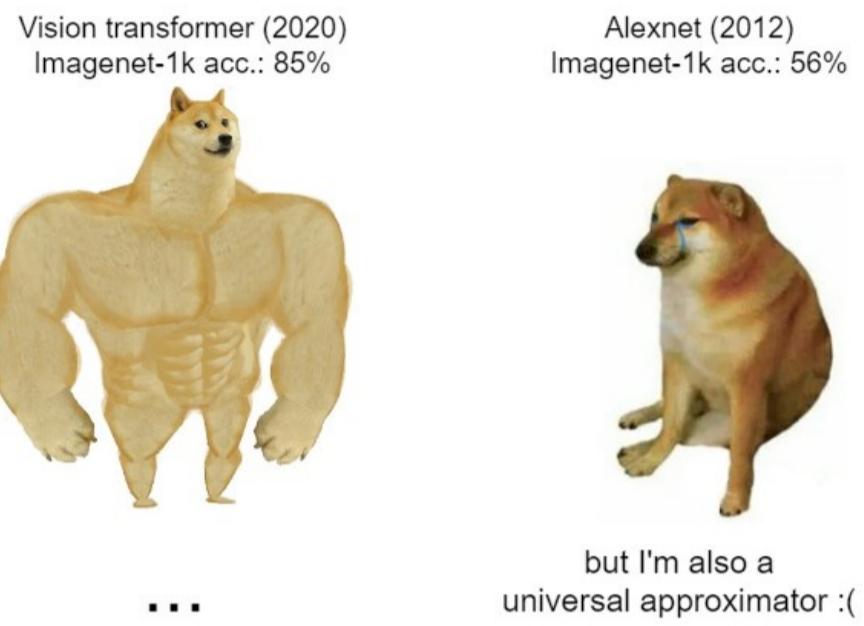


Motivation



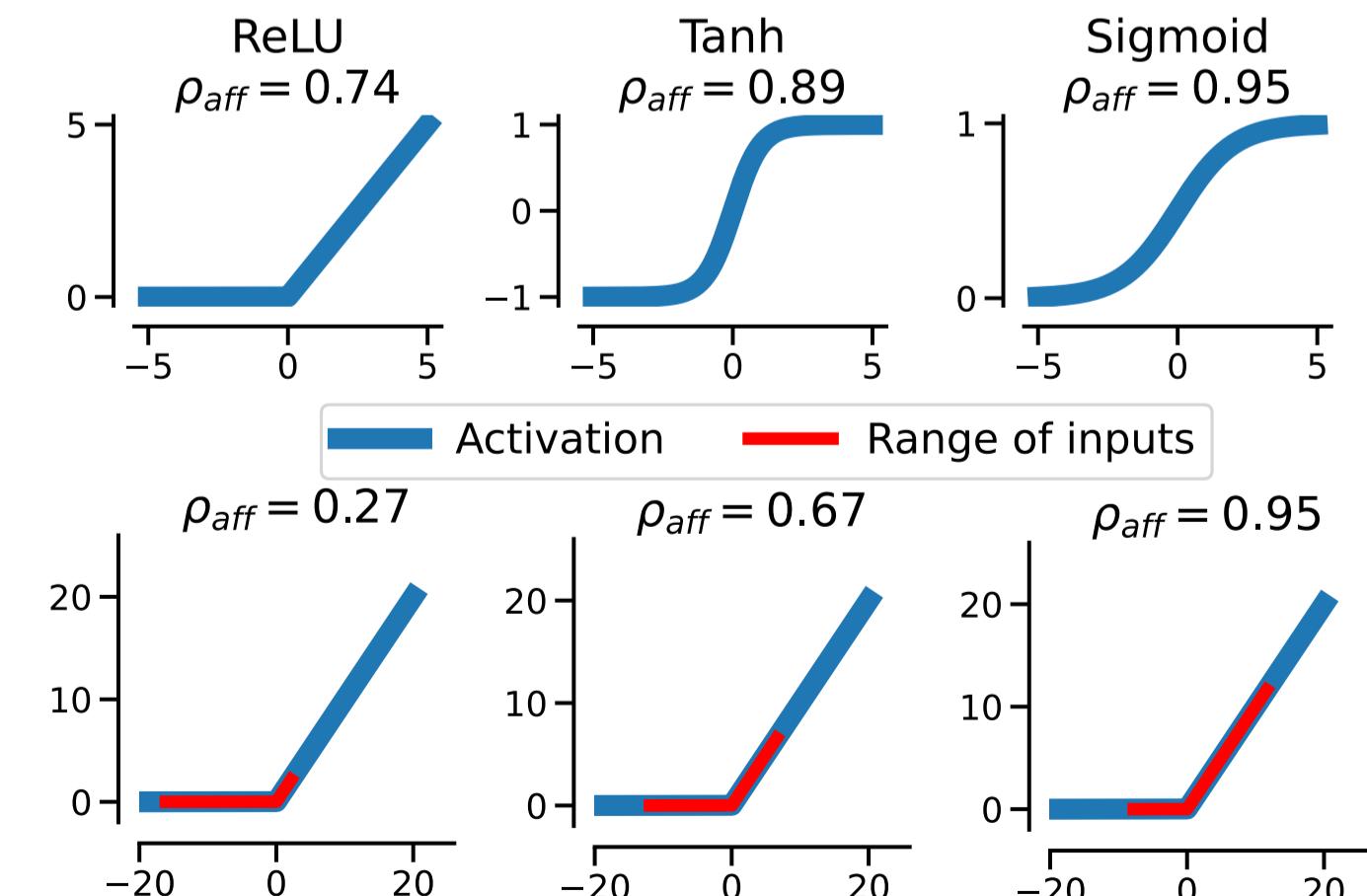
- ▶ How to compare neural architectures proposed over the years?

Our idea: better understand the intrinsic capacity of DNNs by measuring their non-linearity

Take-home message

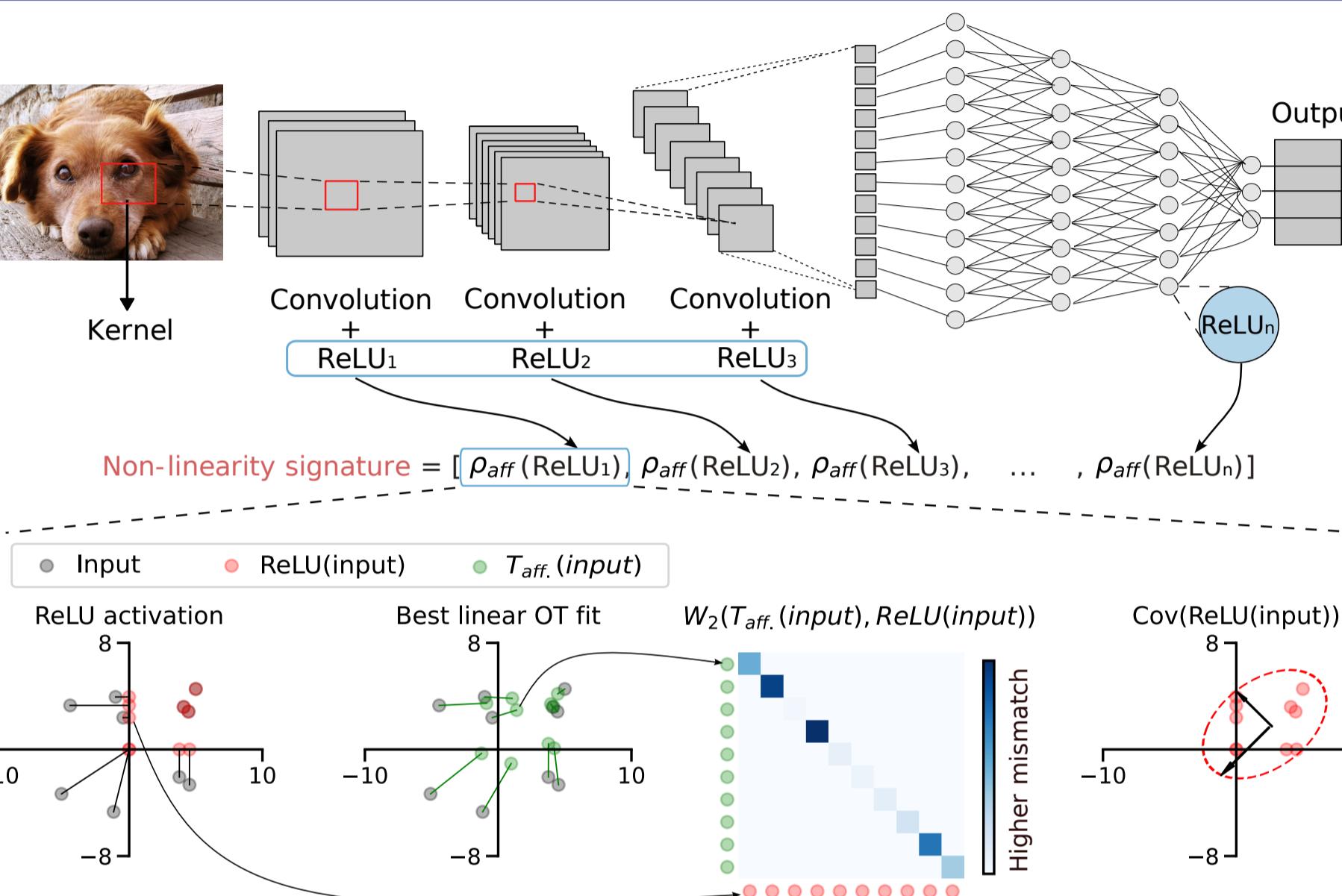
- ✓ Theoretically sound non-linearity measure for activations
- ✓ Landmark DNNs have their distinct non-linearity signature
- ✓ Potential applications: adversarial robustness, detection of novel disruptive models

What makes an activation more non-linear?



- ▶ Shape of the activation function affects its non-linearity
- ▶ Domain of pre-activations affects the non-linearity as well

Non-linearity signature of a DNN



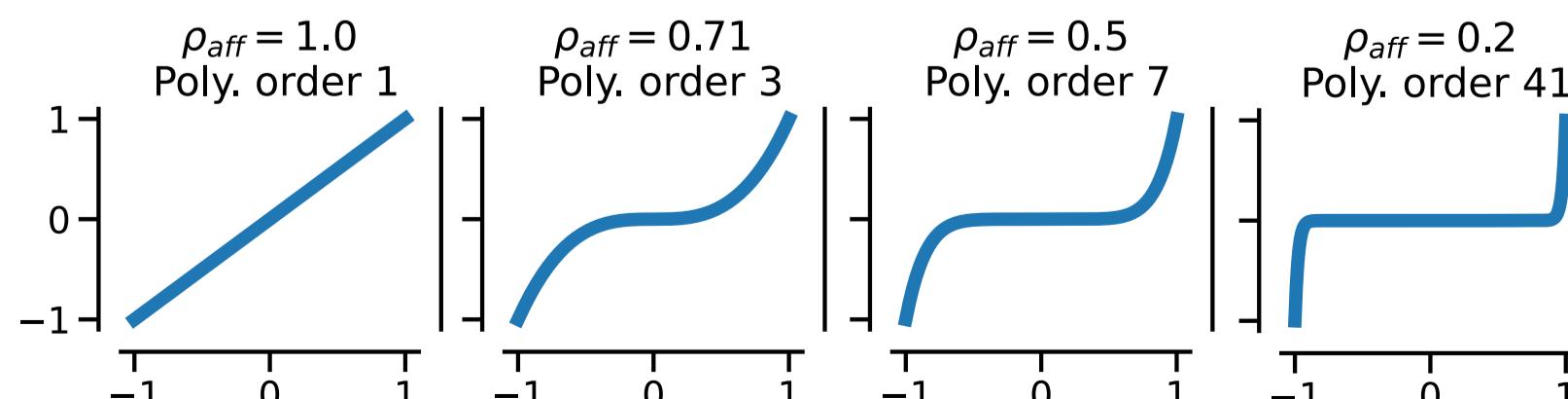
Affinity score: principle tool for measuring non-linearity

▶ Let \mathbf{X} be pre-activations within a DNN ▶ Let $\mathbf{Y} = f(\mathbf{X})$ be output of an activation function f

$$\rho_{\text{aff}}(\mathbf{X}, \mathbf{Y}) = 1 - \frac{W_2(T_{\text{aff}}(\mathbf{X}), \mathbf{Y})}{\sqrt{2 \text{Tr}[\Sigma(\mathbf{Y})]}}$$

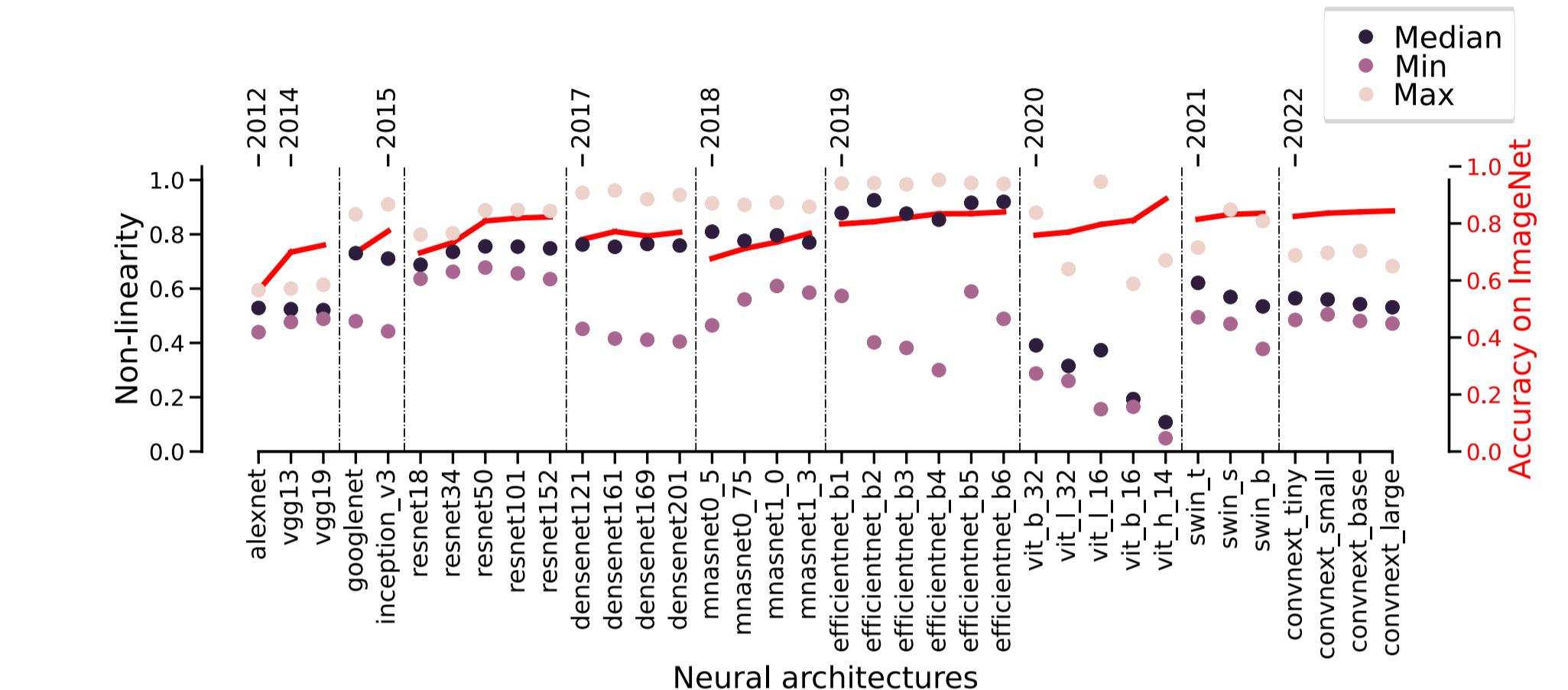
Wasserstein distance Best linear fit
Covariance of post-activations

- ▶ ρ_{aff} = how much \mathbf{Y} differs from being a PSD affine transformation of \mathbf{X}
- ▶ T_{aff} is a globally optimal linear fit, unlike least-squares solution



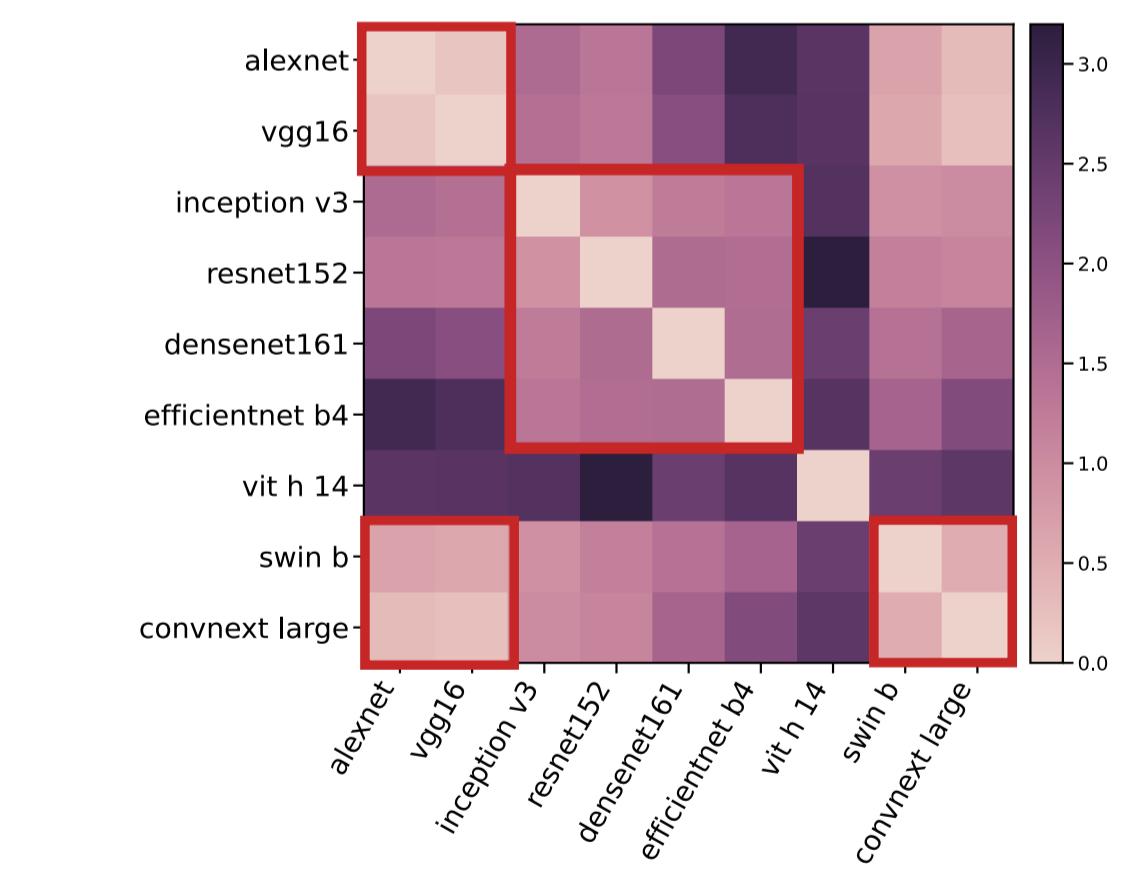
▶ Lower affinity score = transformation is more non-linear

Walking through DNNs' history

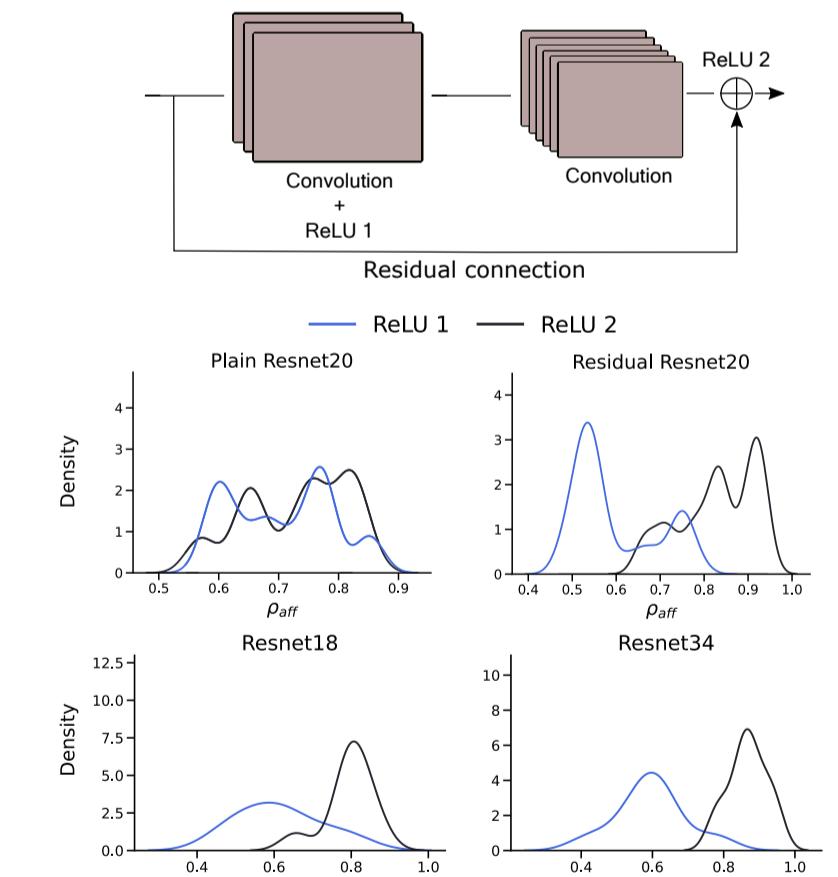


Neural architectures

Pairwise distances between DNNs



Impact of residual connections



- ▶ Before ViTs = more linear activations, more spread of ρ_{aff} .
- ▶ ViTs and after = higher non-linearity for better performance
- ▶ Residual connections = linearization of post-residual activations

