

Table 1: Additional experiments in terms of classification Accuracy (%).

(a) Results of WRN28-10@CIFAR and ResNet101@ImageNet.						(b) Parameters sensitivity to different datasets and noise rates.						
Noise Rate	CIFAR10		CIFAR100		ImageNet	α	CIFAR10 (80% Noise)			CIFAR100 (40% Noise)		
	0.4	0.8	0.4	0.8	0.0		0.8	0.9	0.95	0.8	0.9	0.95
ERM	75.41	30.00	55.68	13.99	78.2	Fix $E_s=60$	75.60	78.58	75.44	70.36	71.38	68.57
SELF[24]	93.34	67.41	72.48	42.06	-	E_s	40	60	80	40	60	80
Ours	93.23	80.13	72.60	44.17	78.7	Fix $\alpha=0.9$	68.27	78.58	78.65	70.30	71.38	67.32

1 We thank the reviewers for their valuable comments. We summarize major concerns from reviewers and respond to
 2 them appropriately as follows. We will add suggested experiments, references, and fix typos in the updated version.

3 **To Reviewer #1:** *Q1:* Insights on setting the hyper-parameter E_s . *A1:* The
 4 optimal value of E_s is related to the begining epoch of overfitting. Since the
 5 baseline TRADES uses a step learning-rate-decay schedule and we observe that it
 6 starts to overfit around the epoch of the first learning rate decay (the 75-th epoch,
 7 see the blue curve in Figure 1), we simply set the E_s to a slightly smaller value
 8 70. However, we also find that our method is not sensitive to E_s : using $E_s = 60$
 9 has similar performance as $E_s = 70$ (see Figure 1, the red and green curves).
 10 This phenomenon is consistent with Table 2b in the main body, where we show
 11 our method is not sensitive to various hyper-parameters (see also Table 1b).

12 *Q2:* Extra experiments on other model structures. *A2:* Besides ResNet34@CIFAR10/100 and ResNet50@ImageNet
 13 in main body, we also report the results of WRN28-10@CIFAR10/100 and ResNet101@ImageNet in Table 1a, where
 14 our approach outperforms ERM and the reported numbers of SELF[24] in most entries, sometimes by a large margin.

15 **To Reviewer #2:** *Q1:* Claim of applicability to any deep supervised learning task; validity of premise in new tasks.
 16 *A1:* We point out a potential misunderstanding about our paper, where we did NOT claim that our method is applicable
 17 to *any* deep supervised learning task. Though not for any task, our method is robustly and reliably effective for *a wide*
 18 *range of tasks:* classification with label noise, selective classification, adversarial learning, vanishing double descent,
 19 etc. As a new task, we run experiments on OOD generalization task (see Table 2 and A4@R#3), where ours has better
 20 performance. The premise that the model can guess the right predictions follows the observation that deep models fit the
 21 clean samples first, which is justified by [14,18,32] and our extensive experiments in the above-mentioned broad tasks.

22 *Q2:* Extreme failure cases. *A2:* In Figure 3 of main body, ERM performs better only in the extreme case where
 23 the model capacity is more than $10\times$ smaller than standard ResNet-18, where ERM’s test accuracy is poor ($\leq 78\%$,
 24 i.e., significant underfitting occurs). For other extreme cases when the data is complex or when there are few training
 25 samples, though our method might fail, ERM will perform poorly too, due to the information-theoretical limit. We
 26 argue that studying models with enough capacity, realistic (amount of) input data, and reasonable performance (e.g.,
 27 $\geq 90\%$ test accuracy) might be of more interests to the community, where our method consistently outperforms ERM.

28 **To Reviewer #3:** *Q1:* Cost of maintaining probability vectors. *A1:* The cost is
 29 not high. Take the large-scale ImageNet as an example. The storage of such vectors
 30 in single precision format for the entire dataset requires $1.2 \times 10^6 \times 1000 \times 32$ bit
 31 ≈ 4.47 GB, which is acceptable since modern GPUs usually have no less than 11GB
 32 memory. Moreover, the vectors can be stored on CPU memory or even disk and
 33 loaded along with the images to further reduce the cost.

34 *Q2:* Robustness to other corruption percentages and datasets. *A2:* In Table 1b, we
 35 conduct extra experiments on two datasets with different noise rates. The results indicates that our approach is robust to
 36 the values of hyper-parameters in various settings.

37 *Q3:* Additional large-scale experiments. *A3:* Please refer to Table 1a and A2@R#1 for additional results on ImageNet.

38 *Q4:* OOD generalization. *A4:* In Table 2, we report the average accuracy using ResNet-34 on CIFAR10-C over 15
 39 corruptions. Under various corruption levels, our method consistently outperforms ERM by a considerable margin,
 40 indicating that self-adaptive training provides implicit regularization for OOD generalization.

41 **To Reviewer #4:** *Q1:* Results on uncorrupted CIFAR. *A1:* On CIFAR10/100,
 42 the test accuracy is 95.32%/78.42% for ERM, and 95.17%/78.69% for ours.

43 *Q2:* Contradiction between L30 and L103. *A2:* We will make it clear: in
 44 the first few iterations, though the model learns to fit the correct labels in a
 45 *progressive* manner (as in L30), its predictions are very unstable, especially in
 46 the very beginning of the training procedure (as in L103). The instability is due
 47 to the use of regularization such as data augmentation (as in L104).

48 *Q3:* Further investigation of sample weights. *A3:* The minimum value of w_i
 49 is not bounded by α due to the moving-average scheme that accumulates the
 50 predictions. Following the procedure in Figure 5 of main body, we display the
 51 average sample weights in Figure 2, where the white areas represent the case that
 52 no sample lies in the cell. We see that the weights on the diagonal are higher.

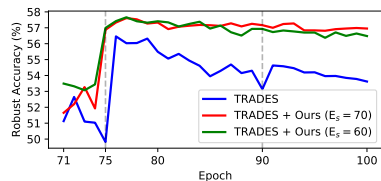


Figure 1: Sensitivity of E_s in the adversarial learning.

Table 2: Average Accuracy (%) on CIFAR10-C at various corruption levels.

Level	1	3	5
ERM	88.44	77.26	58.91
Ours	89.41	78.83	60.77

Figure 2: Average sample weights w_i under various labels.

