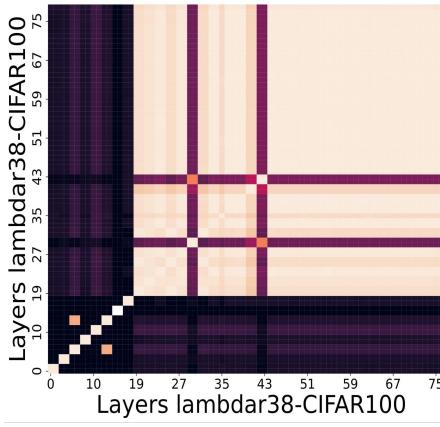
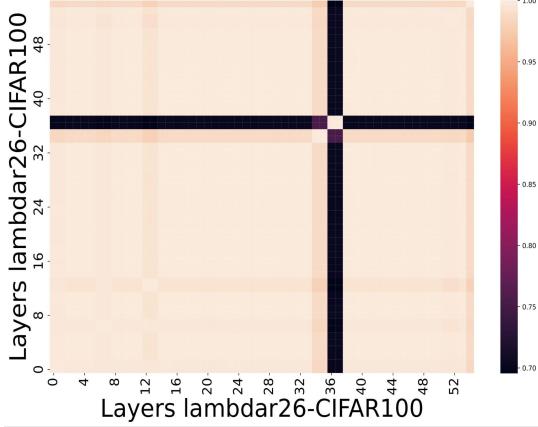


# Visualizing Lambda Networks

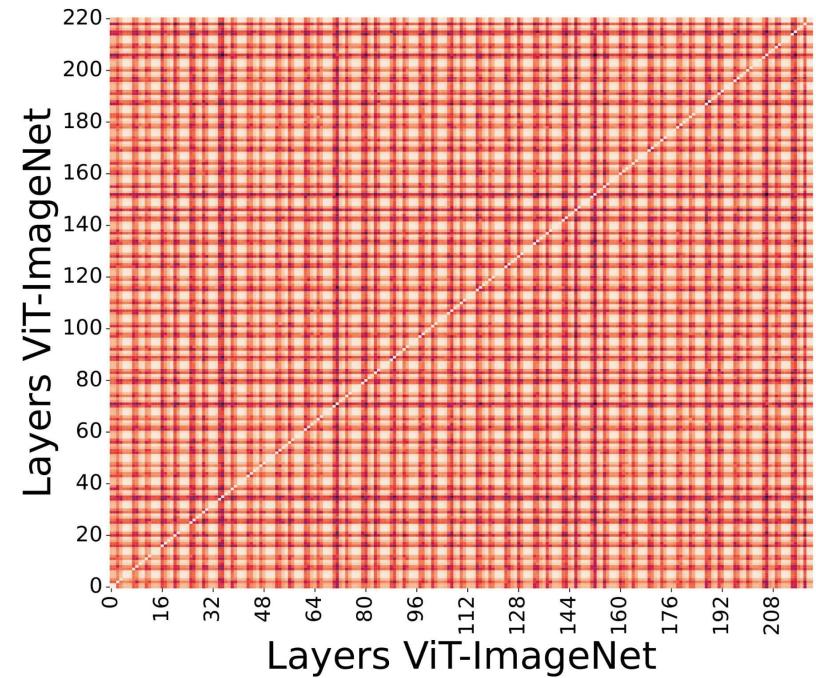
Random weights



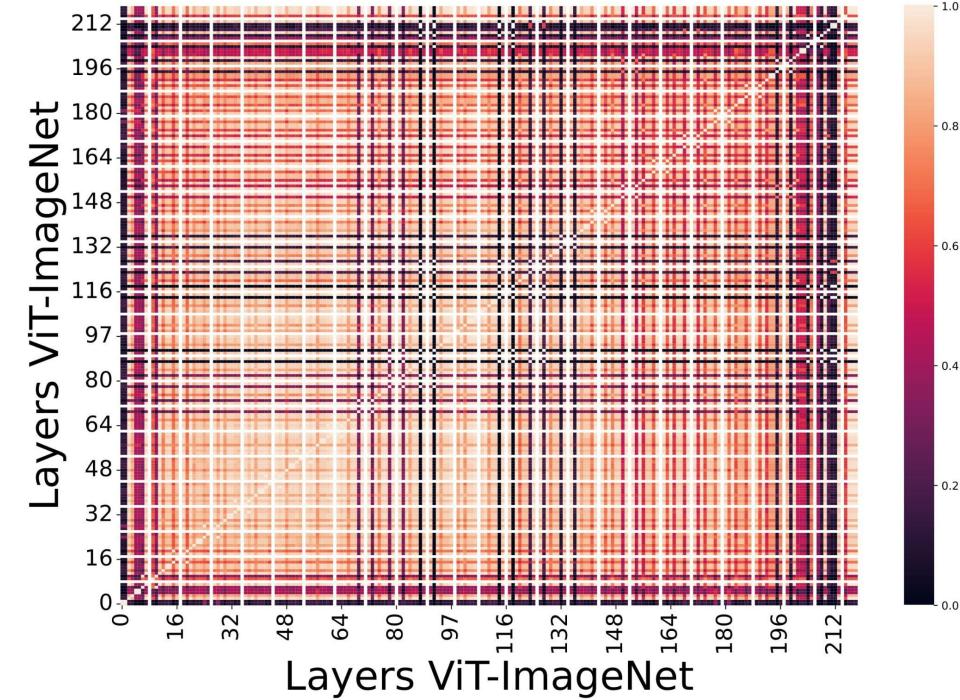
Pre-trained weights



# Visualizing Vision Transformers

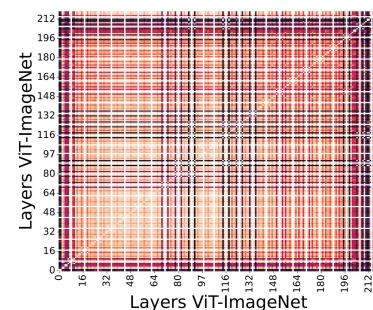
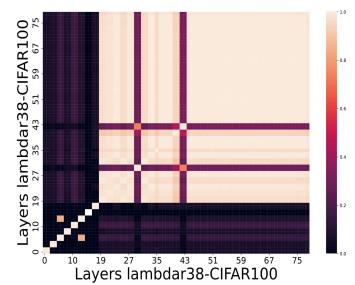
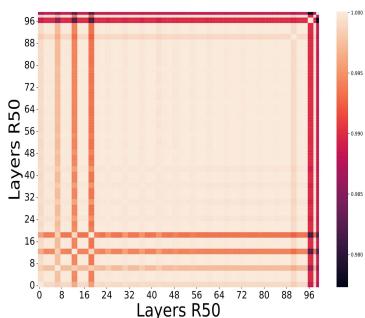
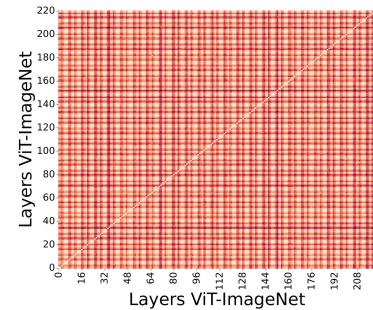
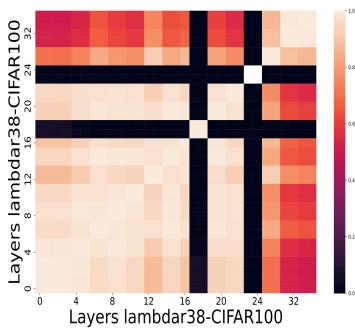
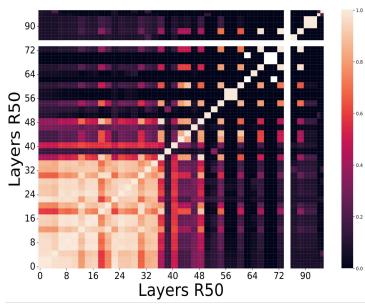


Random weights



Pre-trained weights

# Analyzing representation from CNNs to ViT



**CNNs:**  
High similarity scores  
between layers

**Lambda Networks:**  
High similarity scores  
in layers at later stages

**ViTs:**  
Symmetrical Structures. Certain  
layers are highly uncorrelated

# Conclusion

- 1) Visualization using **CKA** similarity metric captures “**crucial**” information about learned representations of hidden layers.
- 2) Our visualization has revealed an important **insight**: CNNs fails to effectively use its learning capacity; whereas self-attention has led to low CKA scores thus learning different representation across layers.
- 3) Representation learned at deeper layers contribute **more** to the predictive performance.
- 4) **Selective** use of self-attention in deeper layers of lambda networks maximizes the representation capacity of deeper layers.

**Github repo:** <https://github.com/animeshbchowdhury/visualizing-vision-models>