

# ResMLP

## From Hardcoded Feature to non-hardcoded

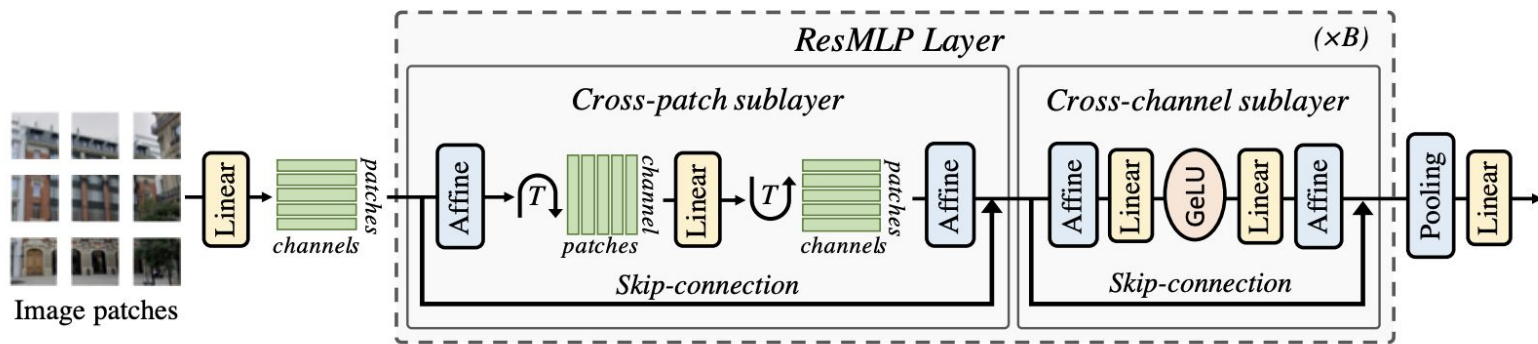
Input Image



Visualization of Learned Gradients



# Architecture



## Mathematical Representation of Architecture

$$\text{Aff}_{\alpha, \beta}(\mathbf{x}) = \text{Diag}(\boldsymbol{\alpha})\mathbf{x} + \beta,$$

$$\mathbf{Z} = \mathbf{X} + \text{Aff} \left( (\mathbf{A} \text{Aff}(\mathbf{X})^\top)^\top \right),$$

$$\mathbf{Y} = \mathbf{Z} + \text{Aff}(\mathbf{C} \text{GELU}(\mathbf{B} \text{Aff}(\mathbf{Z}))),$$

## ResMlp Unique Features

- No batch norm or Layer norm
- Use Linear layer instead of Attention layer Hence claim to have more stability.
- No positional encoding to encode patches position

# Findings During Implementation

Dataset: Cifar 20   Model: ResMlp-s12   patches=16   Embed dim=768

- Direct paper implementation output loss: Nan
- Overfitting at some point

Thanks