

UNIC: Universal Classification Models via Multi-teacher Distillation



Mert Bulent SARIYILDIZ

Philippe WEINZAEPFEL

Thomas LUCAS

Diane LARLUS

Yannis KALANTIDIS

Goal

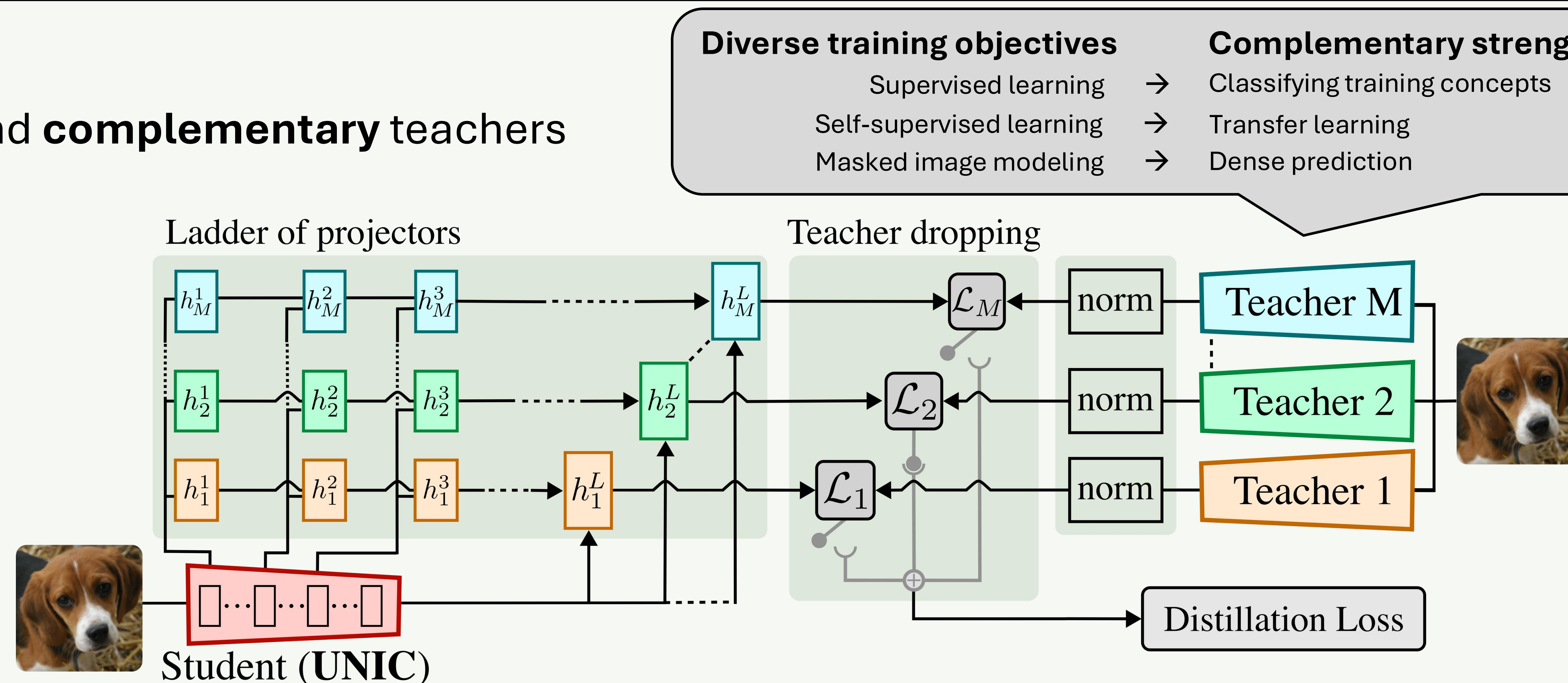
Learn a **single** encoder from **strong** and **complementary** teachers

Approach

Multi-teacher distillation

3 components

- Feature normalization
- Ladder of projectors
- Teacher dropping regularization



Diverse training objectives

- Supervised learning
- Self-supervised learning
- Masked image modeling

Complementary strengths

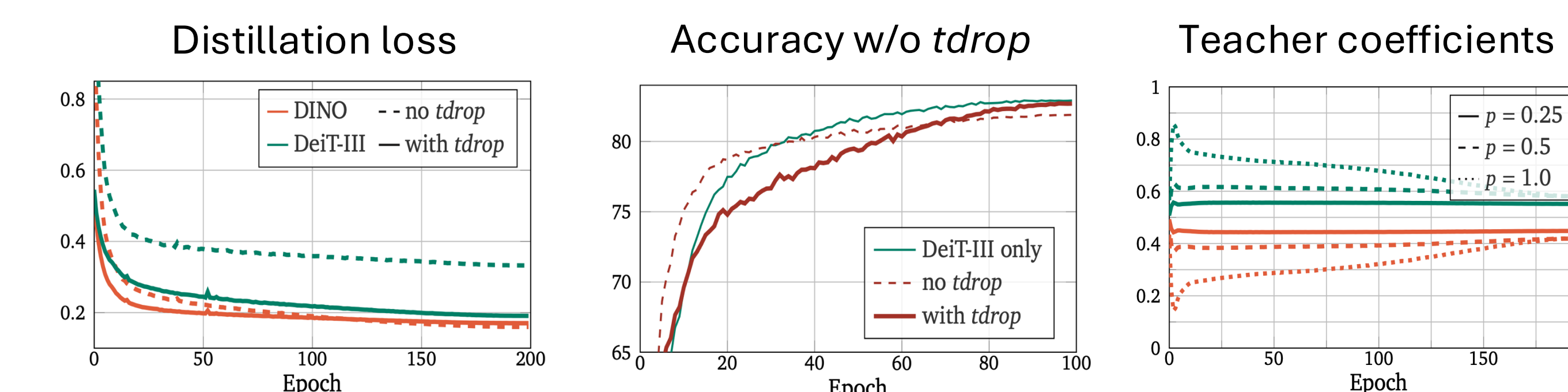
- Classifying training concepts
- Transfer learning
- Dense prediction

Universal encoder: A single model that combines the teacher's strengths

Impact of components

Model	Feature normalization (<i>norm</i>)		Ladder of projectors (<i>LP</i>)		Teacher dropping (<i>tdrop</i>)	
	IN-val top-1 (↑)	Transfer top-1 (↑)	Segmentation mIoU (↑)	Depth RMSE (↓)		
<i>Teacher models</i>						
DINO	77.7	72.4	30.4	0.570		
DeiT-III	83.6	68.5	32.3	0.589		
<i>best teacher</i>	83.6	72.4	32.3	0.570		
<i>Multi-teacher distillation (DINO & DeiT-III teachers)</i>						
basic setup	78.7	73.1	33.9	0.560		
UNIC	81.4	73.8	36.1	0.558		
	82.7	74.2	37.4	0.546		
	83.2	73.5	37.3	0.547		

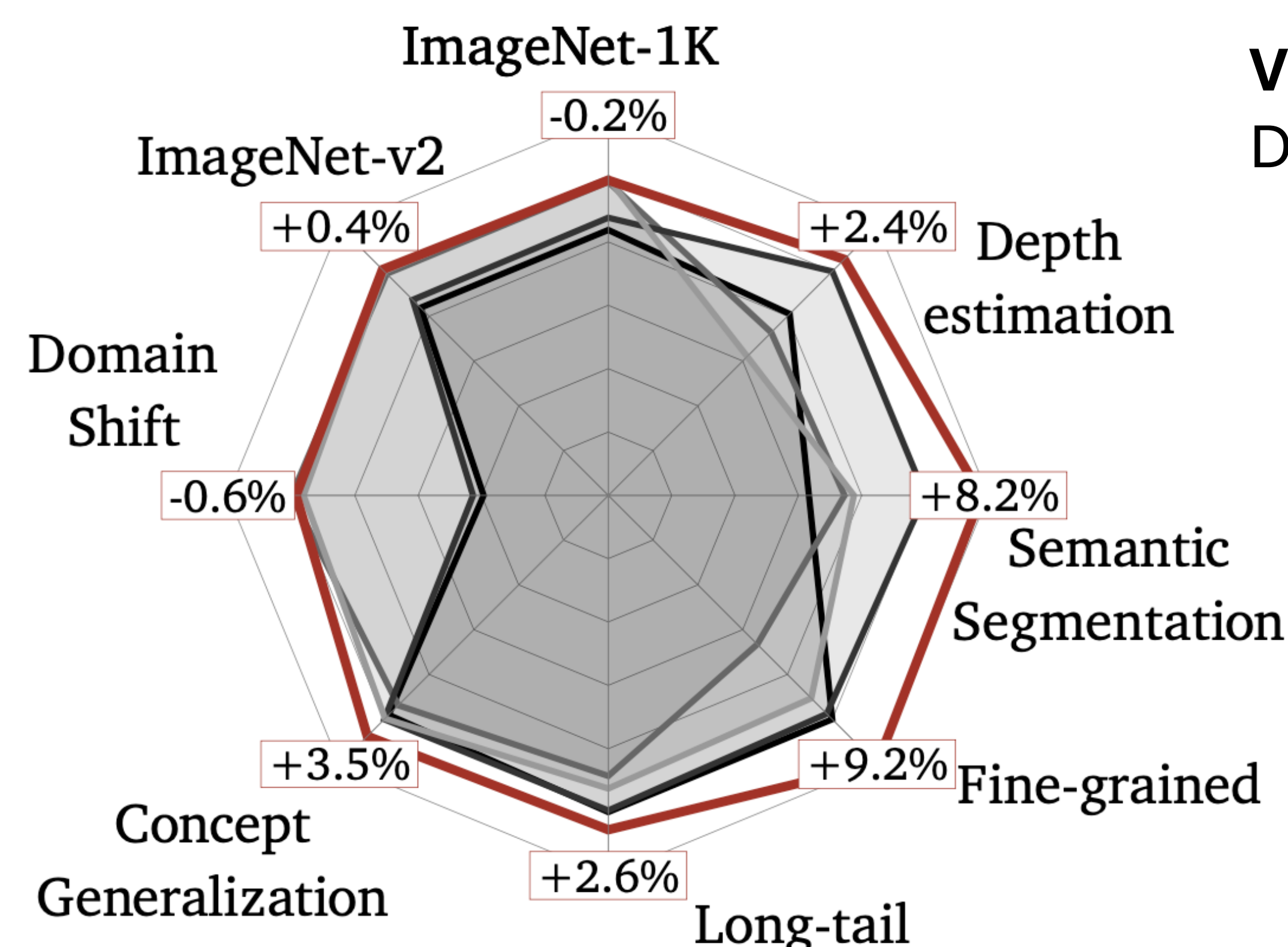
Teacher dropping analysis



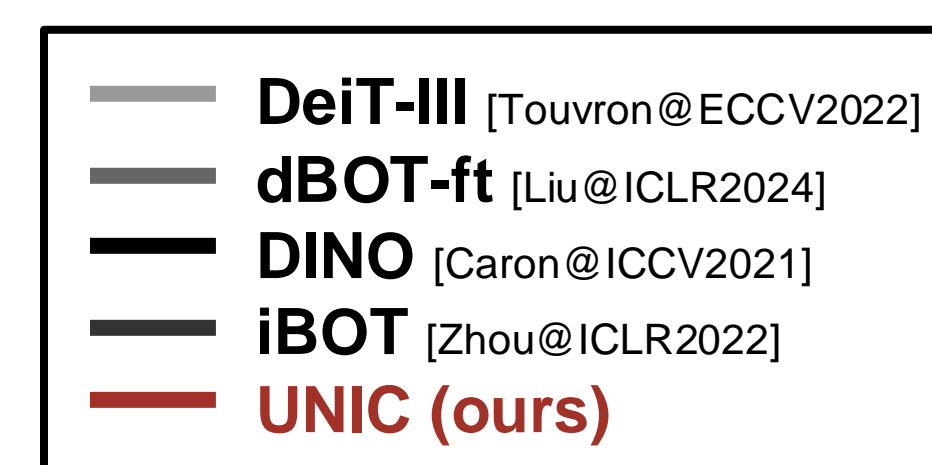
- ✓ Distillation loss is balanced across teachers
- ✓ More discriminative features are learned, albeit slower
- ✓ Teacher importance is dynamically determined

✓ Teacher dropping is effective at **learning all teachers equally well**

Teachers trained on ImageNet-1K



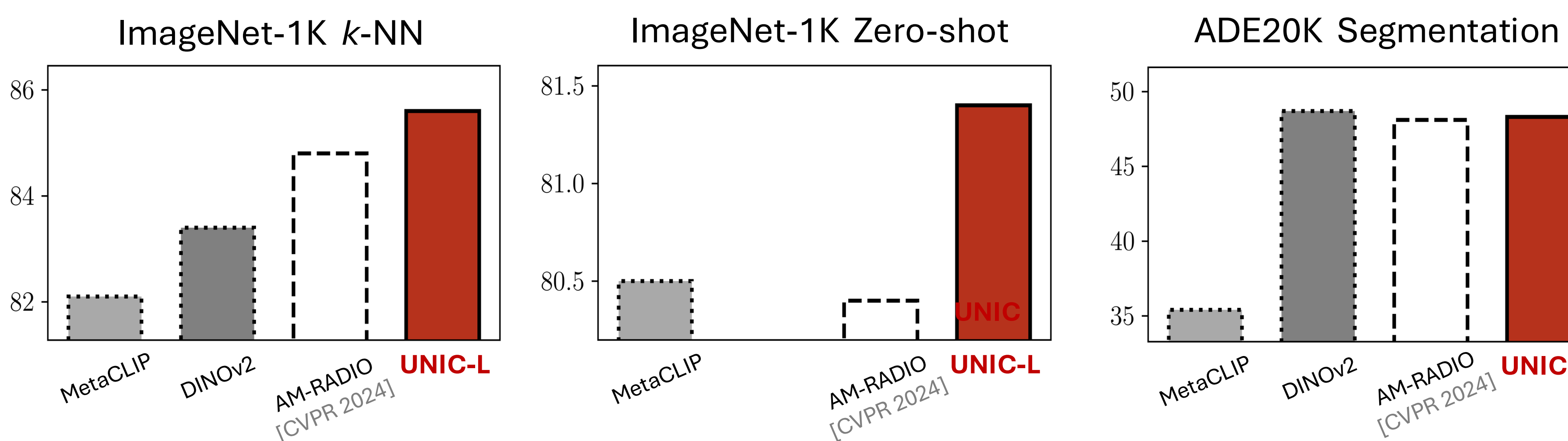
ViT-Base architecture for all Distillation on **ImageNet-1K**



✓ UNIC is on par or better than the best teacher for each task

Comparison to the state of the art

Teachers: MetaCLIP-Huge and DINOv2-Giant, Student: ViT-Large Distillation on **ImageNet-1K**

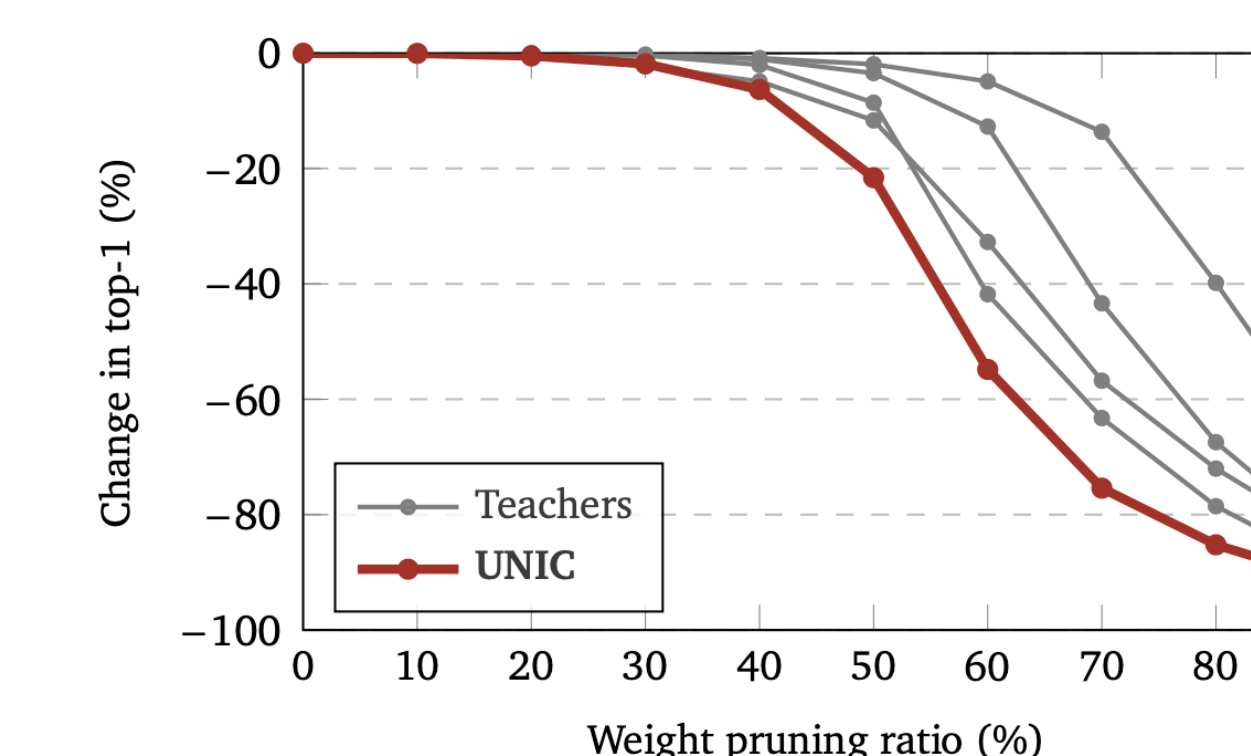


✓ UNIC-L outperforms the best teacher in the vast majority of cases

Model weight and feature space utilization analysis

Weight utilization study

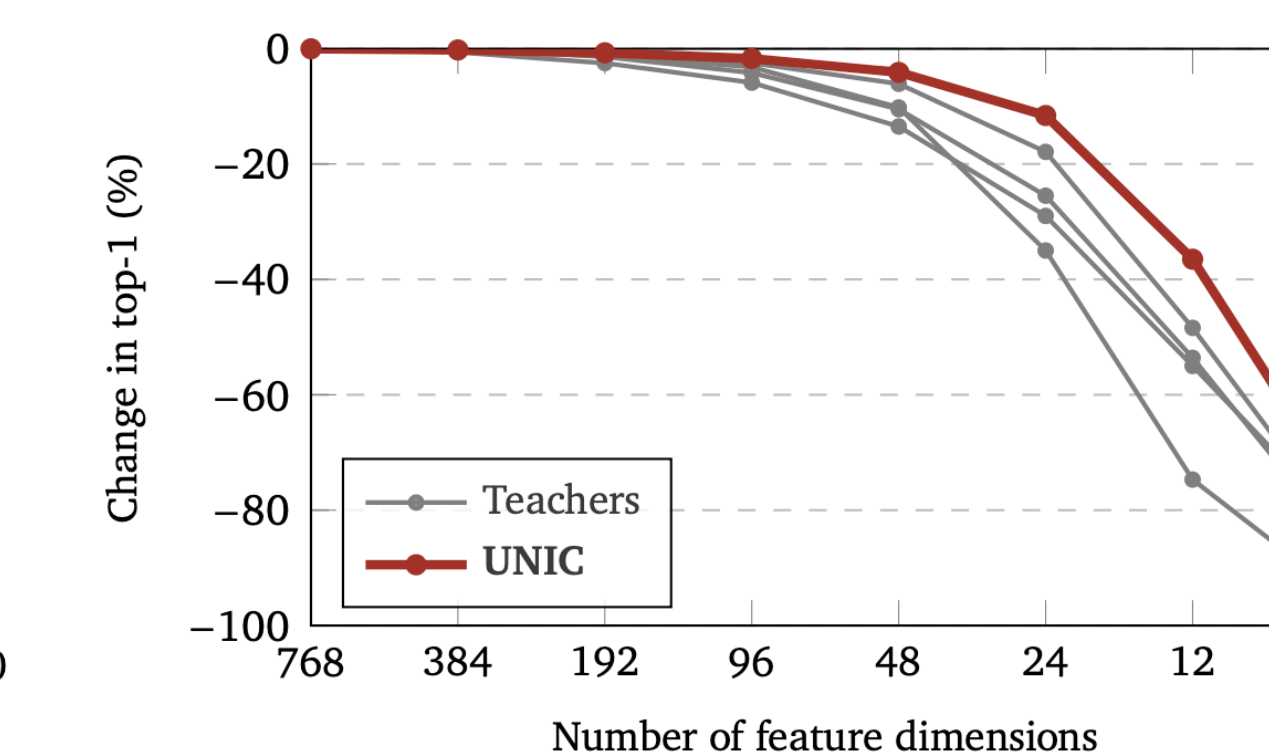
- Prune model weights
- Extract ImageNet-1K features
- Train a linear classifier



✓ UNIC encoders utilize model weights more effectively

Feature space utilization study

- Extract ImageNet-1K features
- Apply PCA
- Train a linear classifier



✓ UNIC features are more resilient to dimensionality reduction