



SAPIENZA
UNIVERSITÀ DI ROMA

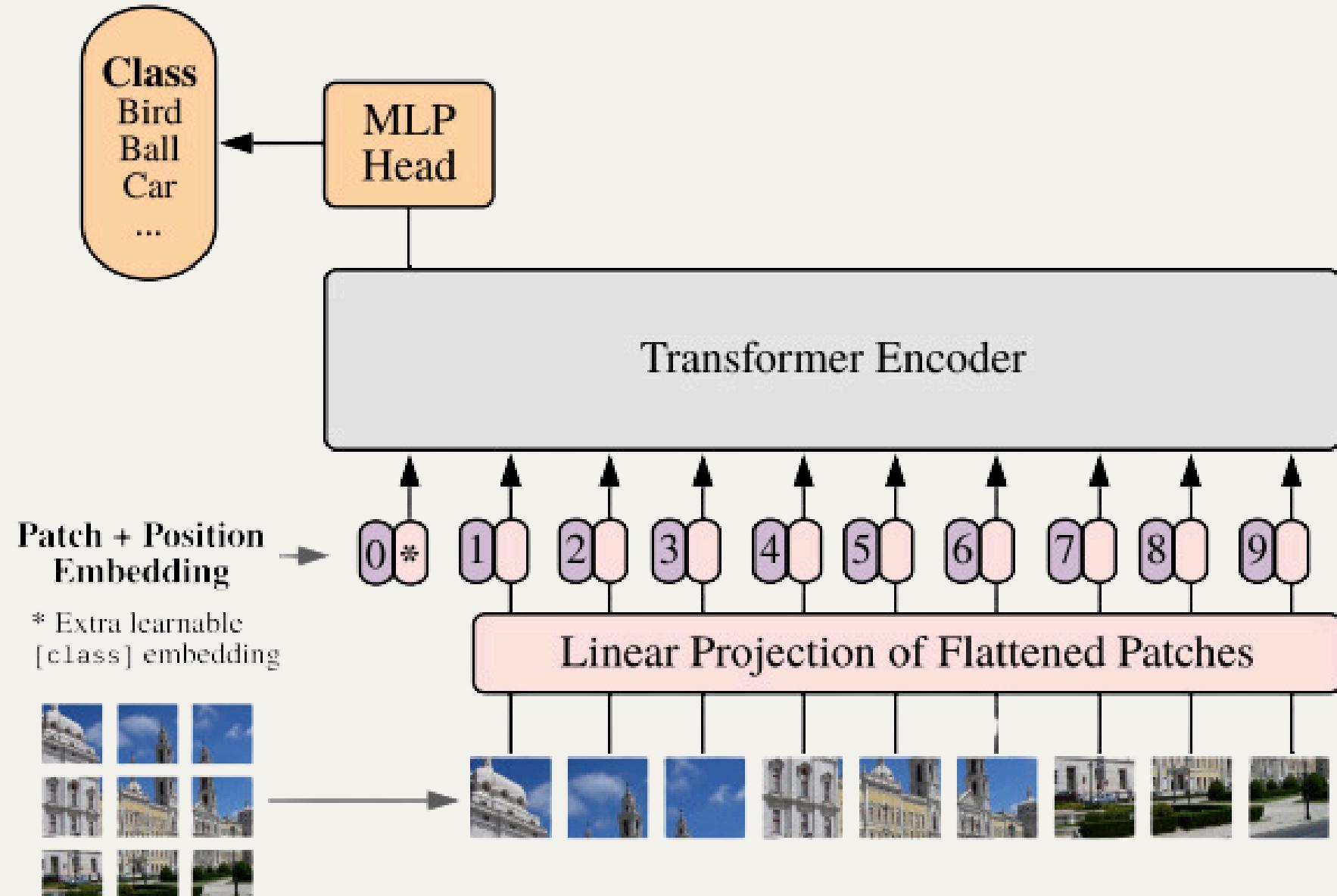
From One to Many: Virtual Subnetworks for Vision Transformer Ensembles

Alessandro Ferrante

Engineering in
Computer Science

Professor Simone Scardapane

Vision Transformers



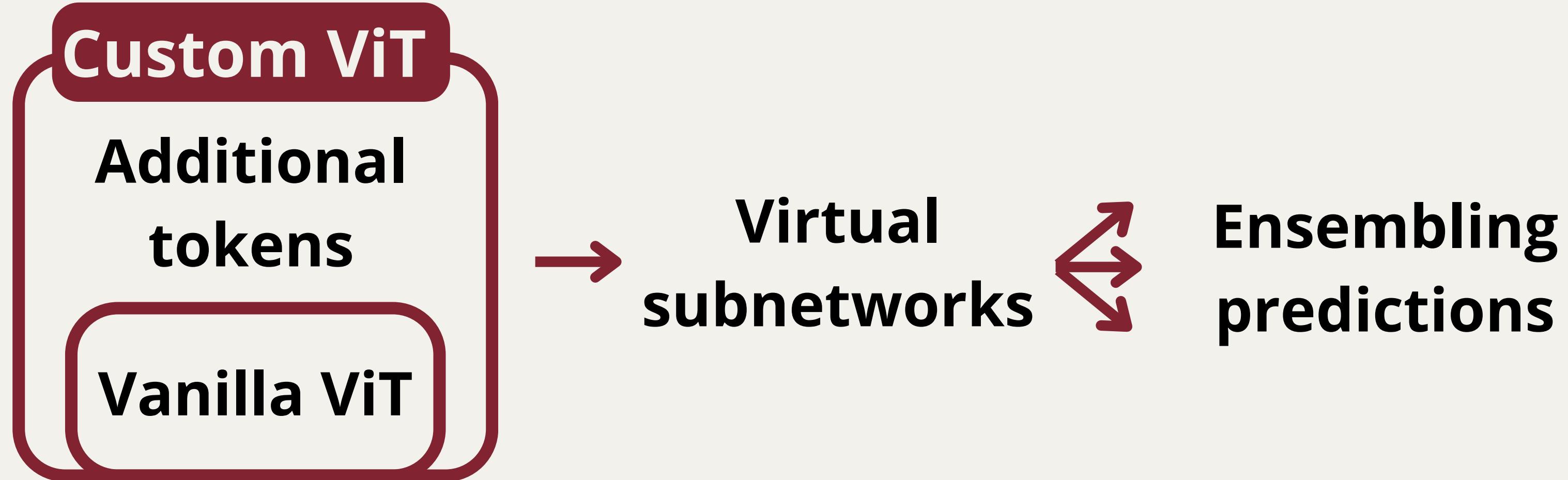
PRO

- Global context awareness
- State of the art performance (>CNN)
- Flexible to many inputs

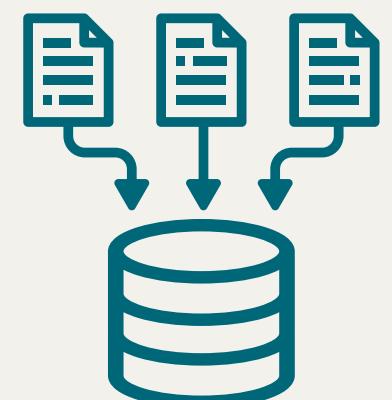
CONS

- Data-hungry
- Computational costs
- Vulnerable to noise and attacks

Extending ViT



Only 1
training

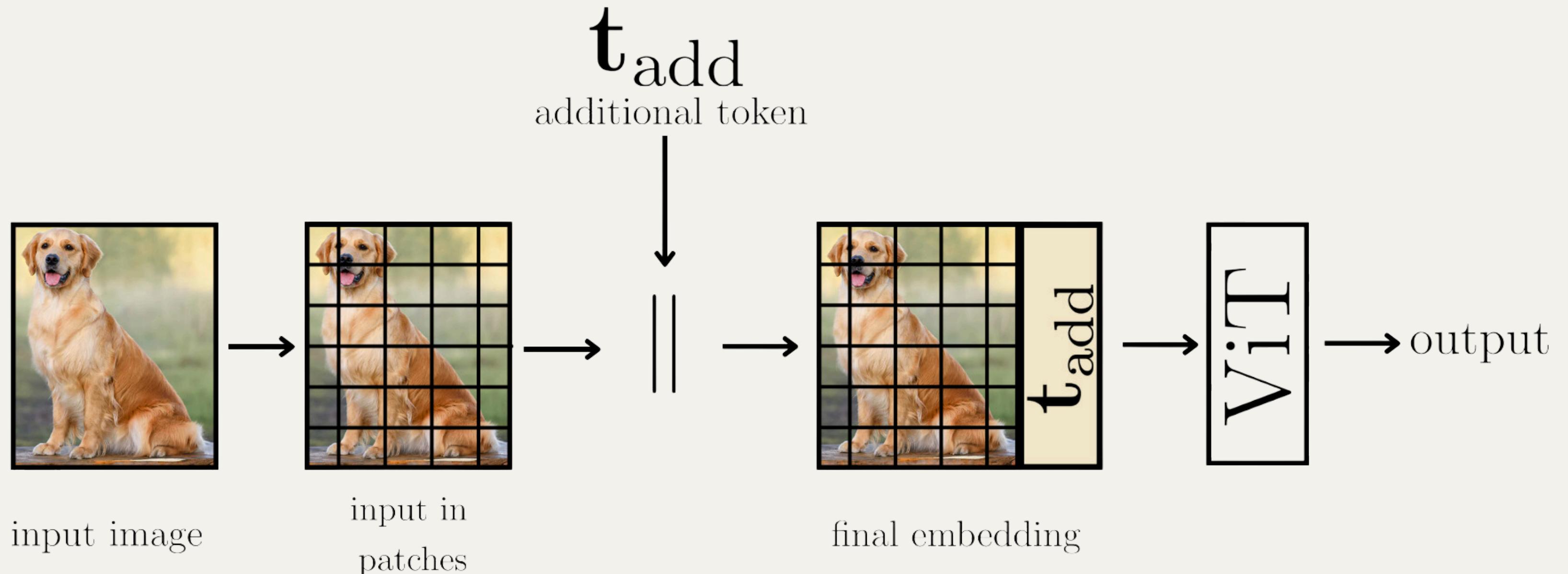


Recycling
data



Robustness

Virtual Subnetworks





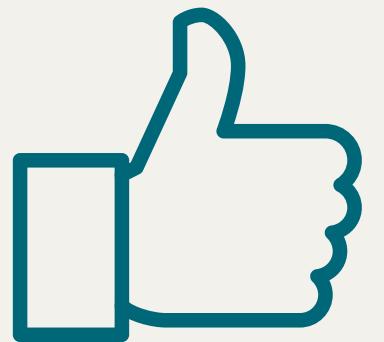
Ensembling



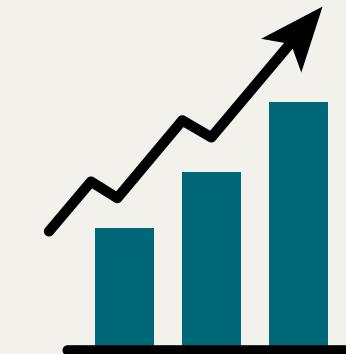
Goals



Maintain accuracy



Increase confidence



Improve generalization



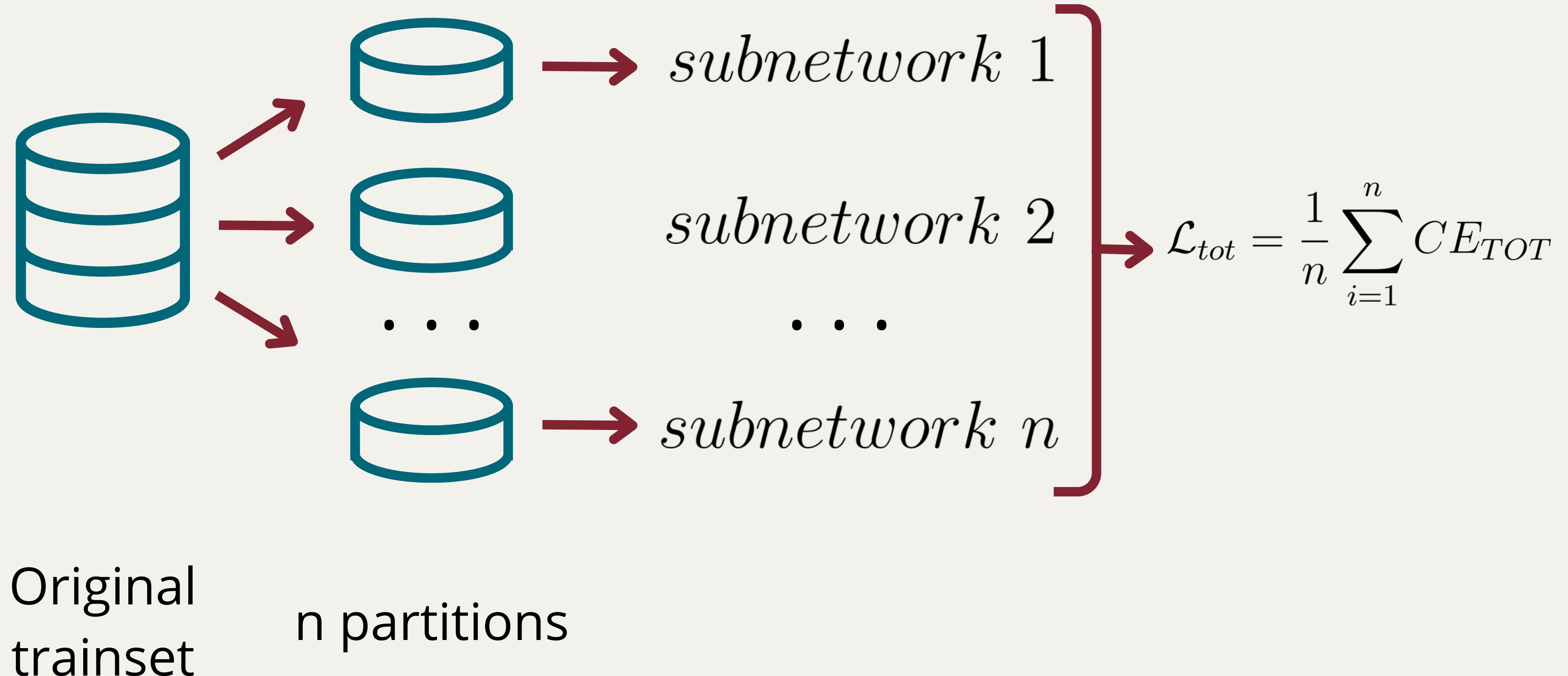
Ensemble by averaging

$$\text{Ensemble Logits} = \frac{1}{n} \sum_1^n \text{Logits list}$$

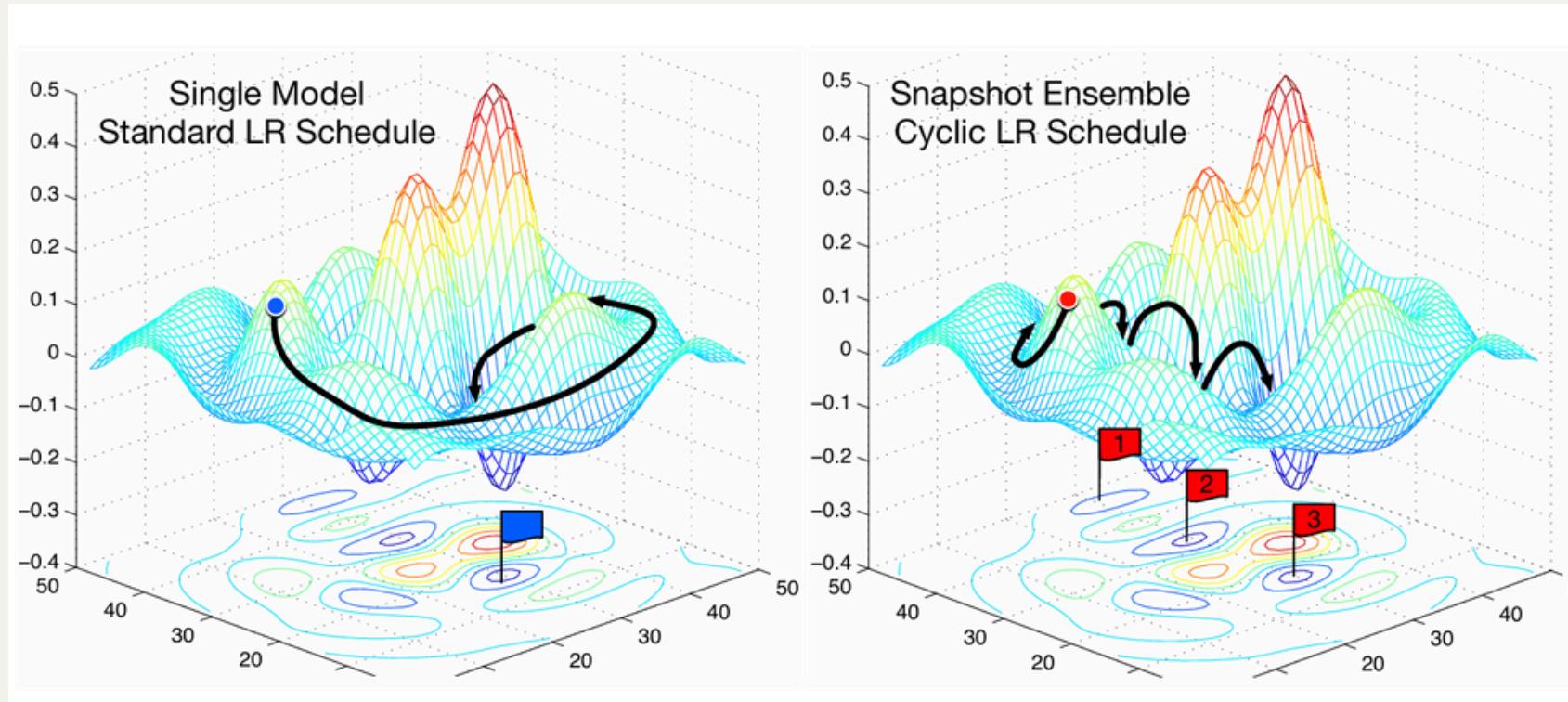
$$\text{KL Divergence} = \frac{1}{n} \sum_{i=1}^n \mathcal{D}_{\text{KL}_i}(P_{\text{Logits}_i} \| P_{\text{Ensemble Logits}})$$

$$\mathcal{L} = \mathcal{L}_{CE_{\text{ensemble}}} - \lambda_{\text{KL}} \cdot \mathcal{D}_{\text{KL}}$$

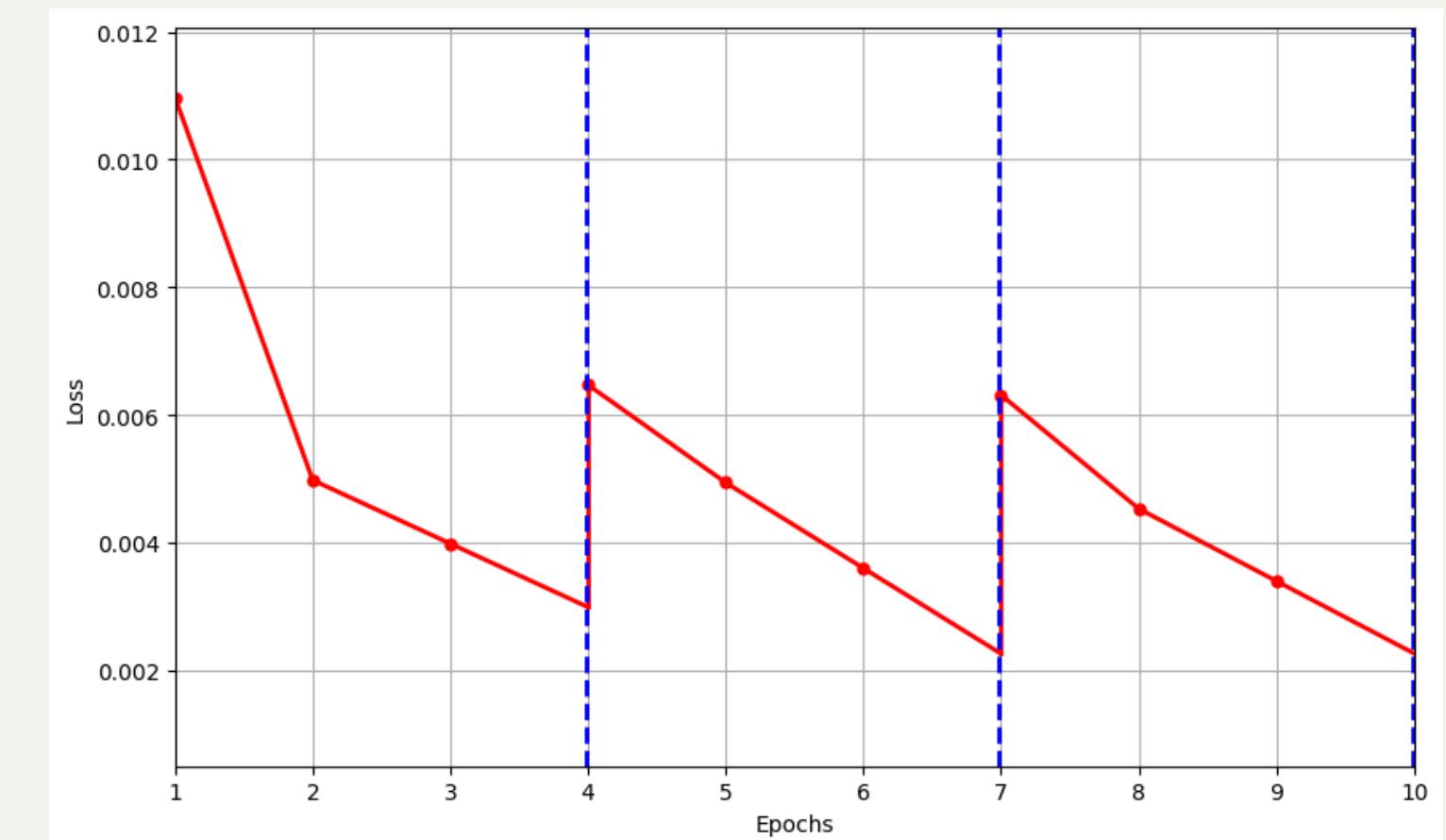
Ensemble by data partitioning



Snapshot Ensemble



- ① Still a **single** training
- ② Up to m snapshots
- ③ Even less computation

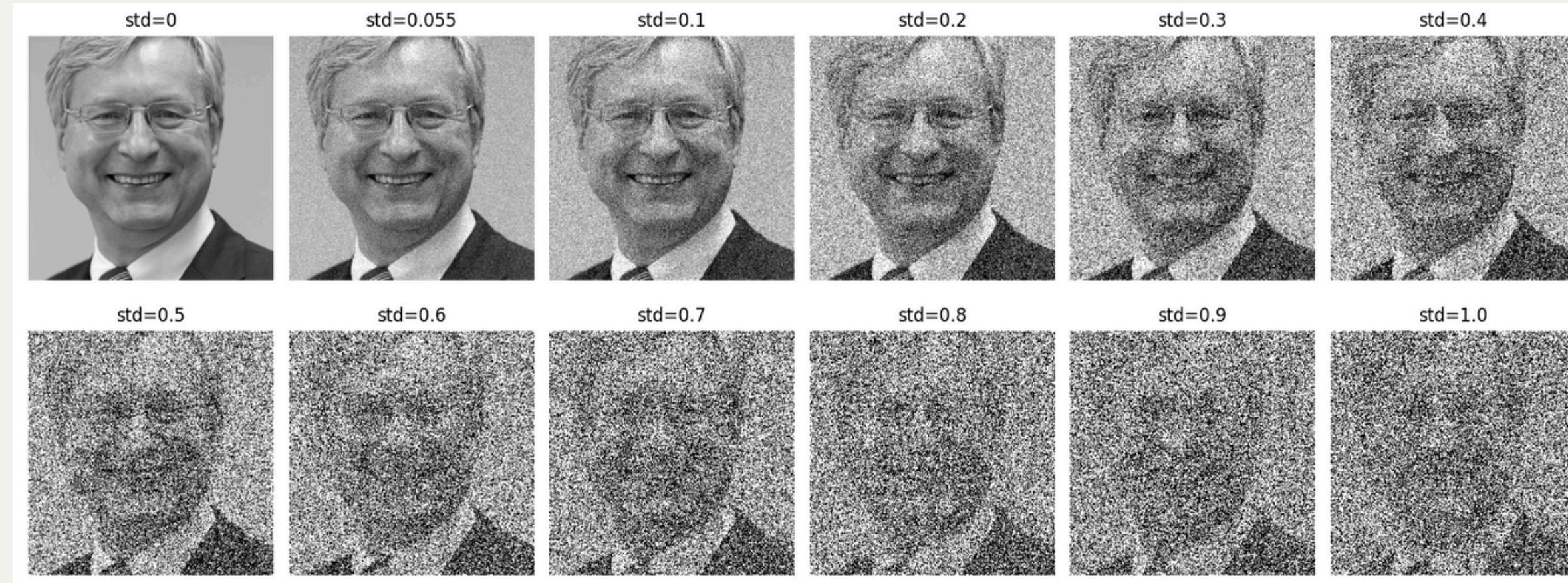




Preliminary results

	Accuracy	ECE
Standard ViT		
ViT Vanilla	0.944	+4.5% 0.042 -66.7%
2x Vanilla	0.948	+4.0% 0.040 -65.0%
3 Subnetworks Averaging	0.986	0.014
5 Subnetworks Partitioning	0.986	0.015
9 Subnetworks Snapshot	0.957	0.013

Gaussian noised input

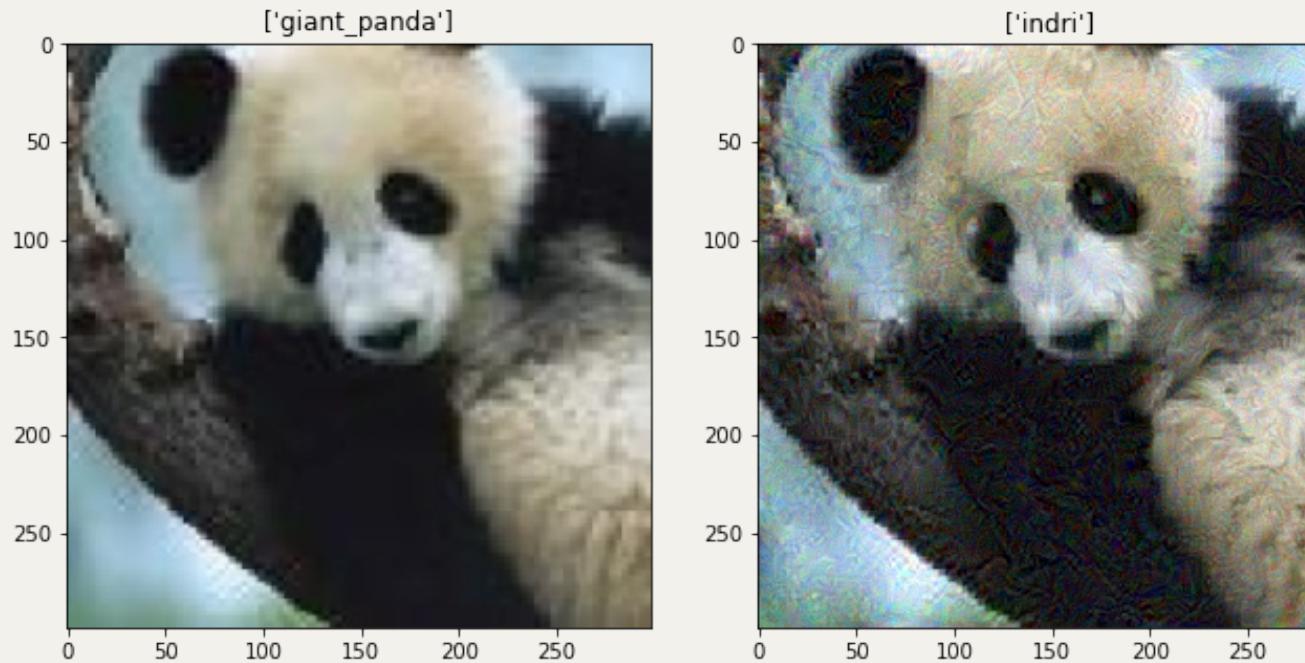


Accuracy		ECE	
3 Subnetworks Averaging	std=0.1 +16.4%	std=1.0 +28.3%	std=0.1 -55.7%
	5 Subnetworks Partitioning	+12.6% +25.8%	5 Subnetworks Partitioning -37.4%
	9 Subnetworks Snapshot	+13.2% +25.8%	9 Subnetworks Snapshot -40.7%
3 Subnetworks Averaging	std=1.0 -56.8%	std=1.0 -33.8%	std=1.0 -41.5%

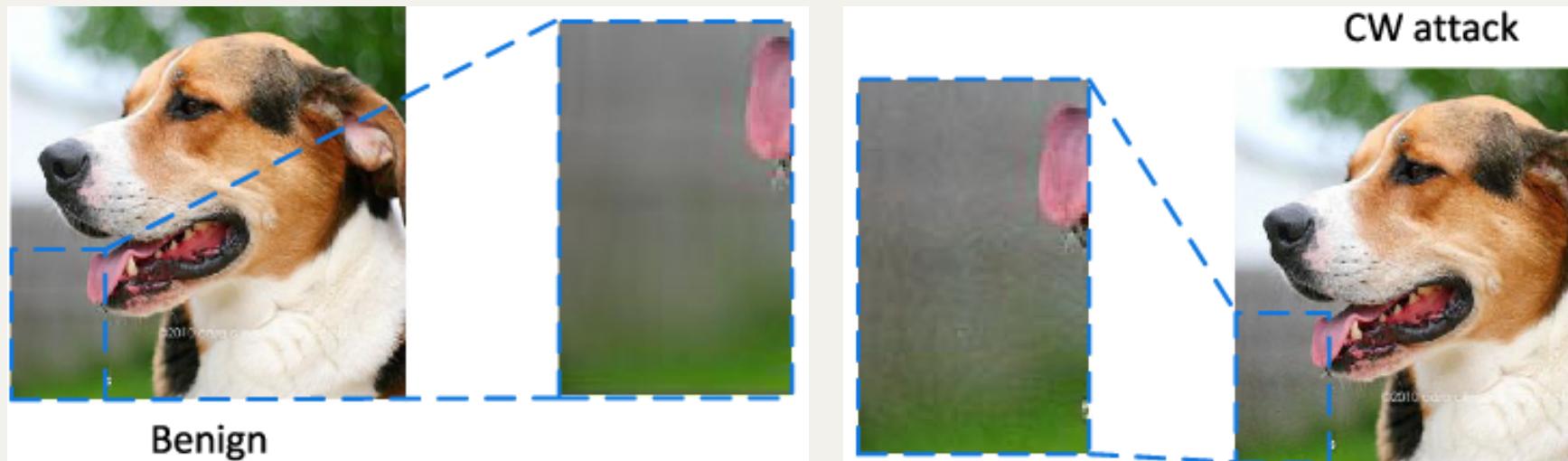


Adversarial attacks

FGSM - PGD



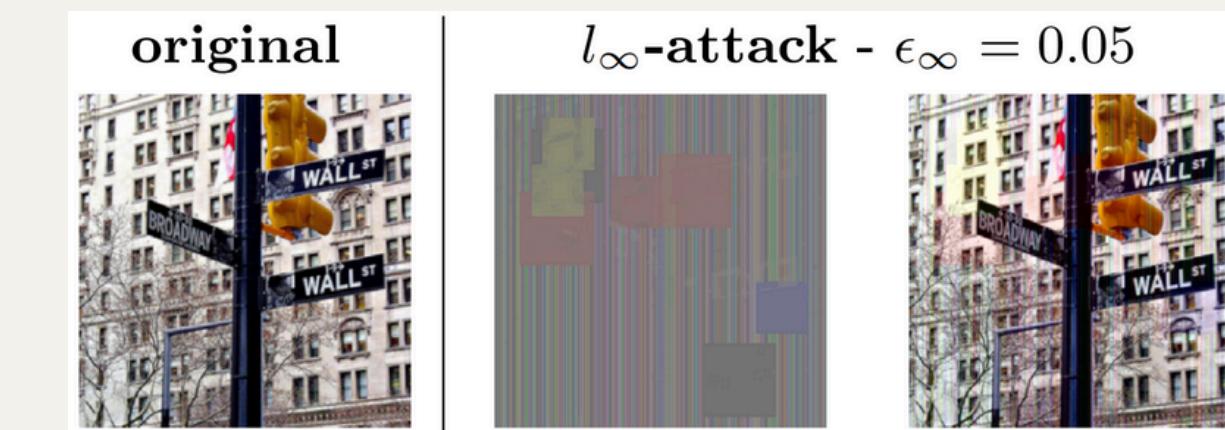
CW



White box

Black box

Square attack



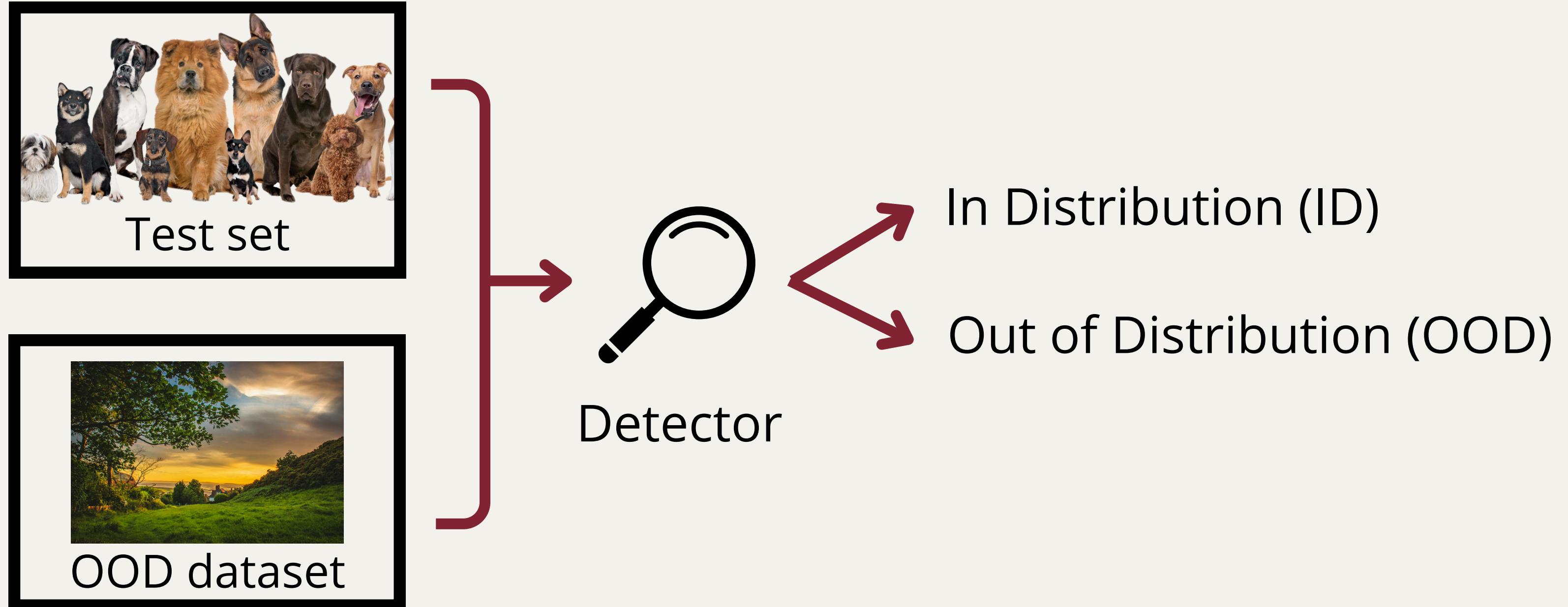


Adversarial attacks (2)

Accuracy and ECE comparison:

	FGSM		PGD		CW		Square		Standard
	0.62	0.29	0.68	0.22	0.68	0.21	0.75	0.18	
ViT Vanilla	0.62	0.29	0.68	0.22	0.68	0.21	0.75	0.18	
2 Vanilla	0.63	0.30	0.70	0.26	0.70	0.20	0.75	0.18	
2 Subnetworks Average	0.79	0.16	0.78	0.16	0.78	0.14	0.88	0.09	
3 Subnetworks Partitioning	0.80	0.18	0.78	0.20	0.81	0.14	0.89	0.08	Custom ViT
9 Subnetworks Snapshot	0.62	0.32	0.71	0.25	0.72	0.20	0.78	0.16	
	+28%	-46%	+13%	-33%	+14%	-32%	+16%	-56%	

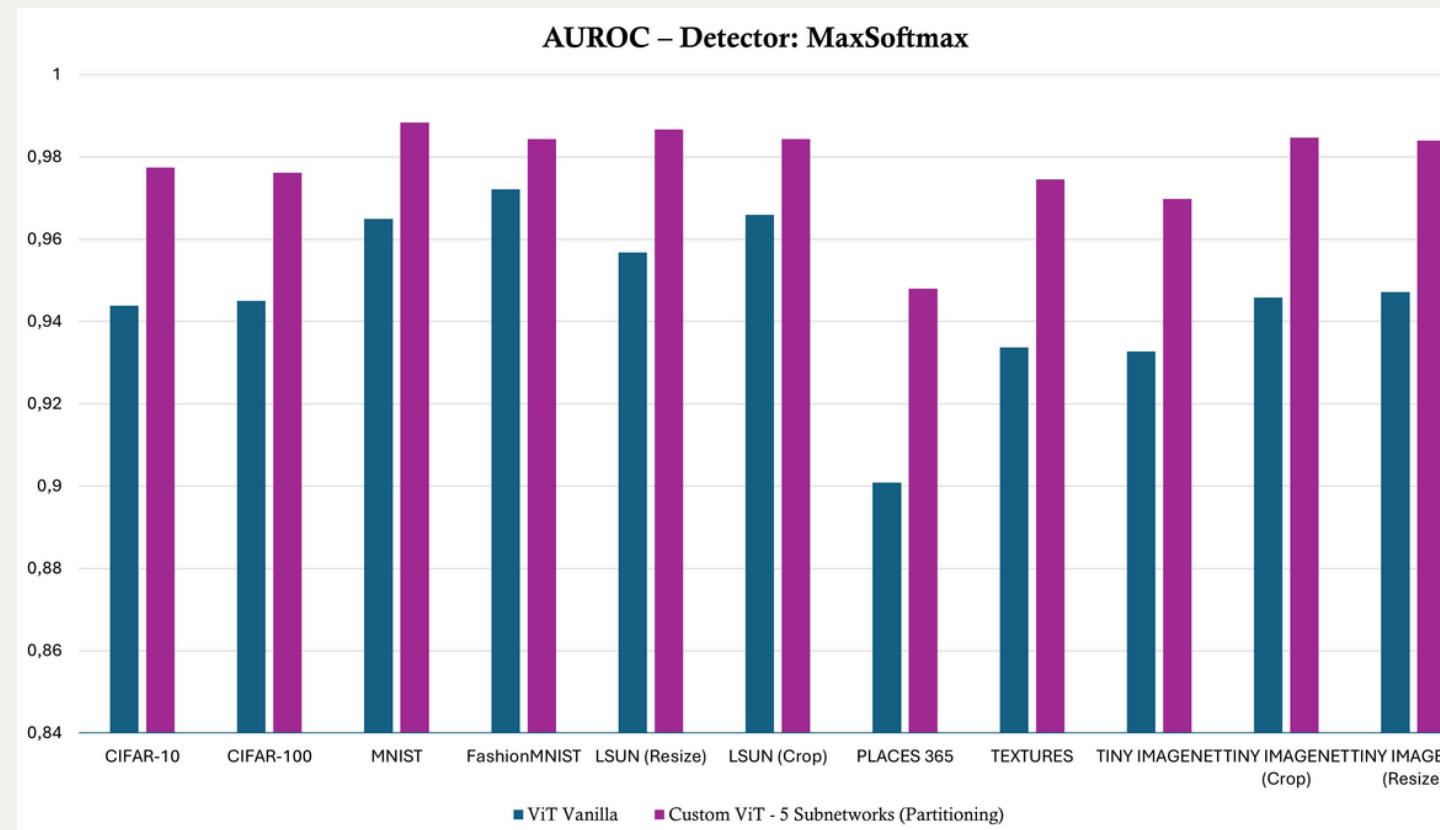
OOD analysis



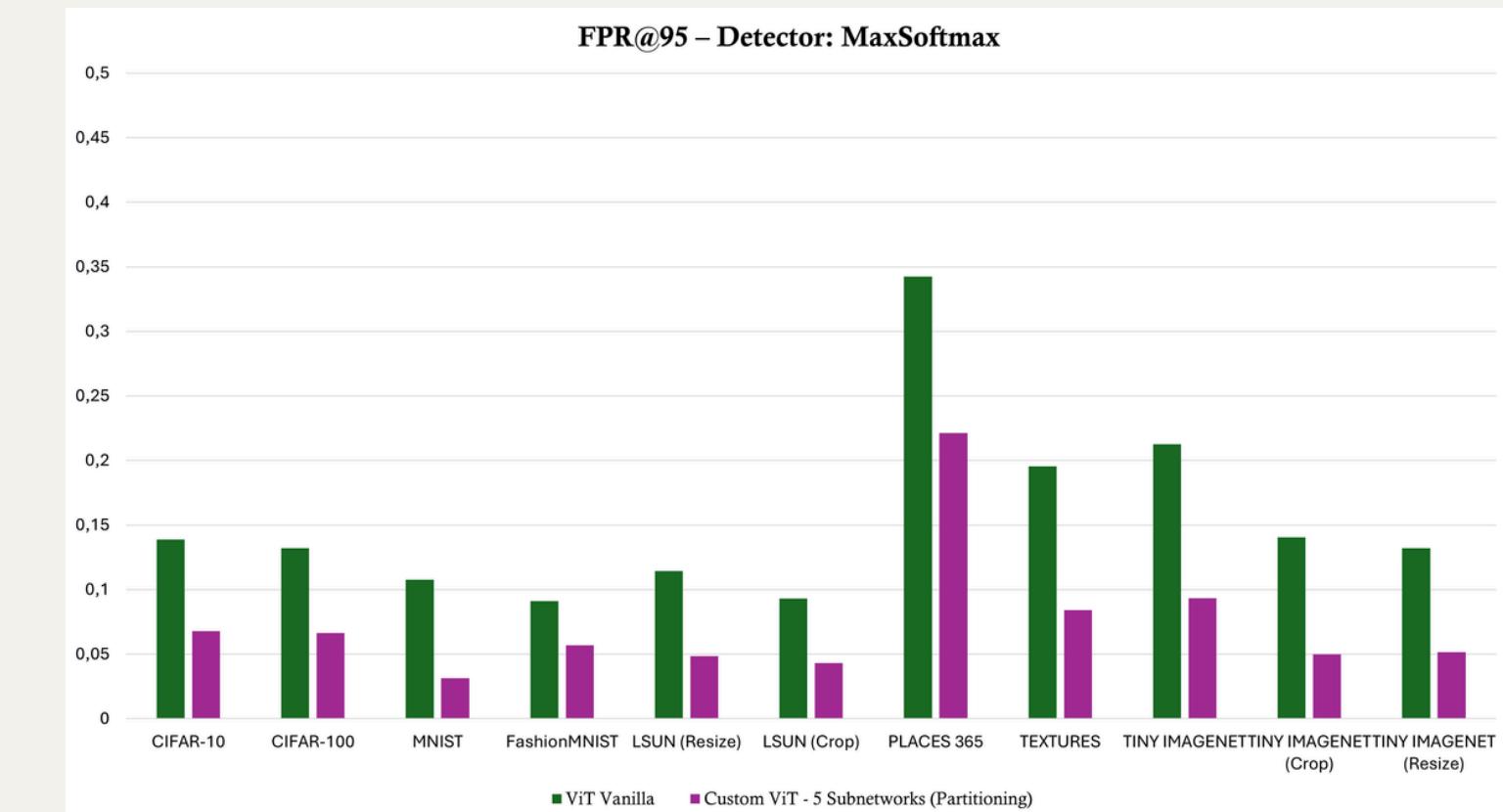


OOD analysis (2)

Detector: MaxSoftmax



Average increment: +3.36%

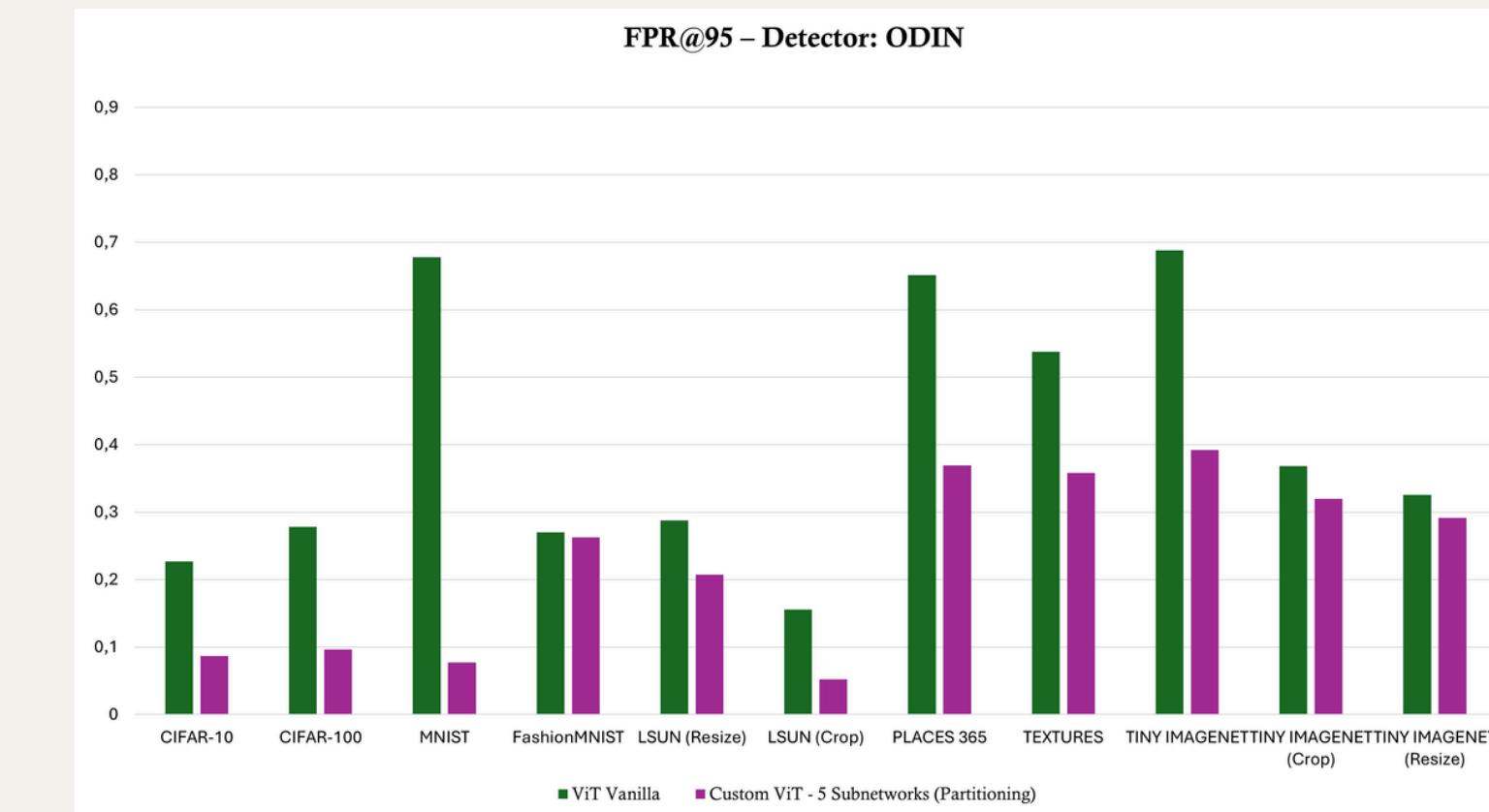
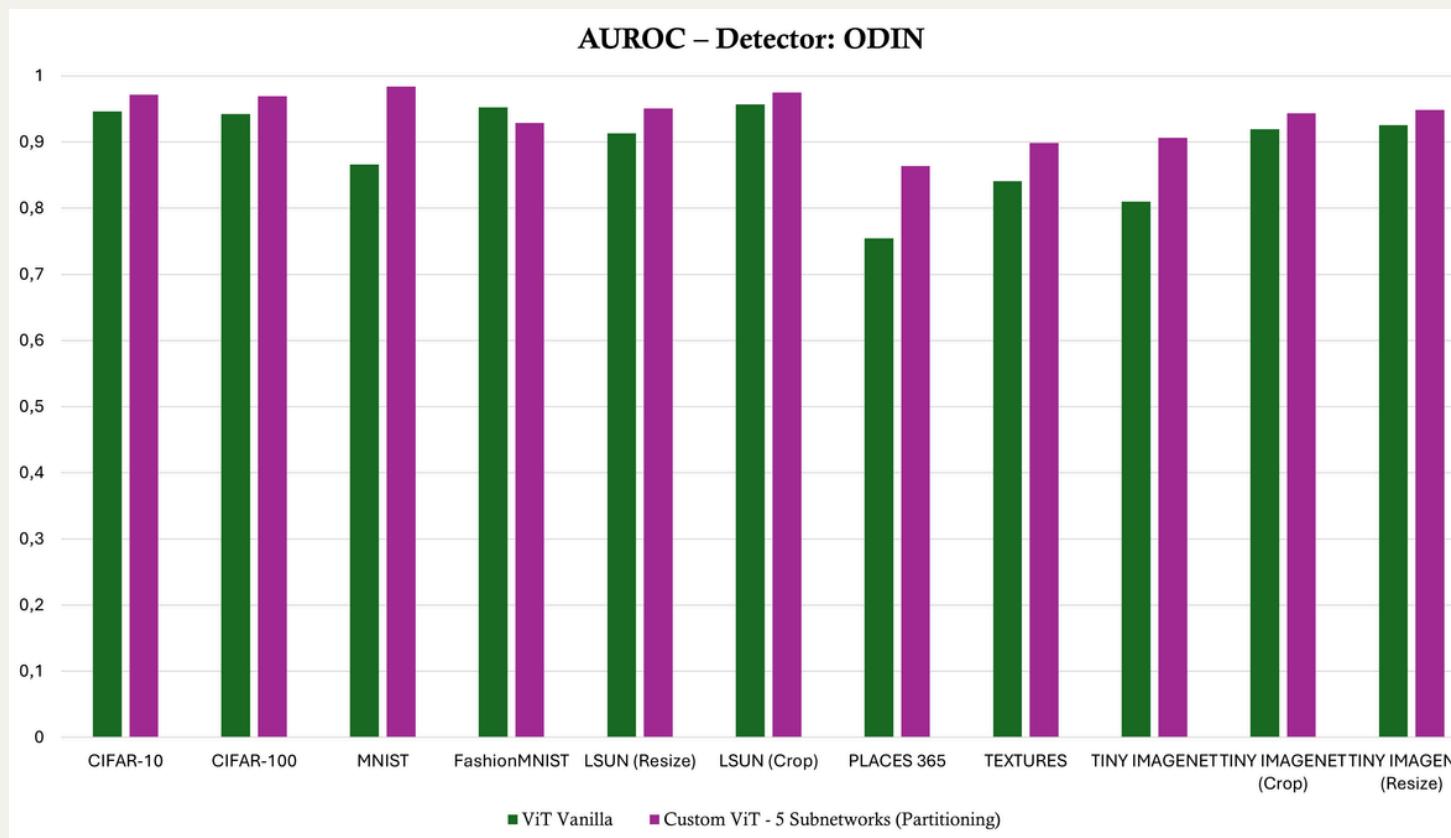


Average reduction: -52.11%

OOD analysis (3)



Detector: ODIN



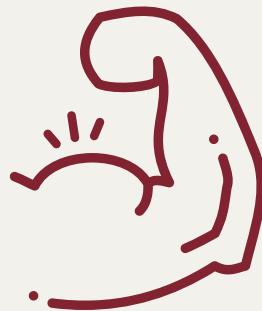
Average increment: +5.23%

Average reduction: -17.76%

Conclusions



The proposed models outperform vanilla ViTs in both accuracy and confidence.



Ensembling boosts robustness against noise, attacks, and OOD tasks.



Custom ViT flexibility allows ad hoc solutions for specific tasks.