



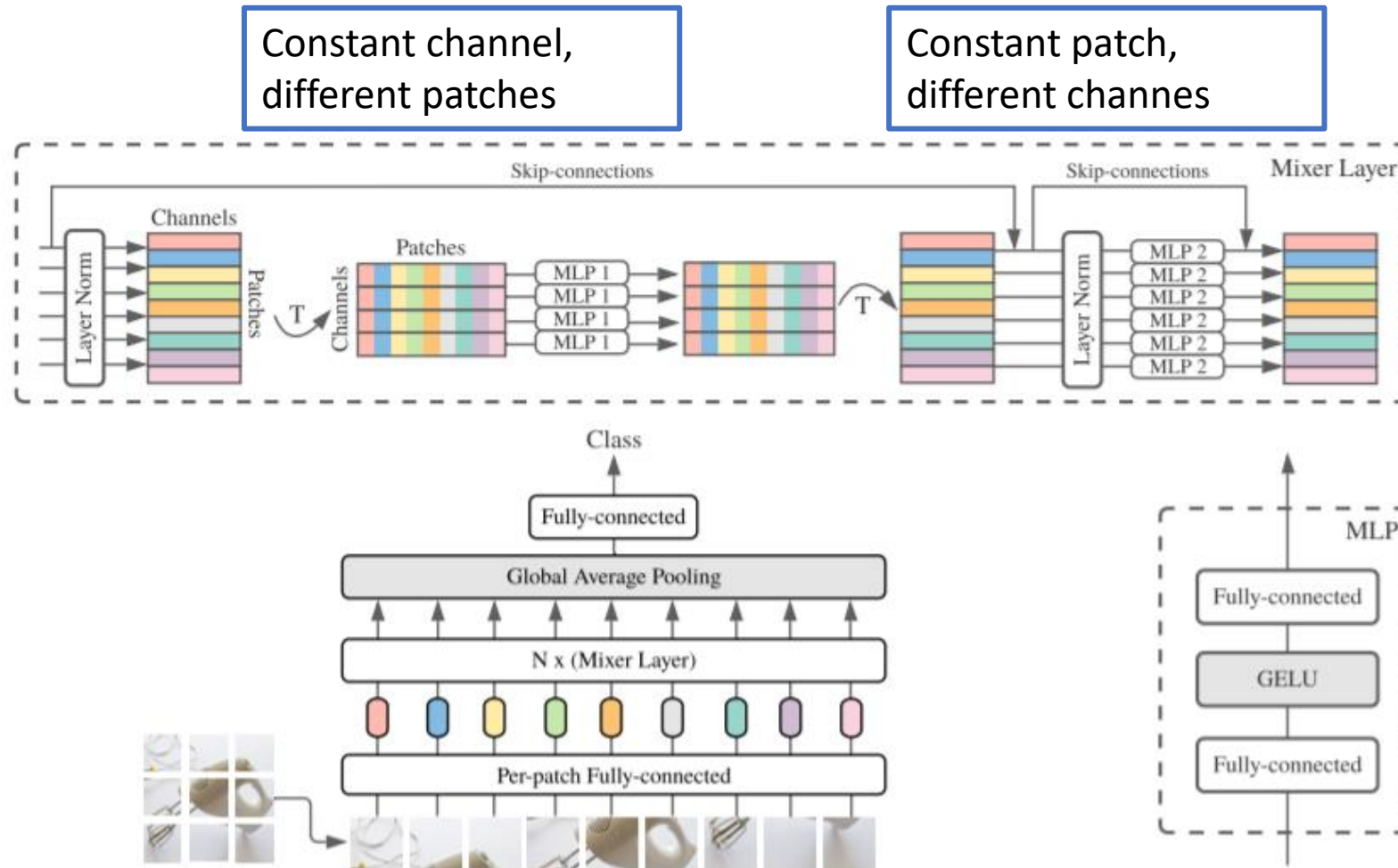
# NeurIPS2021 – key takeaways\*

Paulina Tomaszewska

**\*regarding Computer Vision and Knowledge Transfer**

# MLP-Mixer: An all-MLP Architecture for Vision (Tolstikhin et al.)

\*from May 21' → 164 citations  
(Semantic Scholar)



- In classical MLP, the image would be represented as a vector
- LN: normalize the activations of the previous layer for each given example in a batch independently (mean = 0, std = 1)

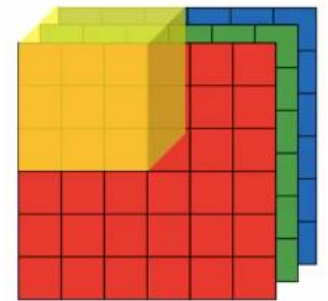


Figure 1: MLP-Mixer consists of per-patch linear embeddings, Mixer layers, and a classifier head. Mixer layers contain one token-mixing MLP and one channel-mixing MLP, each consisting of two fully-connected layers and a GELU nonlinearity. Other components include: skip-connections, dropout, and layer norm on the channels.

# GeLU = Gaussian Error Linear Unit

$$\text{GELU}(x) = xP(X \leq x) = x\Phi(x)$$

if  $X \sim \mathcal{N}(0, 1)$ .

- Uses standard Gaussian cumulative distribution function
- Smoother than ReLU

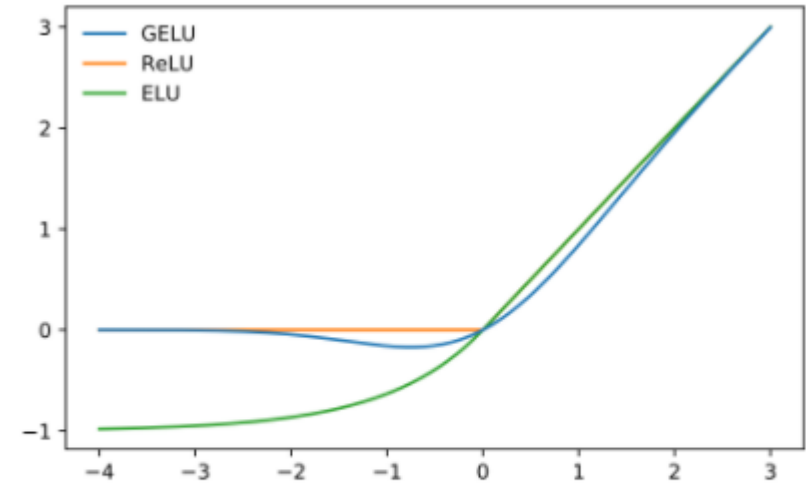


Figure 1: The GELU ( $\mu = 0, \sigma = 1$ ), ReLU, and ELU ( $\alpha = 1$ ).

*Another „degree of freedom“ when searching for the optimal architecture?*

# Invariance to input permutations

- Differences in inductive bias of Mixer and CNN architectures
- Same permutation is used across all images

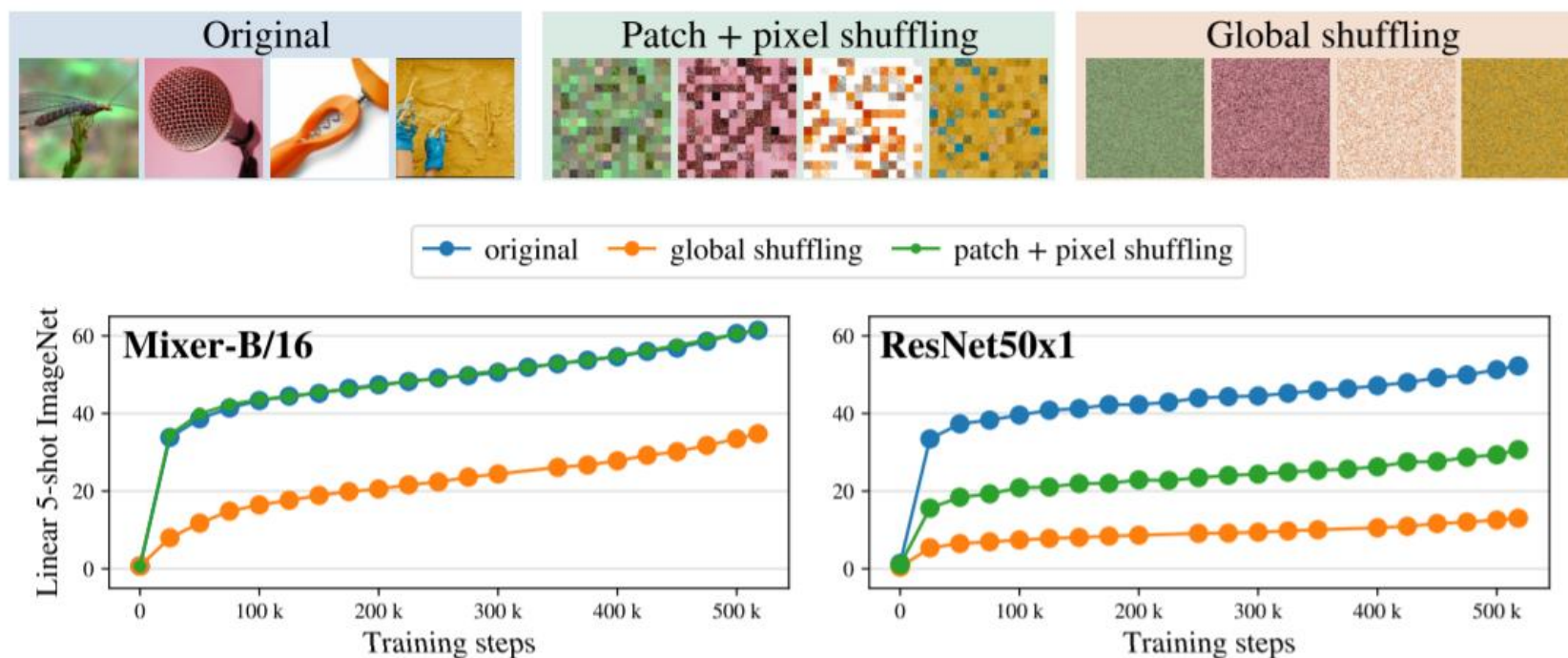
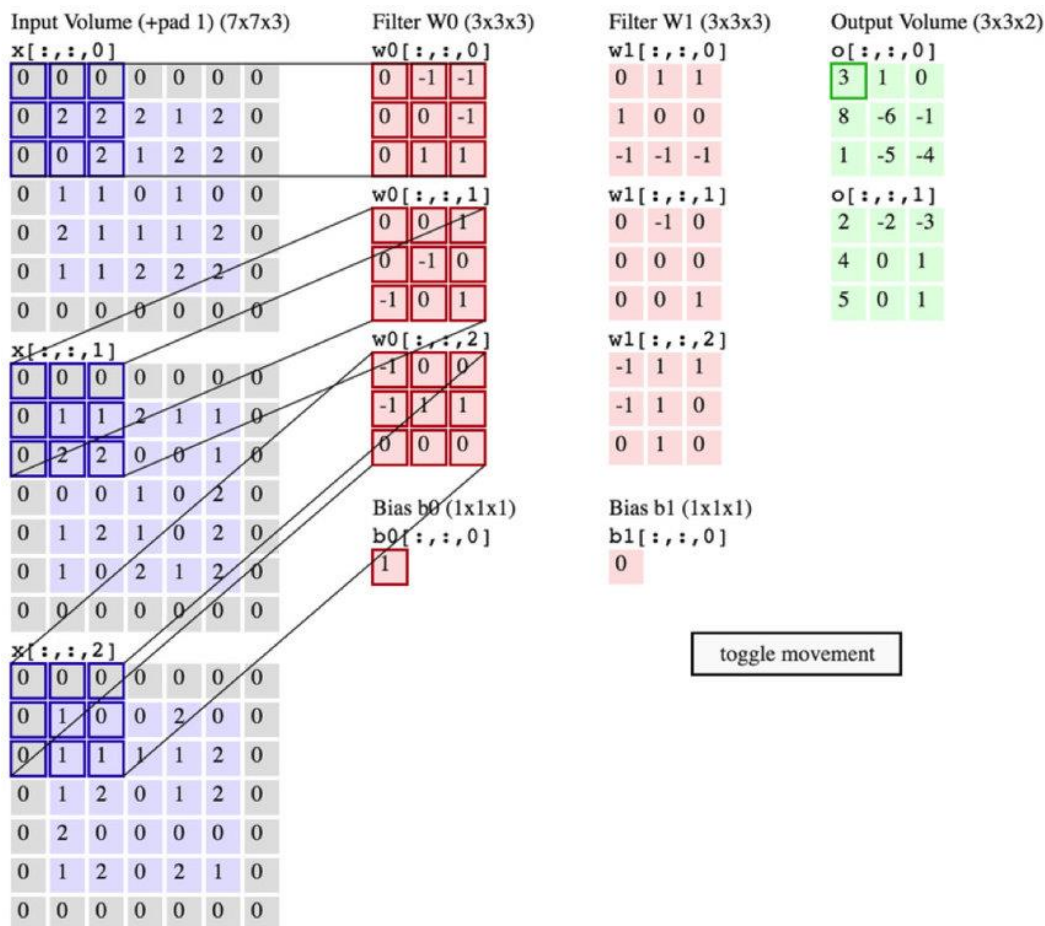


Figure 4: **Top:** Input examples from ImageNet before permuting the contents (left); after shuffling the  $16 \times 16$  patches and pixels within the patches (center); after shuffling pixels globally (right). **Bottom:** Mixer-B/16 (left) and ResNet50x1 (right) trained with three corresponding input pipelines.



# No Conv?



## E MLP-Mixer code

```

1 import einops
2 import flax.linen as nn
3 import jax.numpy as jnp
4
5 class MlpBlock(nn.Module):
6     mlp_dim: int
7     @nn.compact
8     def __call__(self, x):
9         y = nn.Dense(self.mlp_dim)(x)
10        y = nn.gelu(y)
11        return nn.Dense(x.shape[-1])(y)
12
13 class MixerBlock(nn.Module):
14     tokens_mlp_dim: int
15     channels_mlp_dim: int
16     @nn.compact
17     def __call__(self, x):
18         y = nn.LayerNorm()(x)
19         y = jnp.swapaxes(y, 1, 2)
20         y = MlpBlock(self.tokens_mlp_dim, name='token_mixing')(y)
21         y = jnp.swapaxes(y, 1, 2)
22         x = x+y
23         y = nn.LayerNorm()(x)
24         return x+MlpBlock(self.channels_mlp_dim, name='channel_mixing')(y)
25
26 class MlpMixer(nn.Module):
27     num_classes: int
28     num_blocks: int
29     patch_size: int
30     hidden_dim: int
31     tokens_mlp_dim: int
32     channels_mlp_dim: int
33     @nn.compact
34     def __call__(self, x):
35         s = self.patch_size
36         x = nn.Conv(self.hidden_dim, (s,s), strides=(s,s), name='stem')(x)
37         x = einops.rearrange(x, 'n h w c -> n (h w) c')
38         for _ in range(self.num_blocks):
39             x = MixerBlock(self.tokens_mlp_dim, self.channels_mlp_dim)(x)
40         x = nn.LayerNorm(name='pre_head_layer_norm')(x)
41         x = jnp.mean(x, axis=1)
42         return nn.Dense(self.num_classes, name='head',
43                        kernel_init=nn.initializers.zeros)(x)

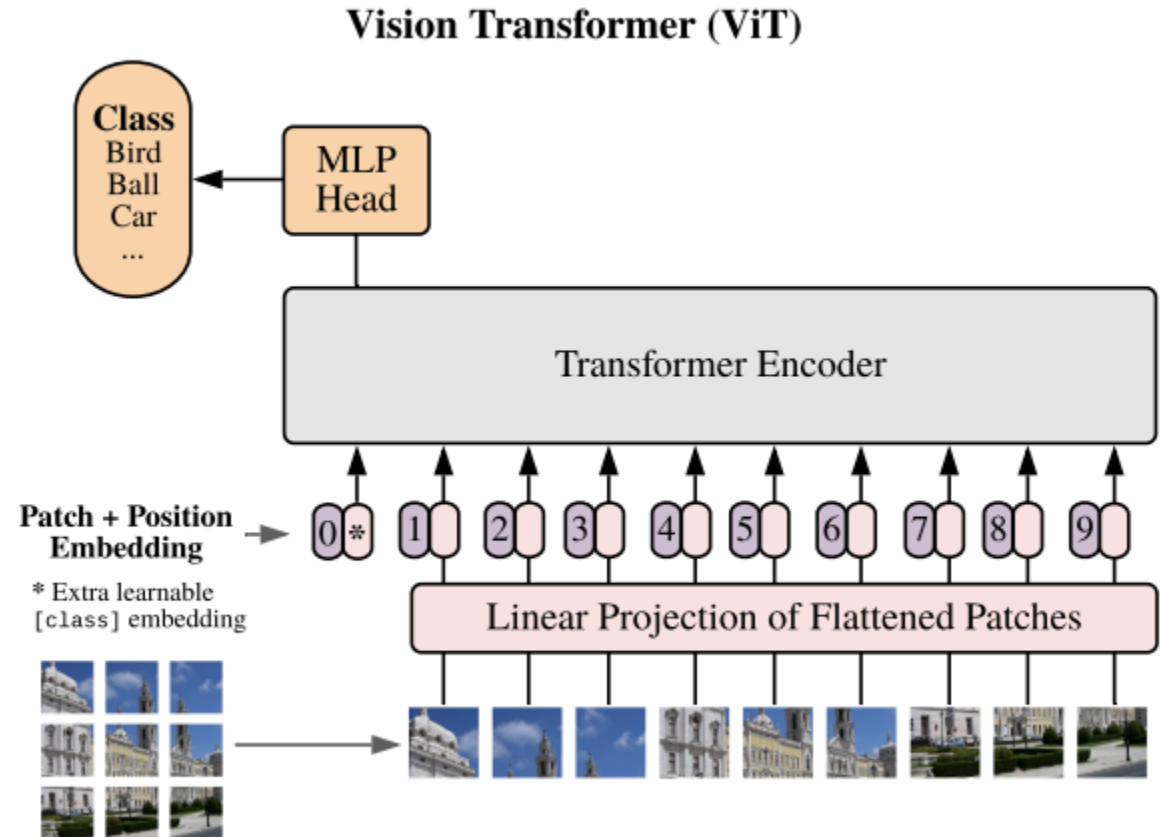
```

Listing 1: MLP-Mixer code written in JAX/Flax.



# MLP-Mixer: summary

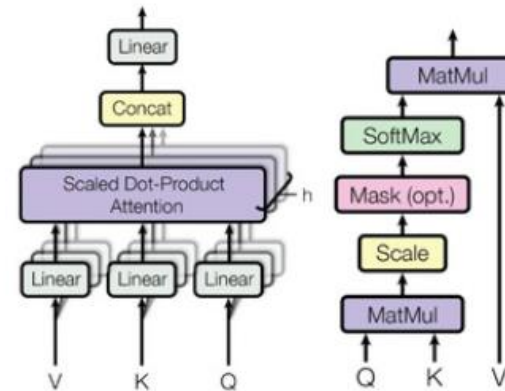
- Patches like in Vision Transformer (ViT)
- 2.5 times faster than ViT but comparable accuracy



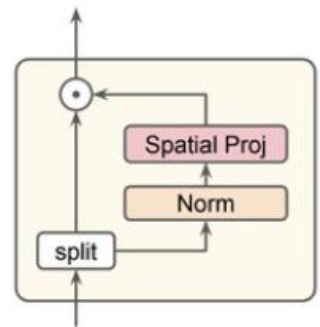
# Pay Attention to MLPs (Liu et al.)

- Gated MLPs
- On par with ViT
- Claim: „self-attention is not the key for model scalability”

## Self-Attention vs Gating



Self-attention:  $q * k * v$  (3rd order)

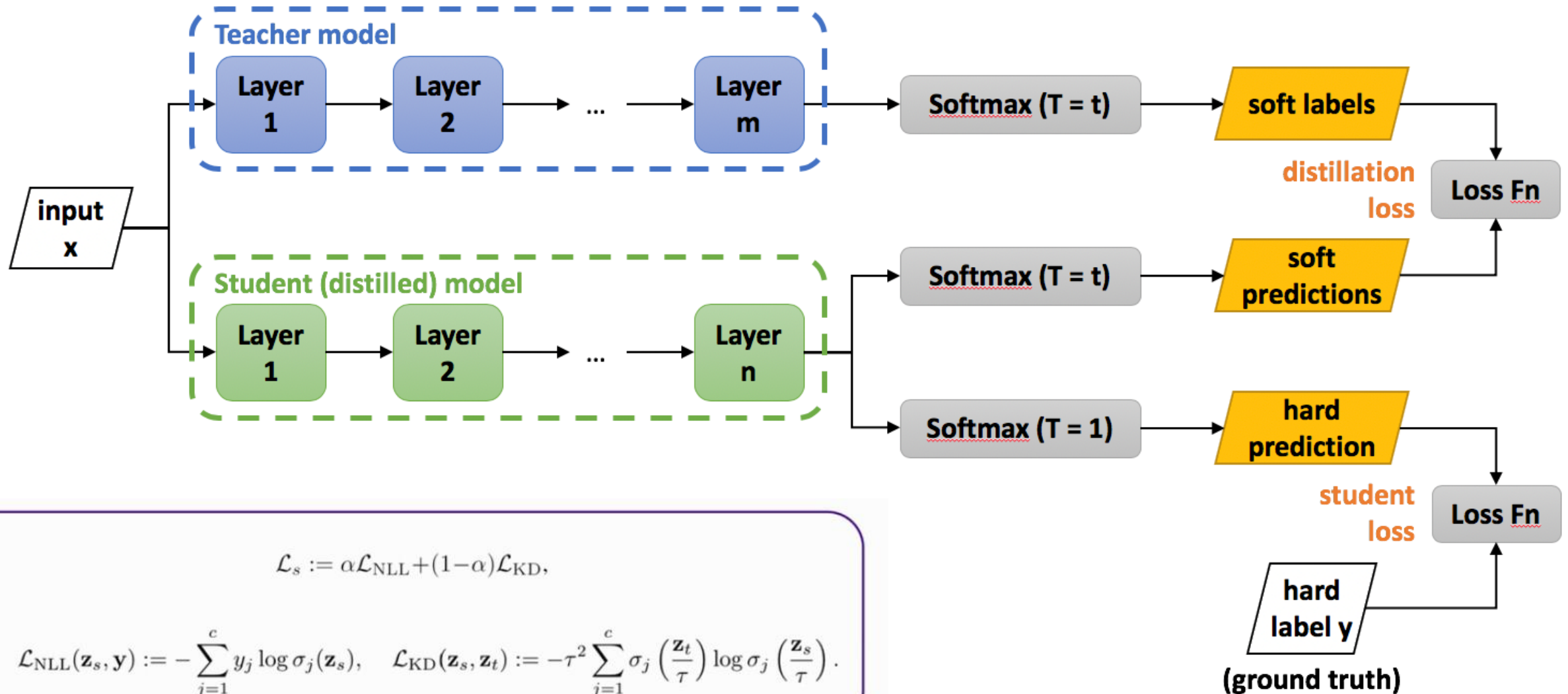


Gating:  $u * v$  (2nd order)

# **Knowledge distillation (KD)**

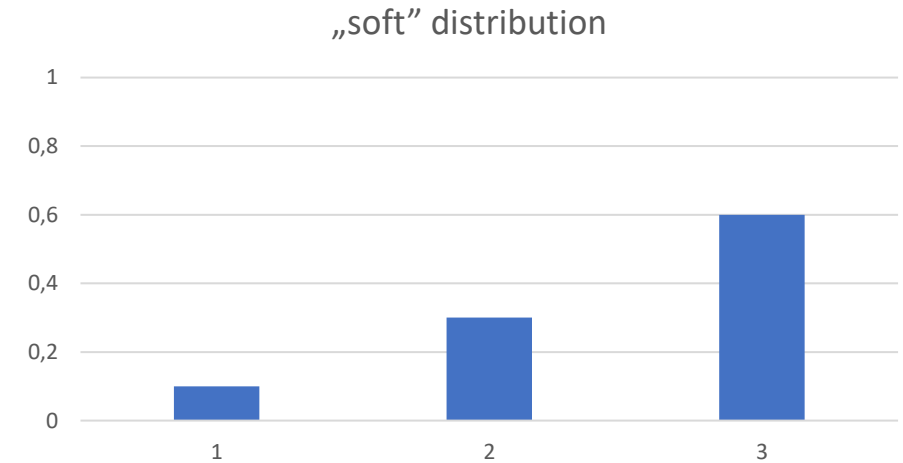
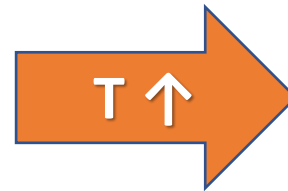
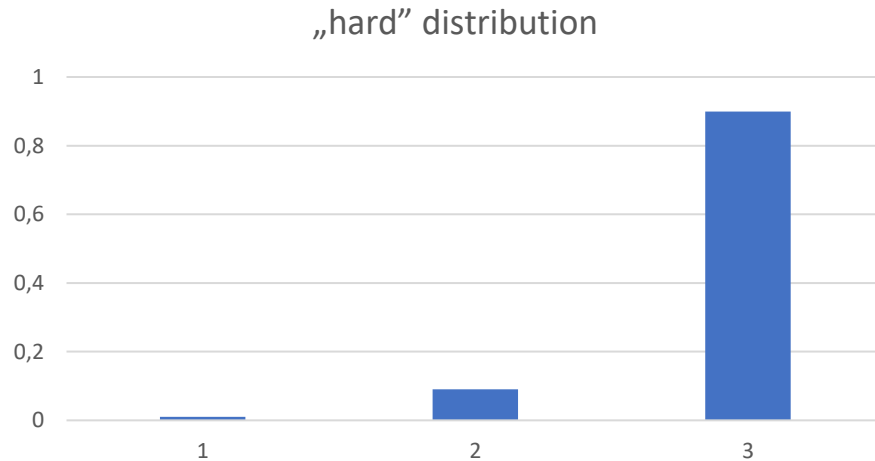


# Architecture



# Temperature scaling

$$P_i = \frac{e^{\frac{y_i}{T}}}{\sum_{k=1}^n e^{\frac{y_k}{T}}} \quad T=1 \rightarrow \text{„normal softmax“}$$



## Example - RNN:

We are sampling from output distribution and choosing the sampled word as your output token (and next input). If the model is extremely confident, it may produce very repetitive and uninteresting text. ☹️

We want it to produce more diverse text which it will not produce because when the sampling procedure is going on, most of the probability mass will be concentrated in a few tokens and thus your model will keep selecting a select number of words over and over again.

In order to give other words a chance of being sampled as well, you could plug in the temperature and produce more diverse text. 😊

# Temperature scaling in KD

- „Hard” distribution doesn't provide much information beyond the ground truth labels already provided in the dataset.
- The „softer” probability provides more information as to which classes the teacher found more similar to the predicted class.
- Previous approach: using rather than the probabilities produced by the softmax as the targets for learning the small model
  - Target: minimize the squared difference between the logits produced by the cumbersome model and the logits produced by the small model

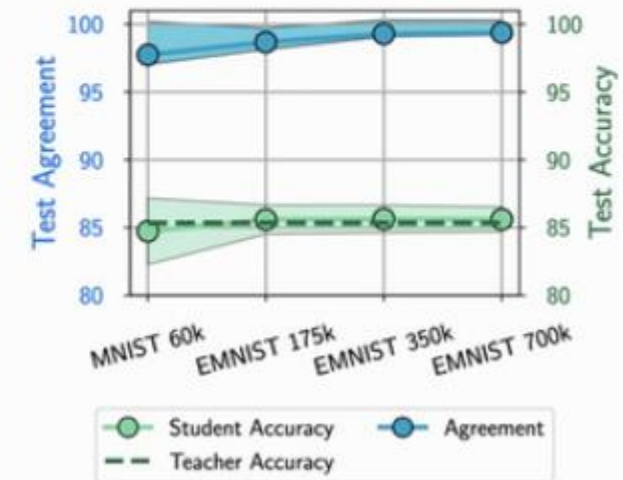
# Does Knowledge Distillation Really Work? (Stanton et al.)

## Decoupling accuracy and fidelity

Expectation: as the size of the distillation dataset increases teacher and student become functionally equivalent.

$$f_{\text{teacher}}(x) = a_0 + a_1x + \dots + a_nx^n \quad f_{\text{student}}(x) = b_0 + b_1x + \dots + b_nx^n$$
$$f_{\text{student}}(x) = f_{\text{teacher}}(x) \quad \forall x \in \{x_1, \dots, x_{n+1}\} \Rightarrow f_{\text{student}} \equiv f_{\text{teacher}}$$

LeNet-5 distillation on MNIST



Source: poster at NeurIPS

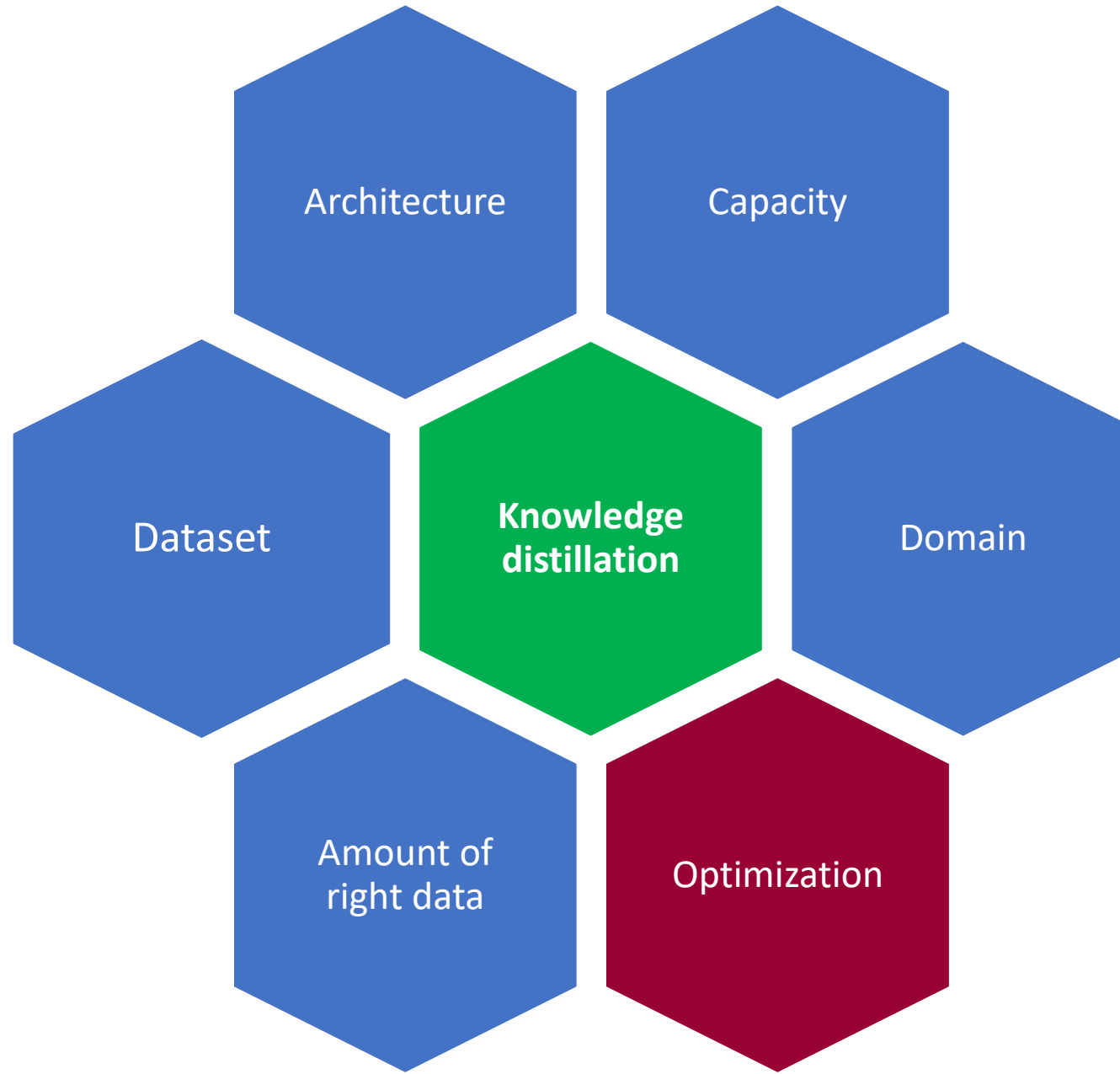


Analogy to self-distillation

# Observations

- Even in case of self-distillation, low fidelity occurs
- Increasing student depth has very little effect on fidelity

# Why does knowledge distillation result in low fidelity?



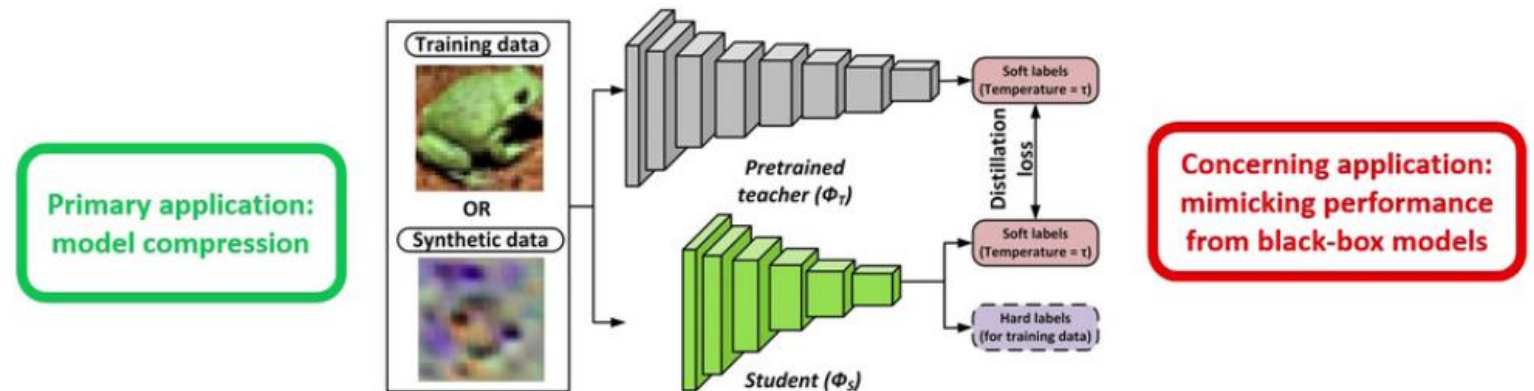


# Analyzing the Confidentiality of Undistillable Teachers in Knowledge Distillation (Kundu et al.)

- Motivation

- Machine Learning as a Service is on a rise
- Models are released as black-boxes APIs so that competitors can't replicate them

## Knowledge-Distillation (KD): A Potential Threat to MLAAS




- KD can transfer the “rich” knowledge of a compute-heavy teacher to a compute-efficient student model under both data-available<sup>[1]</sup> and data-free scenarios<sup>[2]</sup>

[1] Geoffrey Hinton et al., “Distilling the knowledge in a neural network”, NeurIPS 2014 (workshop).

[2] Paul Micaelli and Amos Storkey, “Zero-shot knowledge transfer via adversarial belief matching”, NeurIPS 2019.

# Undistillable Models<sup>[1]</sup>

- 
- A class of models that
    - Perform similar to standard teacher models to maintain their own performance
    - However, act as “**nasty**” teachers to any student model by not allowing it to mimic performance.
  - Core idea
    - Inject **false** sense of generalization to the student<sup>[1]</sup>

Training loss of Undistillable models ( $\Phi_T$ ):

$$\mathcal{L}_N = \underbrace{\mathcal{L}_{CE}(\sigma(g_{\Phi_T}(\mathbf{x}, \mathbf{y})))}_{\text{Cross-entropy (CE) loss}} - \alpha_N * \tau_N^2 * \underbrace{\mathcal{L}_{KL}(\sigma(g_{\Phi_T}(\mathbf{x}, \mathbf{y}), \tau_N), \sigma(g_{\Phi_A}(\mathbf{x}, \mathbf{y}), \tau_N))}_{\text{Self-undermining loss}}$$

Cross-entropy (CE)  
loss

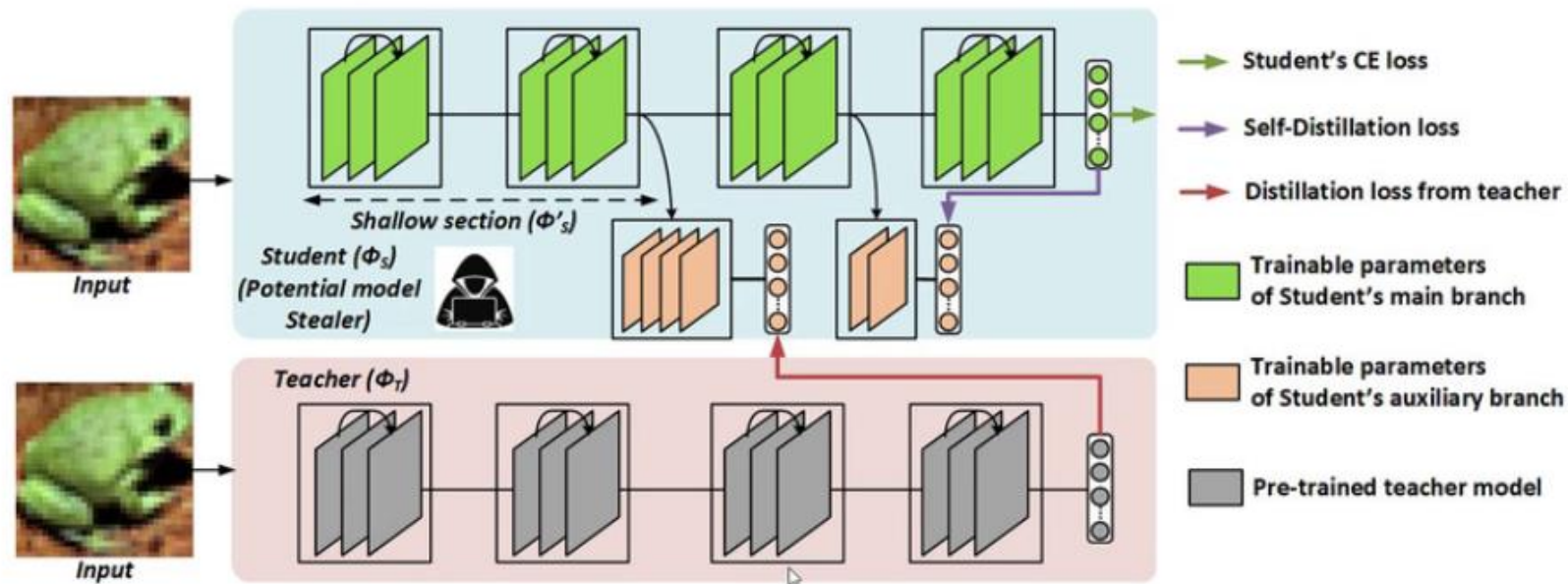
Self-undermining loss

[1] Haoyu Ma et al., “Undistillable: Making a nasty teacher that cannot teach students”, ICLR 2021 (spotlight).

# Proposed architecture – Skeptical Student

Motivation:

- authors observed that impact of nasty teacher on student reduces when the depth of student decreases



- Transfer knowledge to shallow depth ( $\Phi'_S$ ) of a student via aux. classifier (AC)
- Use self-distillation at AC in  $\Phi_S$  -  $\Phi'_S$  to boost performance of student  $\Phi_S$

# Properties

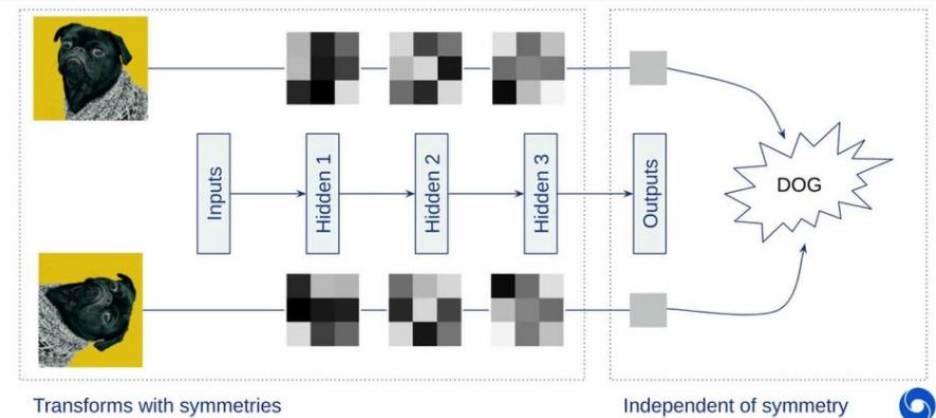
Skeptical students achieve similar to teacher performance even when the teacher is Undistillable (or nasty).

Skeptical students achieve similar to normal students' performance upon distillation from a normal teacher.

# Highlights

- Equivariant/invariant layers
- Data-centric AI
- Demonstrations
- Different approach to Continual Learning

## Respecting symmetries in neural networks

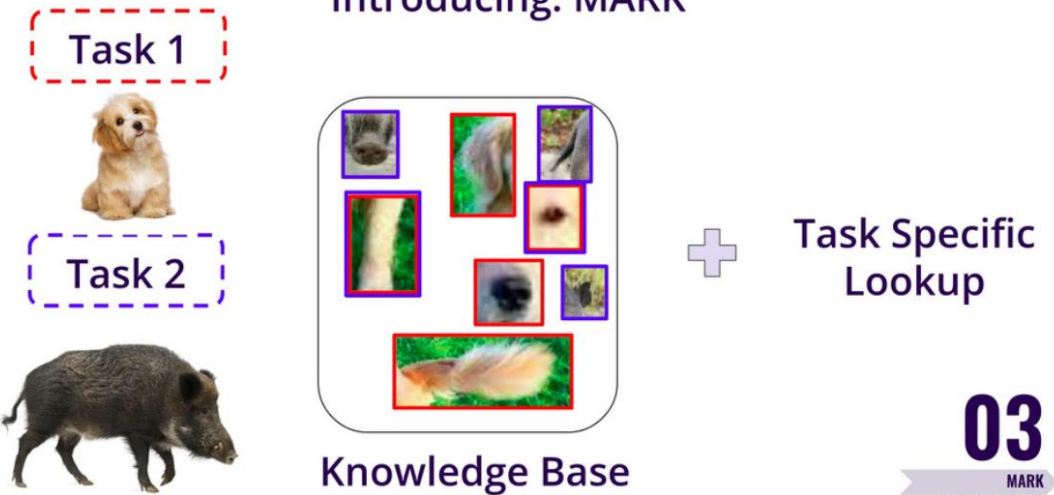


Our idea: Let's maximize knowledge reuse!

Optimizing Reusable Knowledge for Continual Learning via Metalearning  
Julio Hurtado, Alain Raymond-Sáez & Álvaro Soto

14

## Introducing: MARK



**Thank you  
for attention**

