

Good Teachers Explain: Explanation-enhanced Knowledge Distillation

Summary

Knowledge Distillation (KD) is effective for student accuracy.

But

- teacher-student agreement can be low.
- teacher's reasoning might not be distilled.
- teacher's interpretability might get lost.

Agreement := $T: \checkmark \checkmark$

By simply matching explanations (e²KD)

- teacher-student agreement and accuracy improve.
- students remain right for right reasons.
- students remain interpretable.

Experimental Setup

Explanation Methods

Teacher → **Student**





 $RN34 \rightarrow RN18$ $RN50 \rightarrow [RN18, ConvNext, EfficientNet, MobileNet]$ B-cos RN34 \rightarrow B-cos RN18

B-cos DN169 \rightarrow [B-cos ViT_{Tiny}, B-cos RN18]

Datasets ImageNet, Waterbirds, Pascal VOC, SUN397

References Agreement (*Stanton et al., 2021*); B-cos (*Böhle et al., 2021*); GradCAM (*Selvaraju et al., 2017*); RRR (Ross et. al 2017.); Guided Teachers (Rao et al., 2023); EPG (Wang et al., 2020); Waterbirds (Sagawa et al., 2019)

Amin Parchami-Araghi*, Moritz Böhle*, Sukrut Rao*, Bernt Schiele



FOR INFORMATICS

$$\log \sigma_j \left(\frac{z^S}{\tau}\right)$$

), **Explain**
$$(S, x, \hat{y}_T))$$

$IMN \rightarrow SUN$	Acc.	Agr.
Teacher	60.5	_
Using SUN	57.7	67.9
KD	53.5	65.0
+ e ² KD	54.9	67.7

Remain right for right reasons



MAX PLANCK INSTITUTE



SIC Saarland Informatics Campus

EUROPEAN CONFERENCE ON COMPUTER VISION

Similar Predictions, but for Similar Reasons!



Maintain Teacher's Interpretability

Learn to localize (VOC Dataset)

G Score		EPG	IoU	F1
	Teacher	75.7	21.3	72.5
	Baseline	50.0	29.0	58.0
	KD	60.1	31.6	60.1
	$+ \mathbf{e}^2 \mathbf{K} \mathbf{D}$	71.1	24.8	67.6

earn inductive biases (→ViT)
Method	Acc.	Agr.
T: B-cos DenseNet-169	75.2	_
B: B-cos ViT_{Tiny}	60.0	64.6
KD	64.8	70.1
$+ e^2 KD$	66.3	71.8