

CycleMLP

A MLP-like Architecture for Dense Prediction

Marco Benelli

University of Florence

February 14, 2022

Outline

- 1 Introduction
- 2 Method
 - Cycle Fully-Connected Layer
 - Overall Architecture
- 3 Classification Experiments
 - CIFAR10
 - STL10
 - ImageNet-1K

Paradigm Shifts

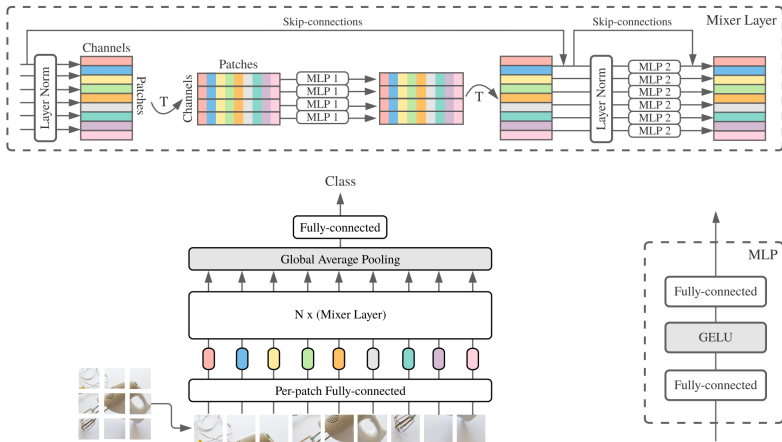
Recent paradigm shifts:

2012 AlexNet

2020 ViT

2021 MLP-Mixer

MLP-Mixer



Mixer Layer

$$\mathbf{U}_{*,i} = \mathbf{X}_{*,i} + \mathbf{W}_2 \sigma(\mathbf{W}_1 \text{LayerNorm}(\mathbf{X})_{*,i}), \quad \text{for } i = 1 \dots C$$

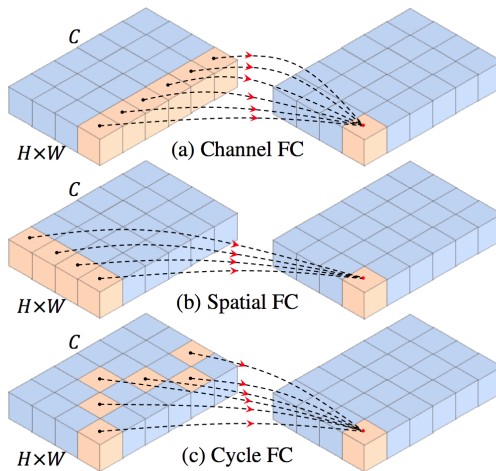
$$\mathbf{Y}_{j,*} = \mathbf{U}_{j,*} + \mathbf{W}_4 \sigma(\mathbf{W}_3 \text{LayerNorm}(\mathbf{U})_{j,*}), \quad \text{for } j = 1 \dots S$$

Challenges

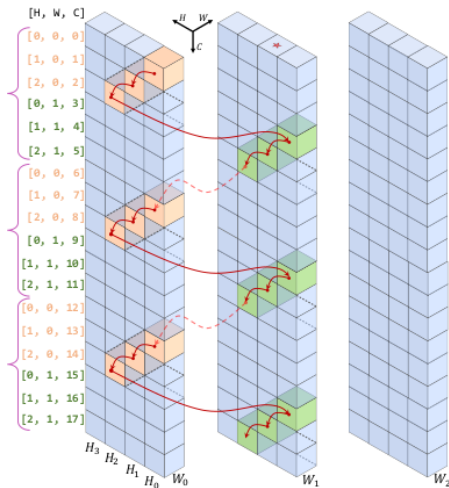
MLP-like models are facing these challenges:

- non-hierarchical architectures
- flexible input scales
- quadratic costs

Cycle Fully-Connected Layer



Stepsize Example



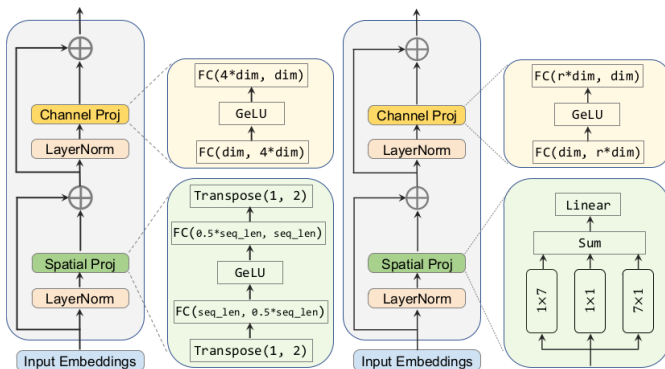
Cycle Fully-Connected Layer

$$\text{CycleFC}(\mathbf{X})_{i,j,:} = \sum_{c=0}^{C_{\text{in}}} \mathbf{x}_{i+\delta_i(c), j+\delta_j(c), c} \cdot \mathbf{W}_{c,:} + \mathbf{b}$$

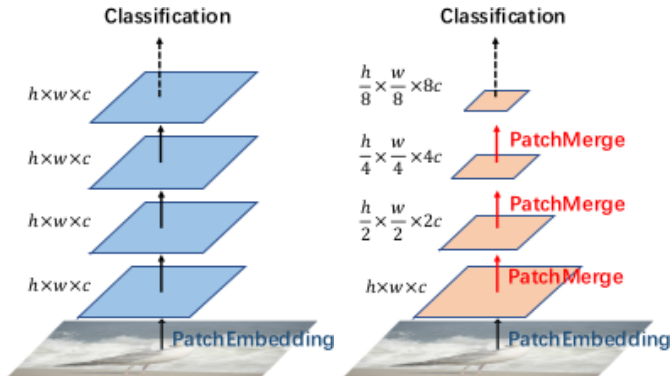
$$\delta_i(c) = (c \bmod S_H) - \left\lfloor \frac{S_H}{2} \right\rfloor$$

$$\delta_j(c) = \left(\left\lfloor \frac{c}{S_H} \right\rfloor \bmod S_W \right) - \left\lfloor \frac{S_W}{2} \right\rfloor$$

Comparison Of MLP Blocks



Hierarchy



Instantiation

	Output Size	Layer Name	B1
Stage 1	$\frac{H}{4} \times \frac{W}{4}$	Overlapping Patch Embedding	$C_1 = 64$
		CycleMLP Block	$E_1 = 4$ $L_1 = 2$
Stage 2	$\frac{H}{8} \times \frac{W}{8}$	Overlapping Patch Embedding	$C_2 = 128$
		CycleMLP Block	$E_2 = 4$ $L_2 = 2$
Stage 3	$\frac{H}{16} \times \frac{W}{16}$	Overlapping Patch Embedding	$C_3 = 320$
		CycleMLP Block	$E_3 = 4$ $L_3 = 4$
Stage 4	$\frac{H}{32} \times \frac{W}{32}$	Overlapping Patch Embedding	$C_4 = 512$
		CycleMLP Block	$E_4 = 4$ $L_4 = 2$

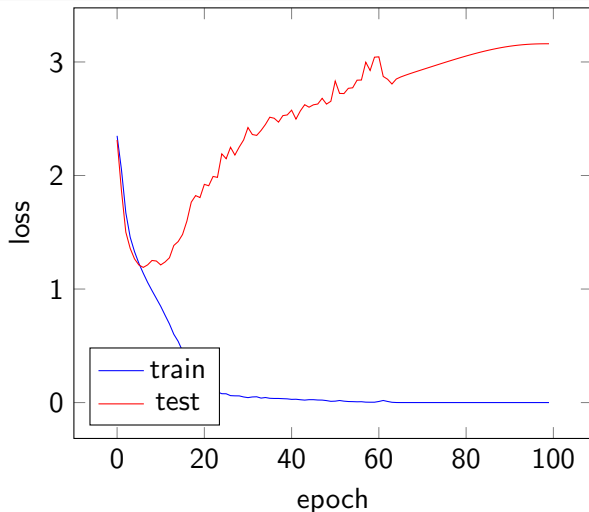
Experimental Setup

- optimizer AdamW
- $\lambda = 5 \times 10^{-2}$
- cosine annealing learning rate schedule
- $\eta_{\max} = 1 \times 10^{-3}$
- $T_{\max} = 100$
- batch size = 256

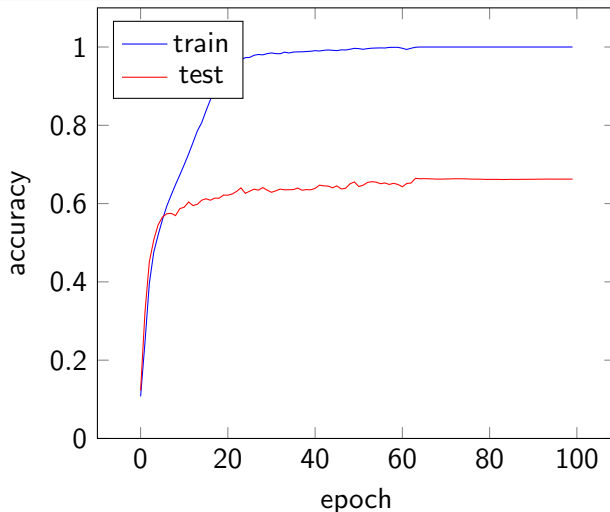
Experiments

Model	STL10	CIFAR10
ResNet	64.9%	77.1%
ViT	44.4%	53.4%
MLP-Mixer	51.4%	55.5%
CycleMLP	49.8%	66.5%

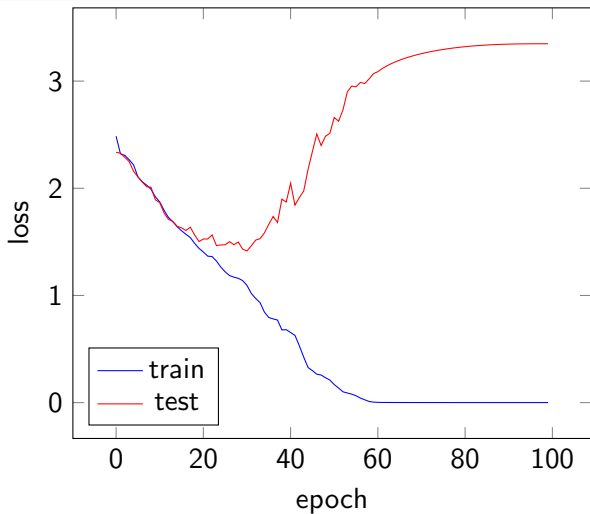
Loss Plot (CIFAR10)



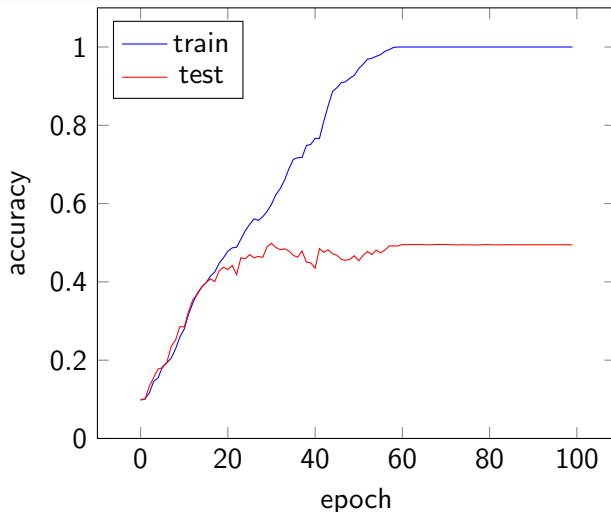
Accuracy Plot (CIFAR10)



Loss Plot (STL10)



Accuracy Plot (STL10)



ImageNet-1K Comparison

Model	Accuracy
ResNet	69.8%
ViT	77.9%
MLP-Mixer	61.4%
CycleMLP	79.1%

Summary

- CycleMLP is built upon the **Cycle FC**.
- Cycle FC is capable of dealing with **variable input scales**.
- The computational cost of Cycle FC is **$O(HWC^2)$** .