

ICCV 2021
International Conference
on Computer Vision

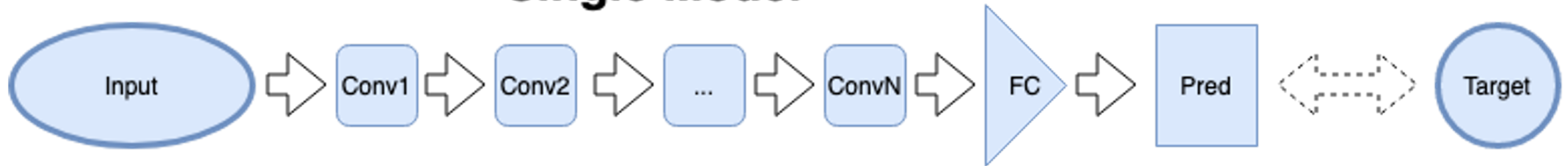
MixMo:

Mixing Multiple Inputs for Multiple Outputs via Deep Subnetworks

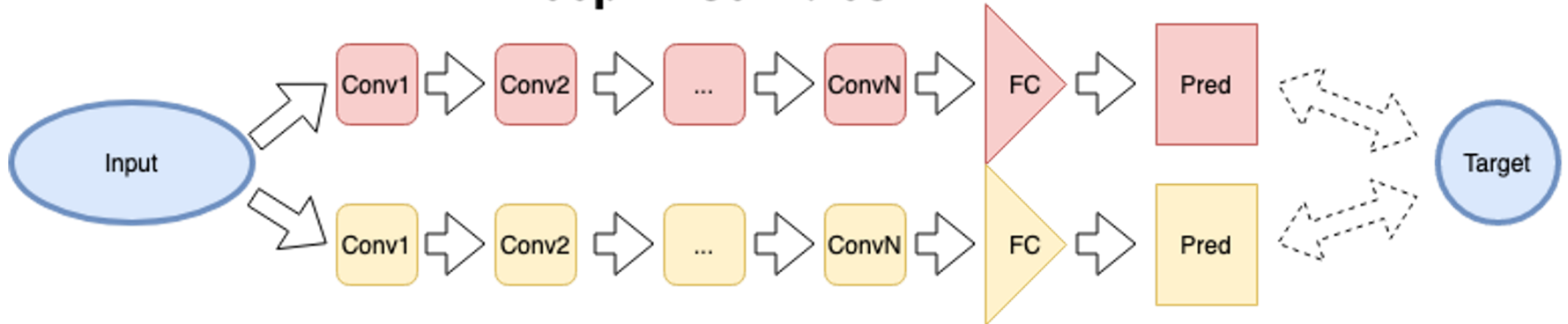
Alexandre Ramé, Rémy Sun
and Matthieu Cord

➤ Deep ensembles

Single Model

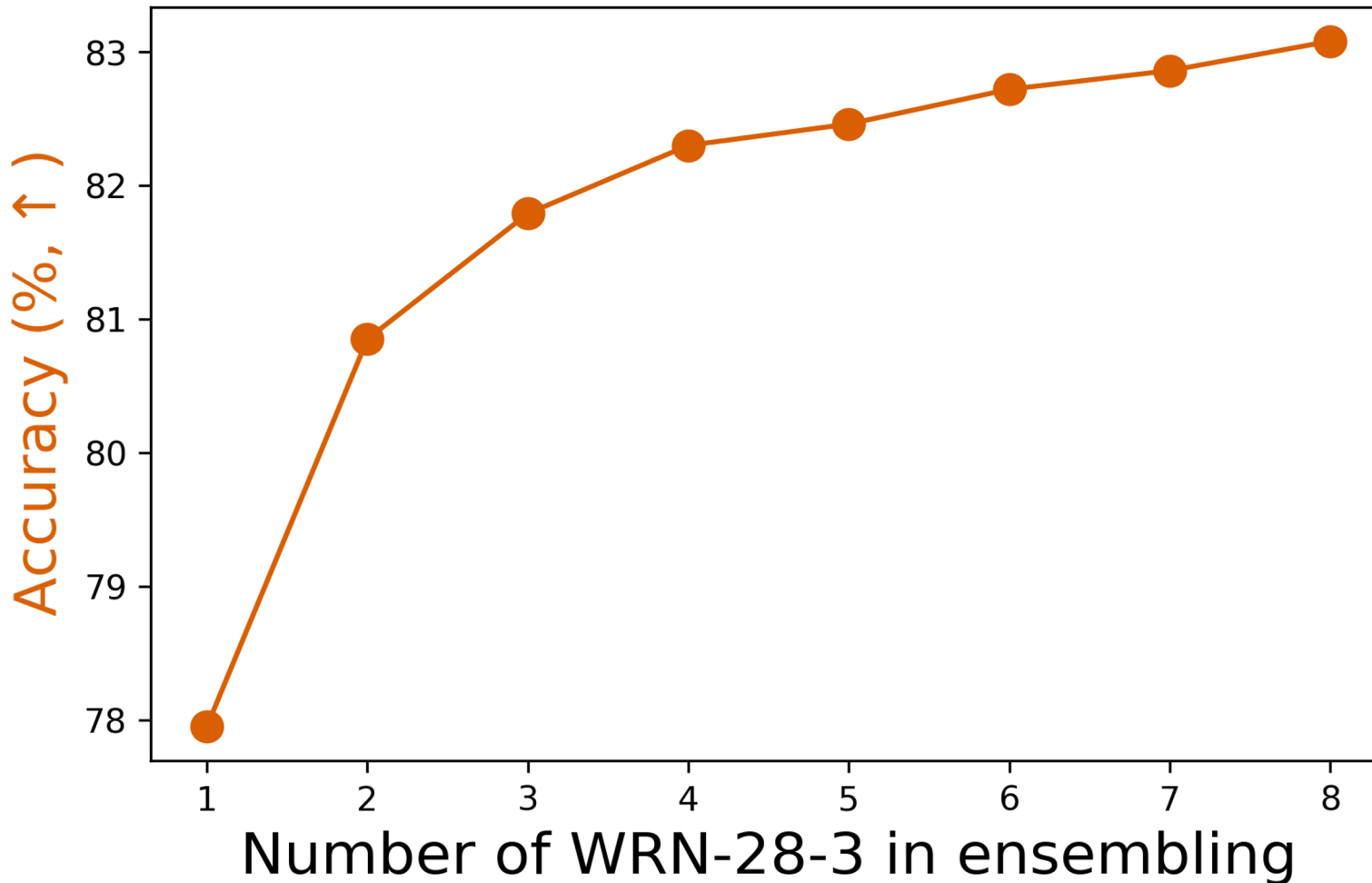


Deep Ensembles



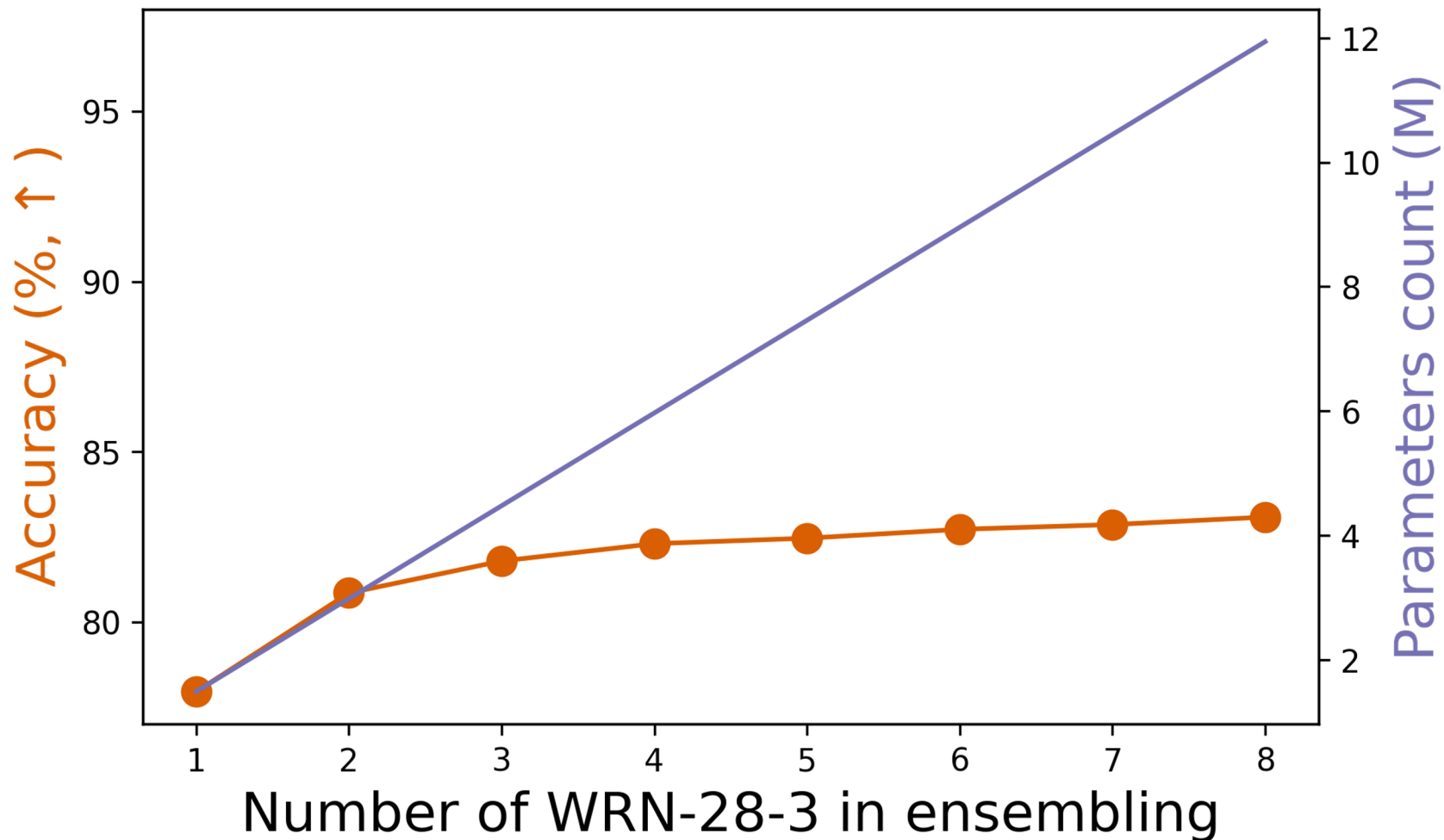


Ensembling improves accuracy ...

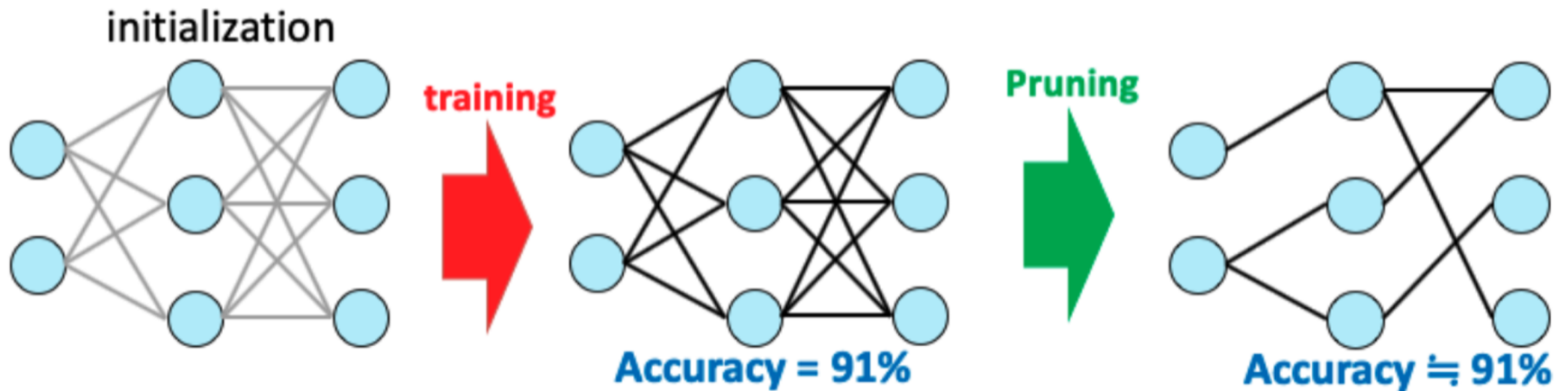




... but ensembling is costly



➤ Idea: leverage network sparsity

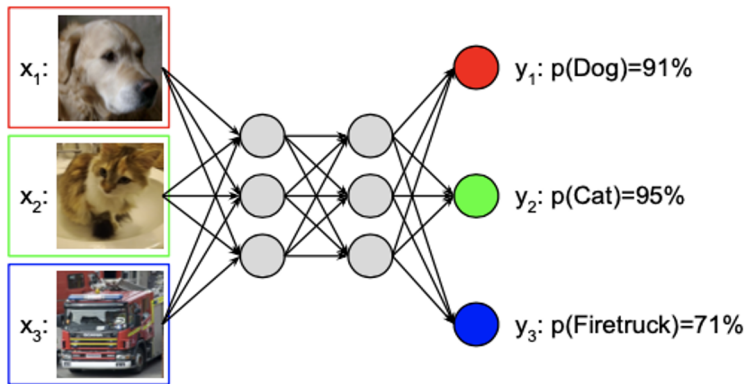


[1] Pruning filters for efficient convnets. Li *et al.*, ICLR 2017

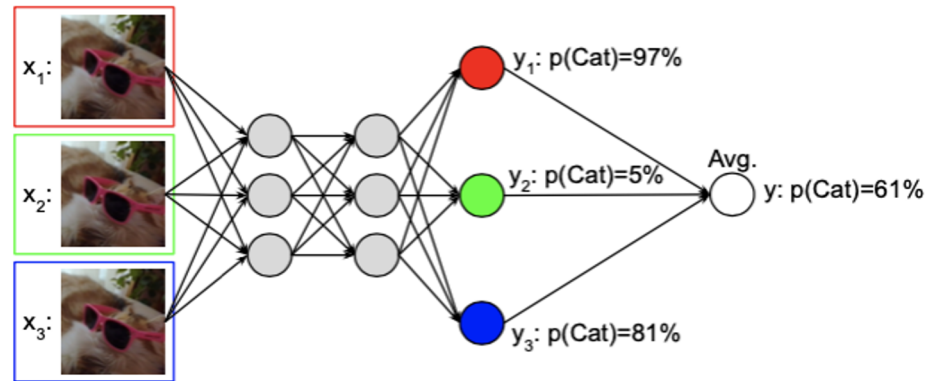
[2] The lottery ticket hypothesis: Finding sparse, trainable neural networks. Jonathan Frankle and Michael Carbin, ICLR 2019

Main idea: multiple subnetworks
inside one base network

➤ Multi-input Multi-output strategy

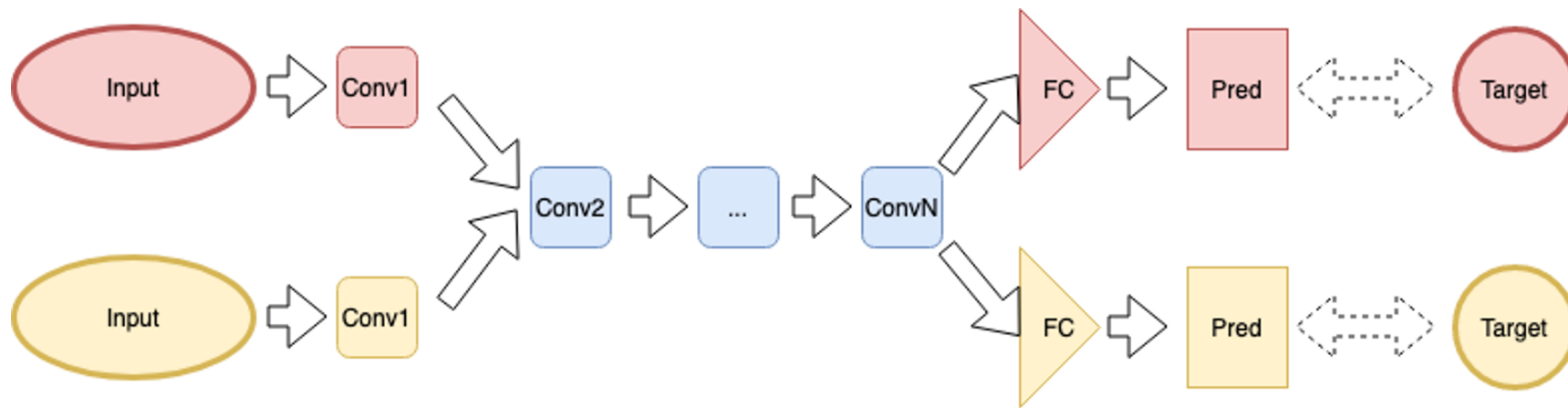


(a) Training

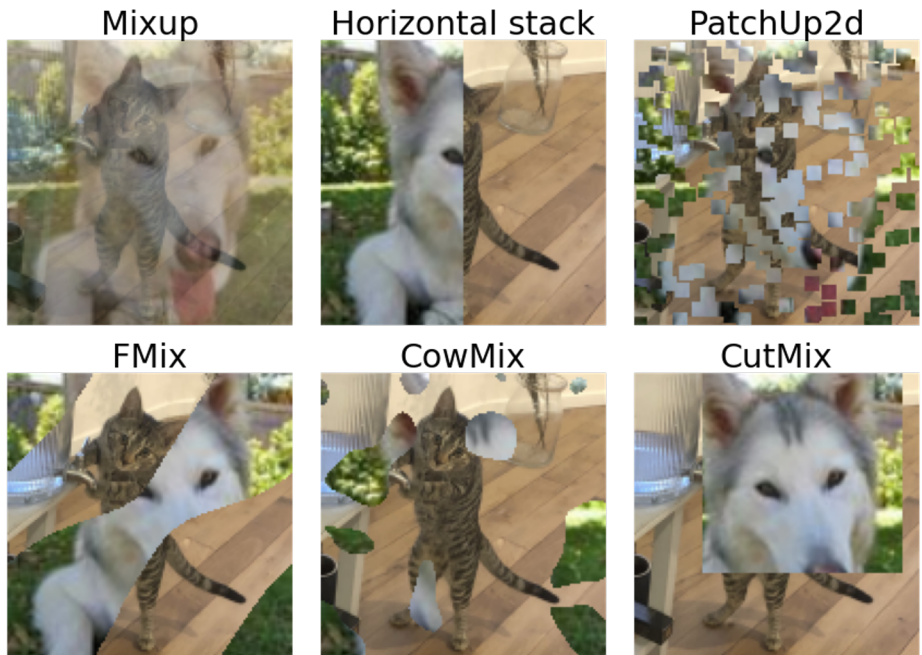
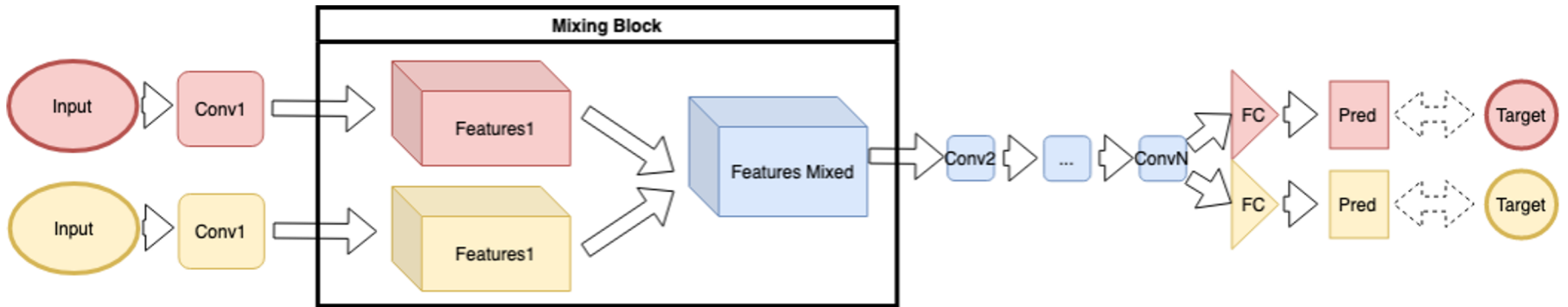


(b) Testing

➤ Architecture: all layers shared except the encoders and the classifiers

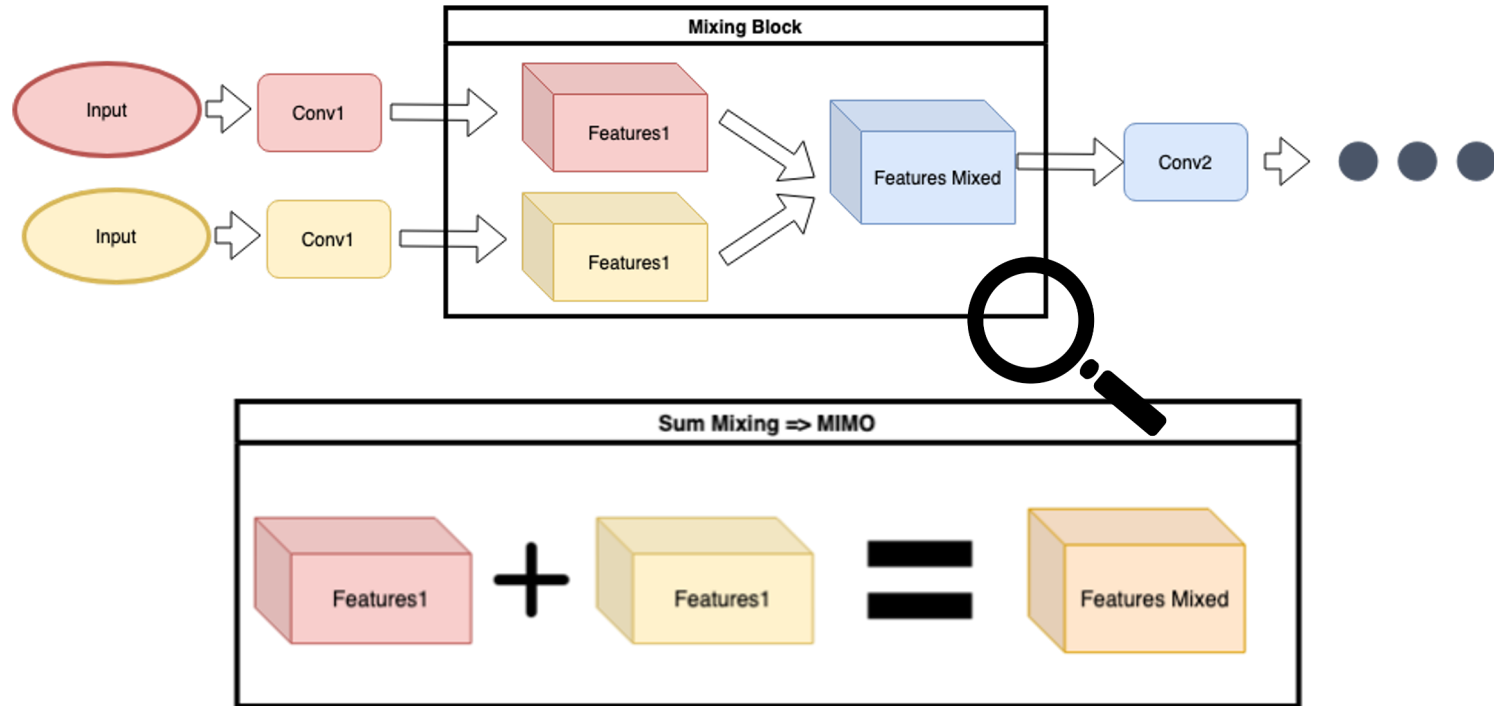


➤ Mixing block combining the inputs' features



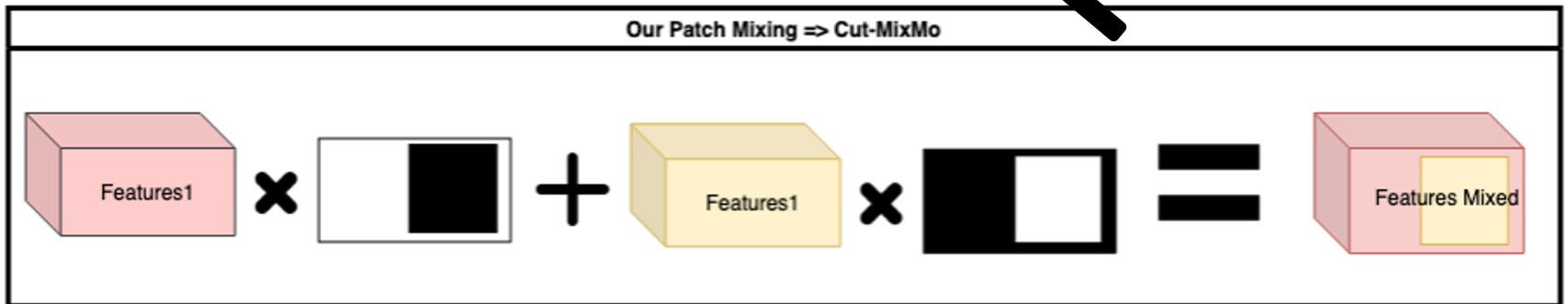
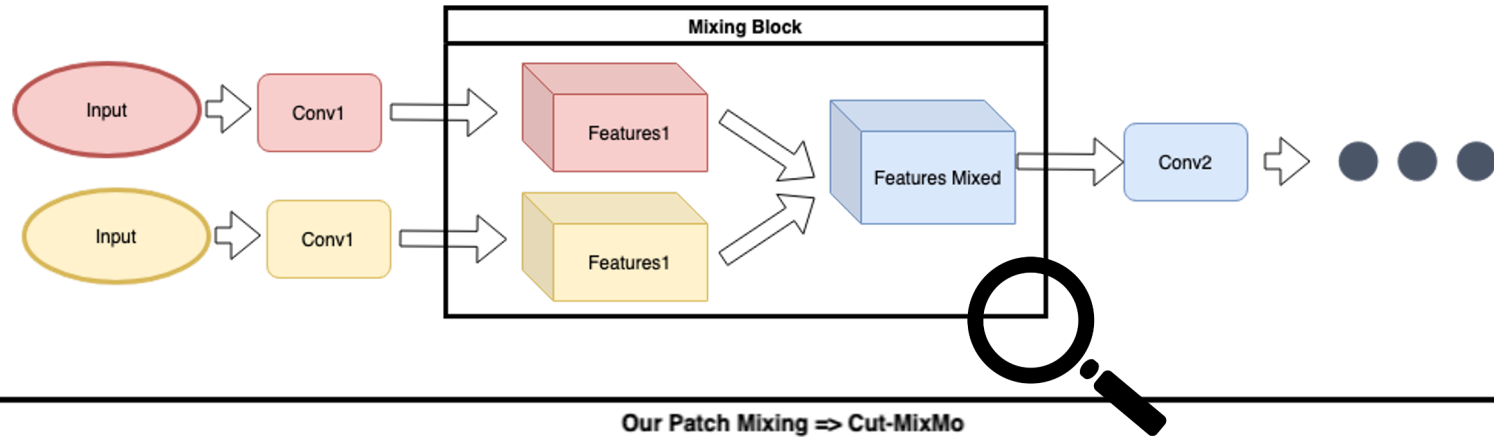


Mixup as mixing block





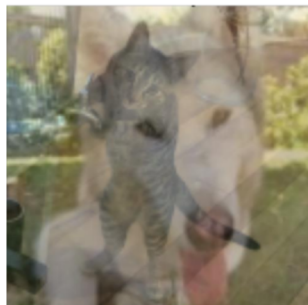
➤ CutMix as mixing block



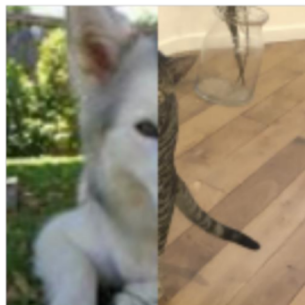


Binary mixing as mixing block improves individual accuracies & ensemble diversity

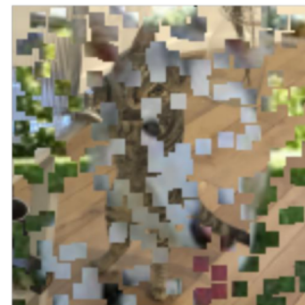
MixUp
82.5%



Concat.
82.78%



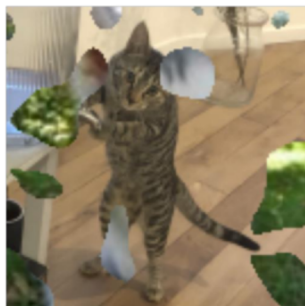
PatchUp
84.16%



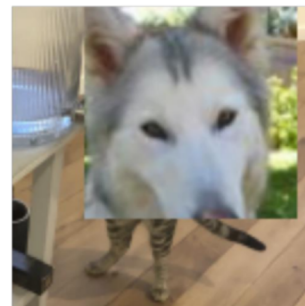
FMix
83.76%



CowMix
84.17%



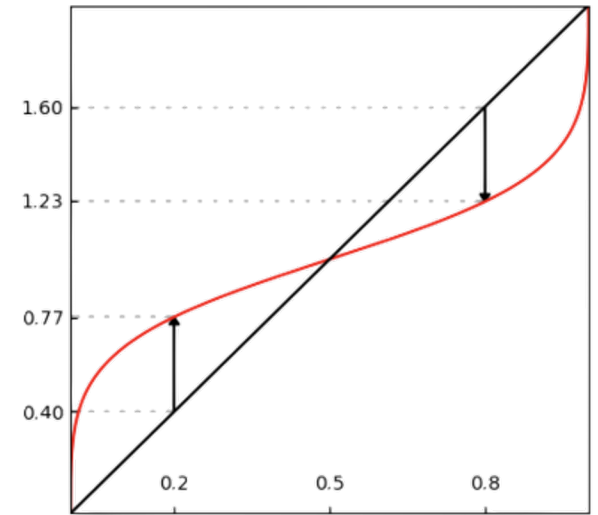
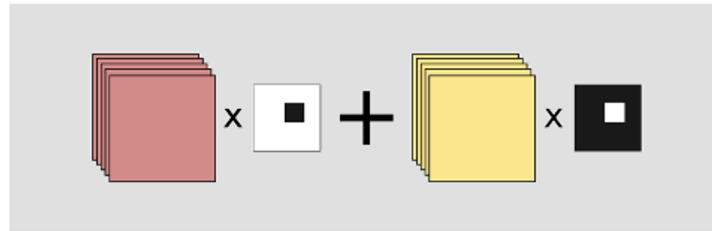
CutMix
84.38%



➤ Balancing the training losses

$$\kappa \sim \beta(\alpha, \alpha) \quad \text{---} \rightarrow \quad \kappa \rightarrow w_r(\kappa)$$

$$\square \sim \mathbf{1}_{\mathcal{M}}(\kappa) \rightarrow$$



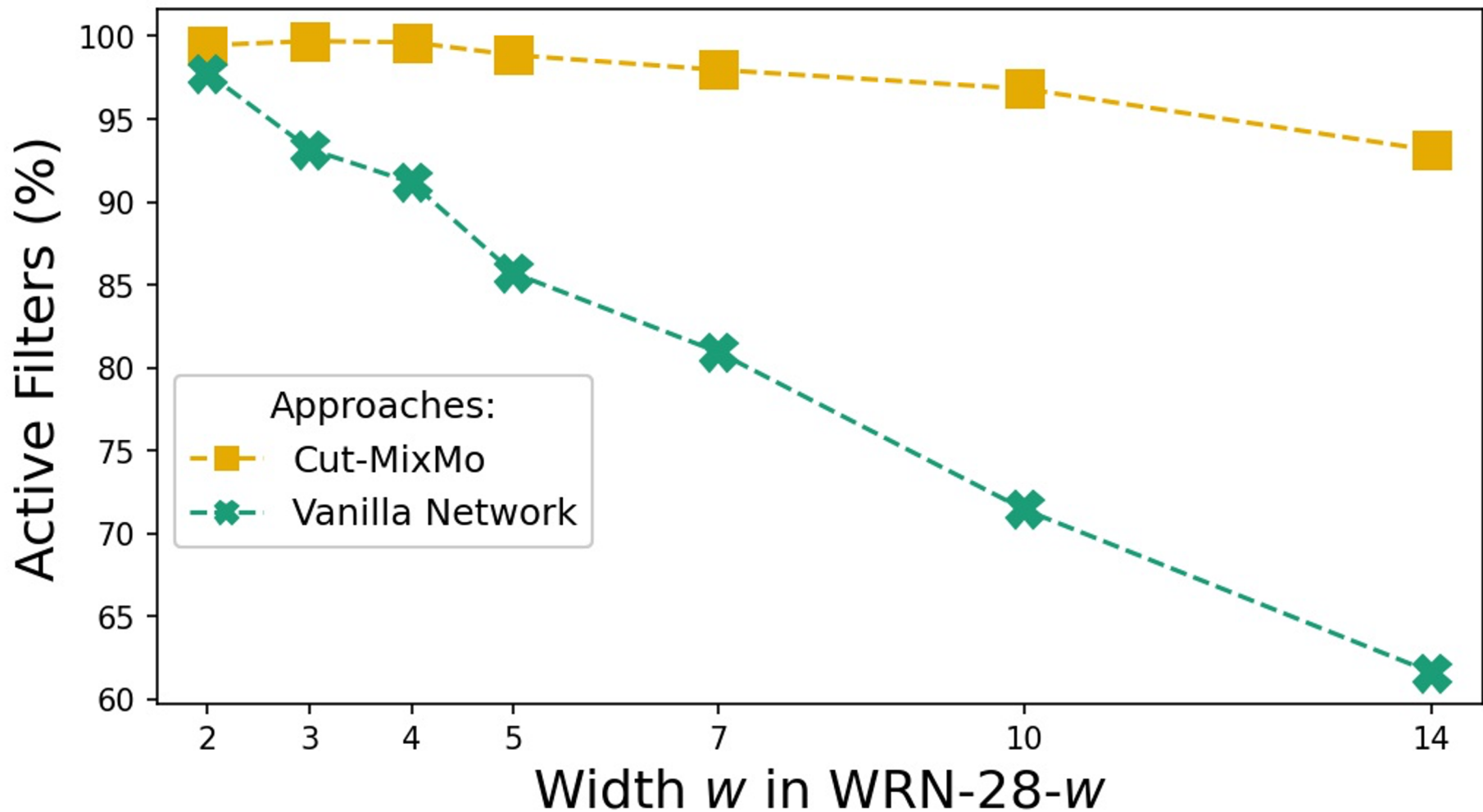
$$\mathcal{L}_{\text{MixMo}} = w_r(\kappa) \mathcal{L}_{\text{CE}}(\text{Cat}, \text{Pencil}) + w_r(1-\kappa) \mathcal{L}_{\text{CE}}(\text{Dog}, \text{Pencil})$$



State of the art on CIFAR and TinyImageNet

| Approach | #Params | WRN-28-10 | | ResNet-18-3 |
|-----------|---------|--------------|--------------|--------------|
| | | CIFAR100 | CIFAR10 | TinyImageNet |
| Vanilla | 1.0 | 81.63 | 96.34 | 65.78 |
| CutMix | 1.0 | 84.05 | 97.23 | 68.95 |
| Deep Ens. | 2.0 | 83.17 | 96.67 | 68.38 |
| MIMO | 1.002 | 83.06 | 96.74 | 68.48 |
| Cut-MixMo | 1.002 | 85.40 | 97.51 | 70.24 |

➤ Better leverages over-parameterization





Contributions

❖ Theoretically

Unifying framework for multi-input multi-output ensembling

Connection with data augmentation

❖ Empirically

State of the art at same inference cost as a vanilla network

More in paper: ImageNet, robustness, memory split advantage

<https://github.com/alexrame/mixmo-pytorch>

Merci !