

Good Teachers Explain: **Explanation-enhanced Knowledge Distillation**

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









Summary

Knowledge Distillation (KD) is effective for student accuracy.

But

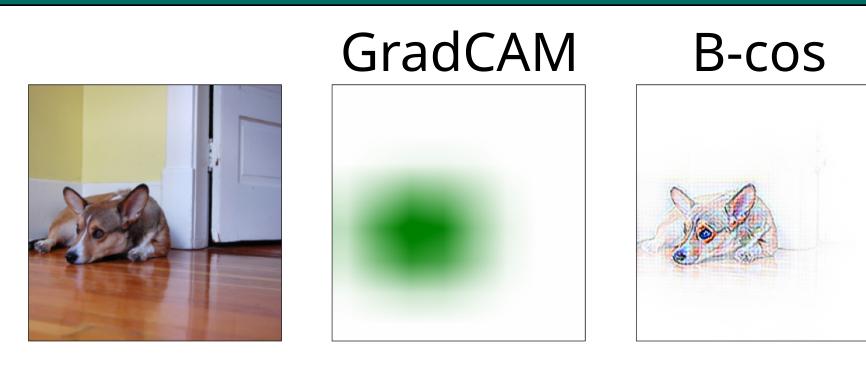
- teacher-student agreement can be low.
- teacher's reasoning might not be distilled.
- teacher's interpretability might get lost.

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 → RN18

RN50 → [RN18, ConvNext, EfficientNet, MobileNet]

B-cos RN34 → B-cos RN18

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

Datasets ImageNet, Waterbirds, Pascal VOC, SUN397

e²KD: Explanation-enhanced KD

In addition to logits

$$\mathcal{L}_{KD} = - au^2 \sum_{j=1}^{C} \sigma_j \left(rac{z^T}{ au}
ight) \log \sigma_j \left(rac{z^S}{ au}
ight)$$

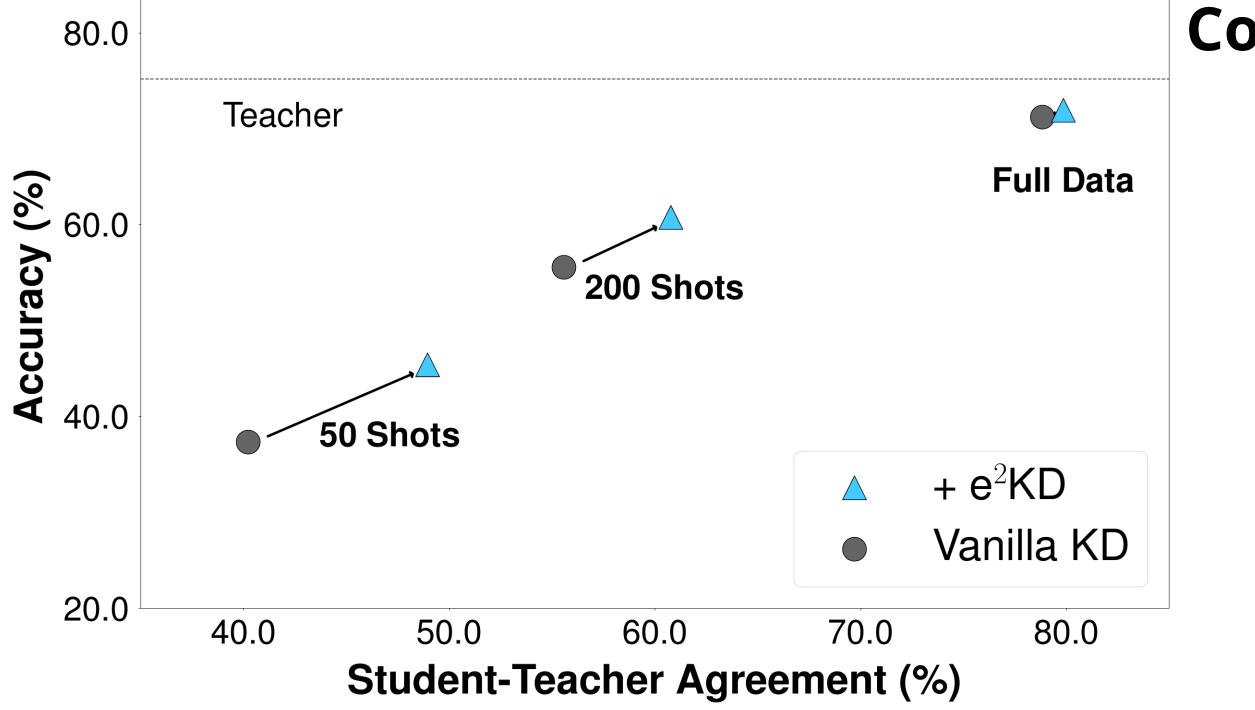
also match explanations:

$$\mathcal{L}_{exp} = 1 - \cos(\mathbf{Explain}(T, x, \hat{y}_T), \mathbf{Explain}(S, x, \hat{y}_T))$$

- Leverages existing explanation methods
- **Model-agnostic Parameter-free**

We discuss three desiderata for faithful distillation

High Agreement with the Teacher



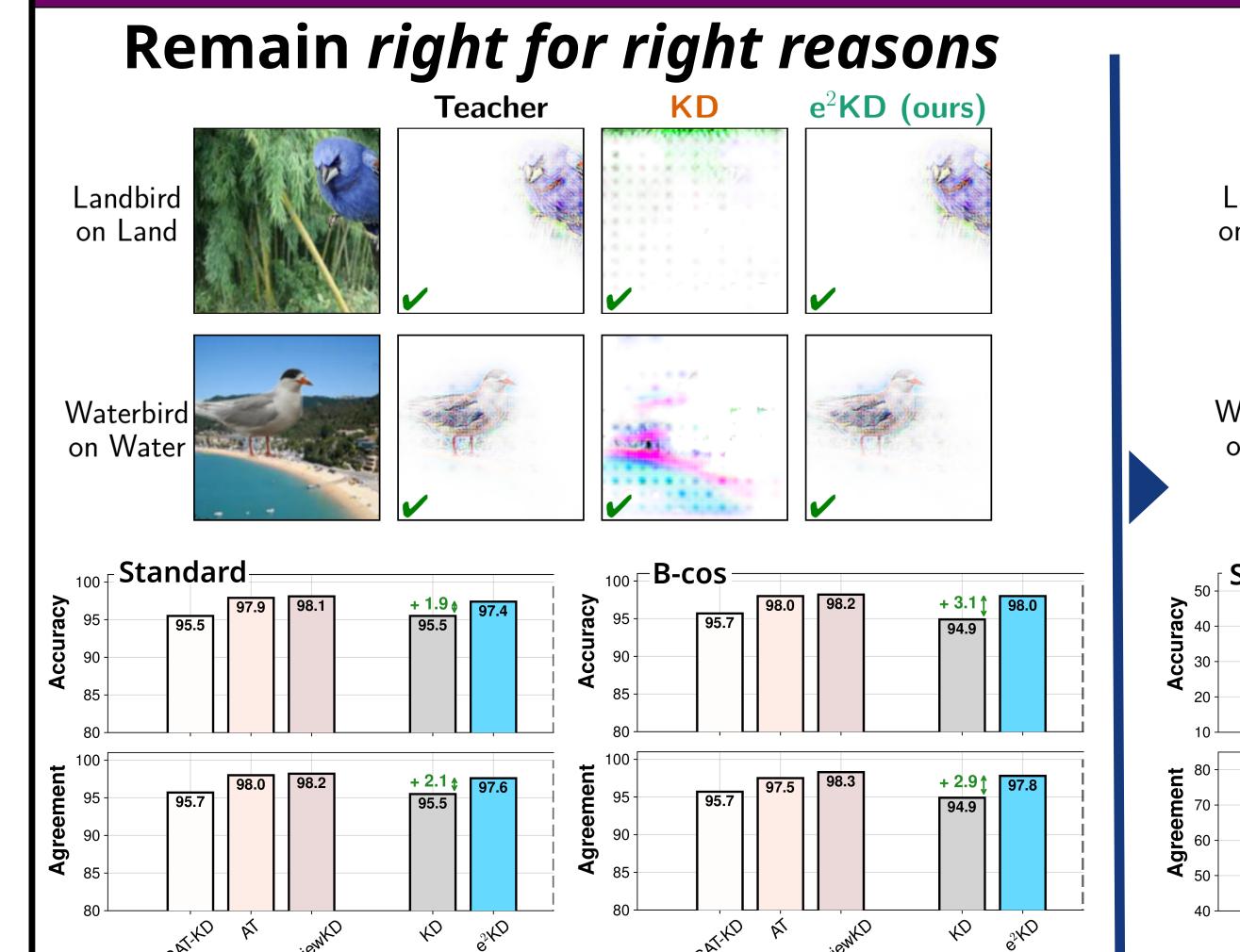
Consistent acc. and agreement gains

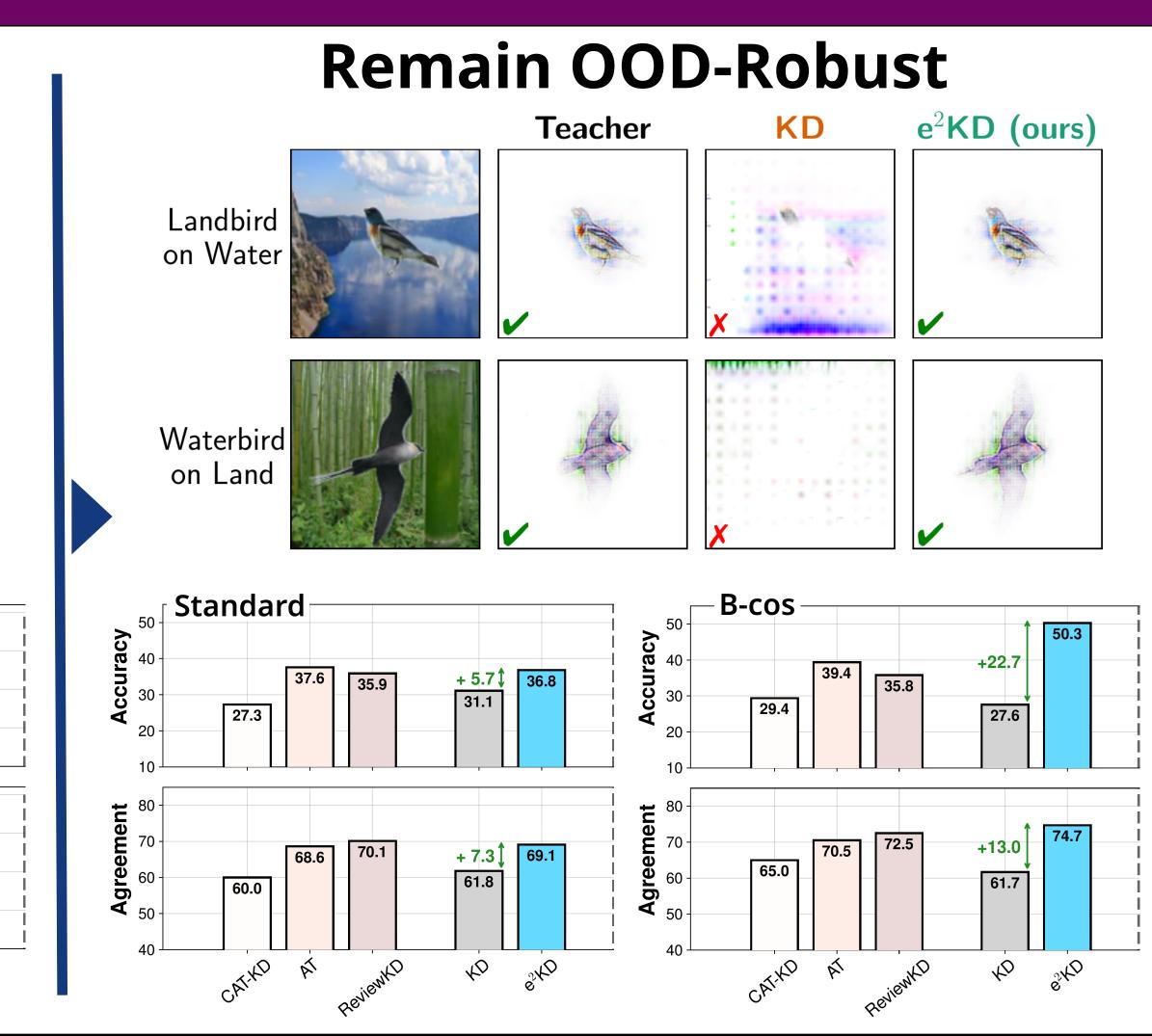
- especially on limited-data
- even on unrelated images

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

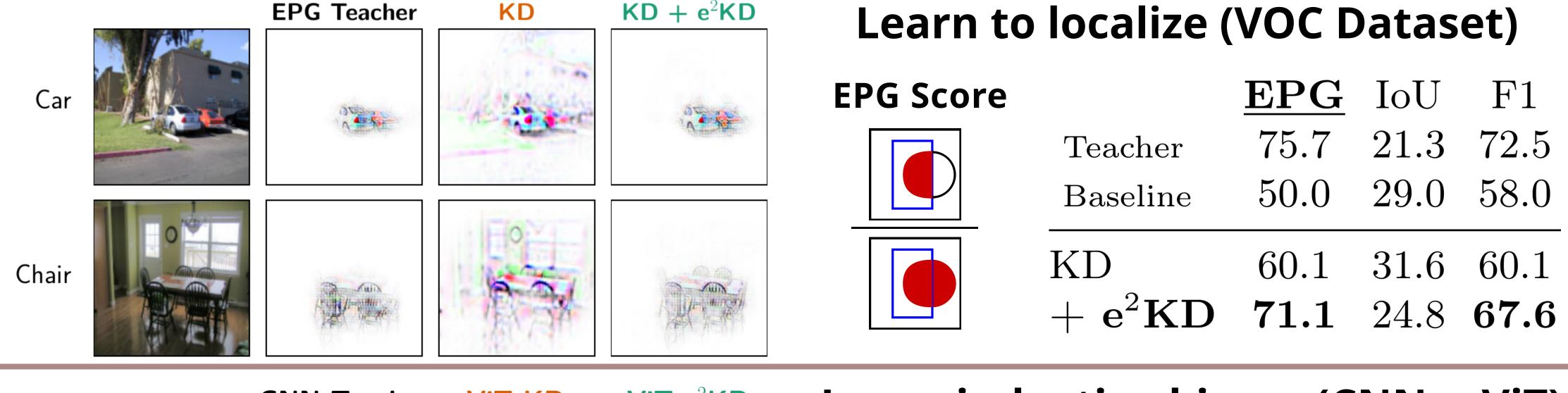
Bonus: 攀e²KD

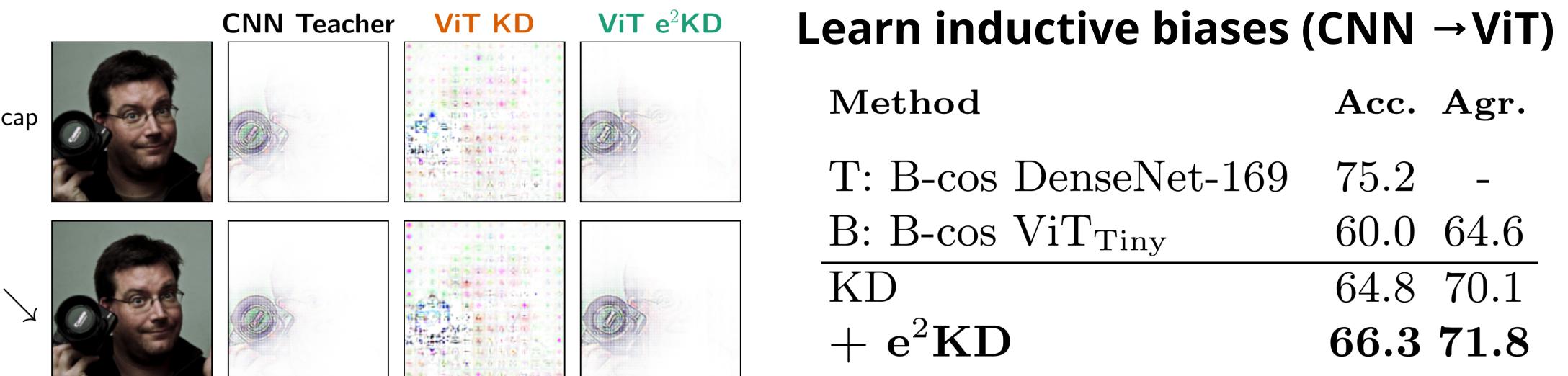
Similar Predictions, but for Similar Reasons!





Maintain Teacher's Interpretability





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)