

Diversity Matters When Learning From Ensembles

Giung Nam^{1*} Jongmin Yoon^{1*} Yoonho Lee^{2,3} Juho Lee^{1,2}

(*Equal contribution)

¹KAIST, South Korea

²AITRICS, South Korea

³Stanford University, USA

NeurIPS 2021

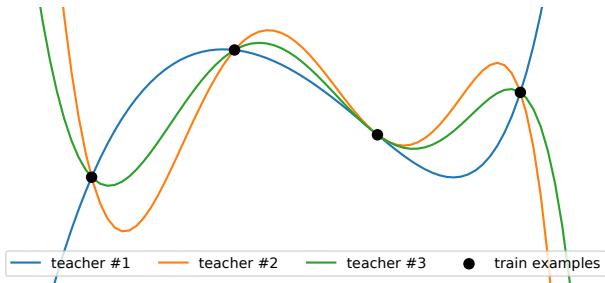
Learning From Ensembles

- ▶ Deep Ensemble (DE):

$$p(y|\mathbf{x}) = \frac{1}{M} \sum_{m=1}^M p_{\theta_m}(y|\mathbf{x}, \theta_m).$$

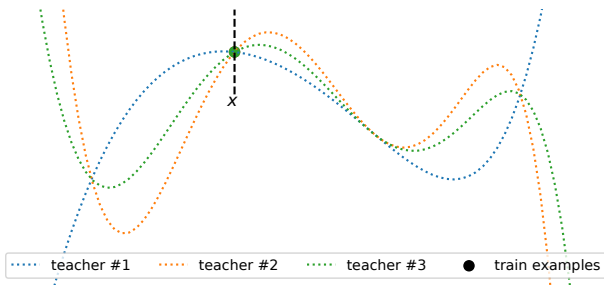
- ▶ How to learn from DE? → Knowledge Distillation (KD).

Motivating Example



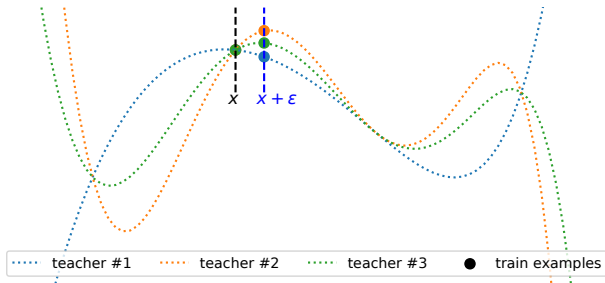
- ▶ Assume that we have three pre-trained teacher models.

Motivating Example



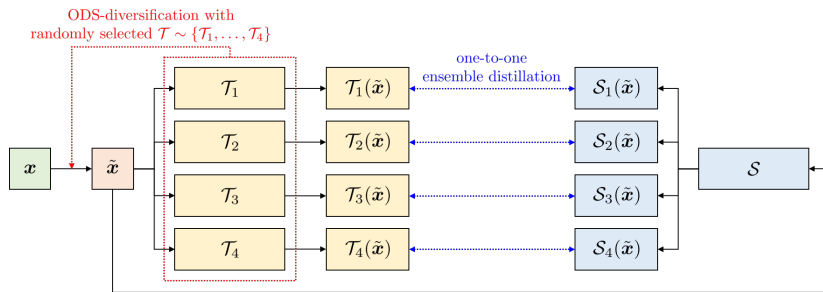
- ▶ Assume that we have three pre-trained teacher models.
- ▶ We cannot capture *diversity* if we reuse the train data x .

Motivating Example



- ▶ Assume that we have three pre-trained teacher models.
- ▶ We cannot capture *diversity* if we reuse the train data x .
- ▶ Idea: use perturbed data $x + \epsilon$ that can capture *diversity*!

Our Approach



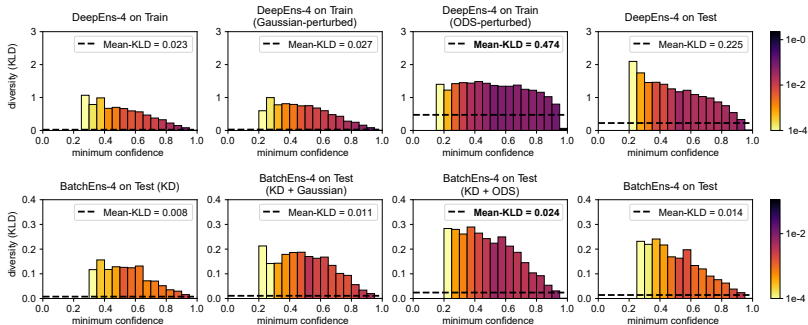
- ▶ Output diversified sampling (ODS):

$$\epsilon_{\text{ODS}}(\mathbf{x}, \mathcal{F}, \mathbf{w}) = \frac{\nabla_{\mathbf{x}}(\mathbf{w}^\top \mathcal{F}(\mathbf{x}))}{\|\nabla_{\mathbf{x}}(\mathbf{w}^\top \mathcal{F}(\mathbf{x}))\|_2} \in \mathbb{R}^D,$$

where $\mathbf{w} \sim \text{Unif}([-1, 1])^K$.

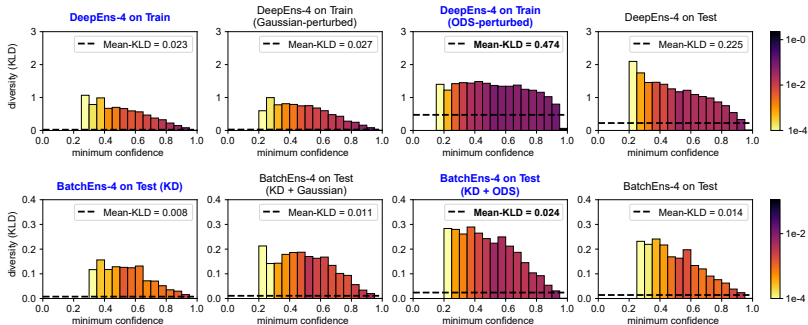
Experimental Results

Impact of ODS on diversities



Experimental Results

Impact of ODS on diversities



Evaluation Results

Image classification tasks

ResNet-32 on CIFAR-10						
Method	# Params	ACC	NLL	BS	ECE	DEE
\mathcal{T} : DeepEns-4	1.86 M	94.42	0.167	0.082	0.008	-
\mathcal{S} : BatchEns-4	0.47 M	93.37 \pm 0.11	0.204 \pm 0.002	0.099 \pm 0.001	0.008 \pm 0.001	1.419 \pm 0.075
+ KD		93.98 \pm 0.20	0.188 \pm 0.003	0.091 \pm 0.002	0.009 \pm 0.002	2.019 \pm 0.174
+ KD + Gaussian		93.93 \pm 0.12	0.187 \pm 0.001	0.090 \pm 0.001	0.009 \pm 0.002	2.042 \pm 0.089
+ KD + ODS		93.89 \pm 0.10	0.181 \pm 0.002	0.090 \pm 0.001	0.006 \pm 0.001	2.486 \pm 0.164
+ KD + ConfODS		94.01 \pm 0.19	0.180 \pm 0.001	0.089 \pm 0.001	0.007 \pm 0.001	2.524 \pm 0.080
\mathcal{T} : DeepEns-8	3.71 M	94.78	0.157	0.077	0.005	-
\mathcal{S} : BatchEns-8	0.48 M	93.47 \pm 0.14	0.202 \pm 0.005	0.098 \pm 0.002	0.006 \pm 0.002	1.494 \pm 0.156
+ KD		94.15 \pm 0.13	0.182 \pm 0.001	0.088 \pm 0.001	0.010 \pm 0.001	2.391 \pm 0.113
+ KD + Gaussian		94.09 \pm 0.08	0.184 \pm 0.002	0.089 \pm 0.001	0.010 \pm 0.002	2.206 \pm 0.175
+ KD + ODS		94.13 \pm 0.08	0.175 \pm 0.003	0.086 \pm 0.001	0.006 \pm 0.002	2.991 \pm 0.322
+ KD + ConfODS		94.18 \pm 0.12	0.174 \pm 0.002	0.086 \pm 0.001	0.007 \pm 0.001	3.064 \pm 0.246

WRN28x10 on CIFAR-100						
Method	# Params	ACC	NLL	BS	ECE	DEE
\mathcal{T} : DeepEns-4	146.15 M	82.52	0.661	0.247	0.022	-
\mathcal{S} : BatchEns-4	36.62 M	80.34 \pm 0.08	0.755 \pm 0.007	0.280 \pm 0.002	0.027 \pm 0.001	1.449 \pm 0.085
+ KD		80.51 \pm 0.22	0.744 \pm 0.003	0.274 \pm 0.001	0.021 \pm 0.003	1.582 \pm 0.035
+ KD + Gaussian		80.39 \pm 0.12	0.761 \pm 0.006	0.277 \pm 0.000	0.022 \pm 0.003	1.379 \pm 0.072
+ KD + ODS		81.88 \pm 0.32	0.674 \pm 0.016	0.257 \pm 0.006	0.026 \pm 0.003	3.303 \pm 0.769
+ KD + ConfODS		81.85 \pm 0.32	0.672 \pm 0.010	0.256 \pm 0.003	0.024 \pm 0.000	3.333 \pm 0.491

Evaluation Results

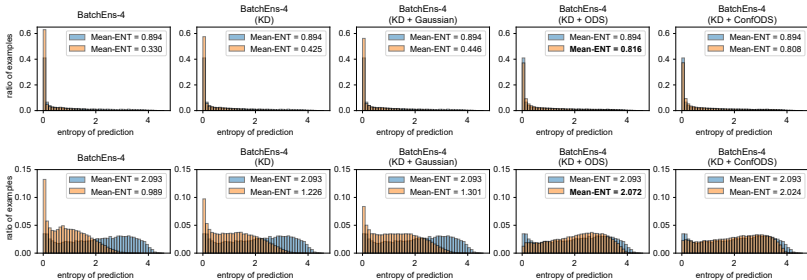
Cross-architecture knowledge distillation for compression

WRN28x5 on CIFAR-100					
Method	# Params	ACC	NLL	BS	ECE
\mathcal{T} : DeepEns-4 (WRN28x10)	146.15 M	82.52	0.661	0.247	0.022
\mathcal{S} : BatchEns-4	9.20 M	78.75 \pm 0.11	0.801 \pm 0.012	0.297 \pm 0.003	0.021 \pm 0.002
+ KD		78.89 \pm 0.10	0.804 \pm 0.012	0.296 \pm 0.003	0.022 \pm 0.001
+ KD + Gaussian		78.80 \pm 0.41	0.815 \pm 0.009	0.297 \pm 0.005	0.020 \pm 0.002
+ KD + ODS		80.24 \pm 0.05	0.742 \pm 0.008	0.279 \pm 0.002	0.028 \pm 0.004
+ KD + ConfODS		80.62 \pm 0.25	0.733 \pm 0.007	0.275 \pm 0.003	0.027 \pm 0.001

WRN28x5 on TinyImageNet					
Method	# Params	ACC	NLL	BS	ECE
\mathcal{T} : DeepEns-4 (WRN28x10)	146.40 M	69.90	1.242	0.403	0.016
\mathcal{S} : BatchEns-4	9.23 M	64.86	1.455	0.464	0.022
+ KD		65.86	1.432	0.456	0.022
+ KD + Gaussian		65.72	1.446	0.457	0.022
+ KD + ODS		65.98	1.408	0.456	0.022

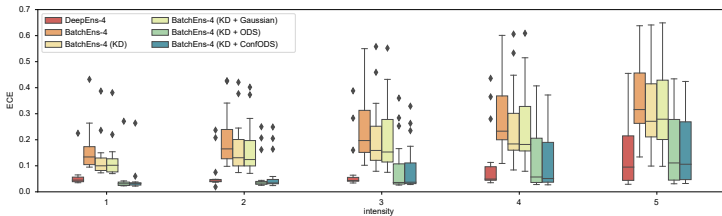
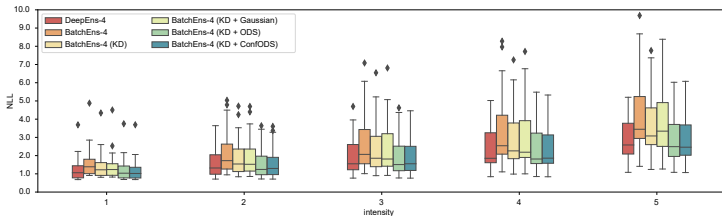
Evaluation Results

Predictive uncertainty for out-of-distribution examples



Evaluation Results

Calibration on corrupted dataset



Summary

- ▶ Diversity matters when learning from ensembles!
- ▶ ODS is one way to capture the *diversity* in KD framework.
- ▶ Such enhanced *diversities* help to reduce the performance gap between student and teacher in KD framework.