

1 We appreciate the reviewers’ time and suggestions! We address them all and report new experimental results below.
 2 *Reviewer 1:* • **Missing citations:...** We will cite the suggested citations and discuss their differences/relations with
 3 our method. Although DIH can be helpful to identify noisy data in noisy-label setting (ref.Middle plot in Fig. 1),
 4 our curriculum is not specifically designed for noisy labeled data and this paper mainly focuses on clean data setting.
 5 DIHCL still achieves 90.34% test-set accuracy under 40% symmetric label noise on CIFAR10 (ref.Top plot in Fig. 1).

6 • **The statement may be revised that “updating in-**
 7 **stantaneous hardness typically requires extra in-**
 8 **ference steps of a model over all the samples” ...**
 9 **extra GPU memory...** SPADE Alg in [ref2] sam-
 10 ples a mini-batch and then selects samples within
 11 the mini-batch using the MentorNet. Comparing
 12 to DIHCL, SPADE incurs the following extra com-
 13 putational costs: (1) the feature extraction on the
 14 sampled mini-batch; (2) training of the MentorNet;
 15 (3) repeated training on well-learned and easy sam-
 16 ples with correct labels. In addition, DIHCL can be
 17 incorporated with [ref2] as a predefined curriculum
 18 to train the MentorNet more efficiently. DIHCL
 19 requires extra GPU memory linear in the number of
 20 samples, which is negligible compared to the
 21 GPU memory for network training. We will add a
 22 discussion of the memory cost in the next version.

23 • **Is the method specific to cyclic learning rate...** DI-
 24 HCL is applicable to other learning rate schedules.
 25 We report the result of DIHCL with a piecewise exponential decay learning rate in Fig. 1. DIHCL improves the test
 26 accuracy from 95.72% to 96.04% in this case. We also visualized the dynamics of samples partitioned by DIH at
 27 epoch-40 and it shows that the properties in Section 2 also hold for different learning rate schedules.

28 • **Clarity** We will simplify the experiment part and use better names for the DIH variants.

29 *Reviewer 2:* • **Compare more baseline methods...** Given the limited time for rebuttal, we compared DIHCL with [2]-[4]
 30 in Table 1 on CIFAR100 when used to train WideResNet-28-10. We will add a complete comparison in the next version.

31 • **“Increasing the subset size (previous studies)” vs “reducing the subset size (as proposed)”** It
 32 depends on one’s preference for easy v.s. hard samples. If the curriculum always selects the
 33 easiest samples (e.g., SPL), the former should be used since only an increasing size can include
 34 harder but more informative samples in later stages. We reduce the subset size in DIHCL for
 35 two reasons: (1) in early stages, training on sufficient samples yields an accurate estimate of
 36 DIH since it is a time-moving average; (2) in later stages, we can reduce unnecessary training
 37 costs on easy/memorized samples, as suggested by observations in Section 2. In Fig. 1, we provide a comparison
 38 between the two and it shows that the former is slightly worse than the latter on the final test accuracy.

39 • **Need variance for in Table 1.** Given limited time, we report the
 40 mean±variance over 5 trials on two datasets in Table 2. We will add complete
 41 variance results in the next version.

42 • **SOTA results with all techniques.** We will try to compete with the baselines
 43 in the SOTA setting in the next version.

44 *Reviewer 3:* • **Effects of T_0 and γ_k ...** T_0 is necessary to get an early estimate
 45 of DIH and is supported by the theoretical analysis in the Appendix. But a
 46 small $T_0 = 3$ suffices to get a stable training because: after the T_0 warm-starting epochs, we start from all samples
 47 (exploration) and gradually reduce the subset size. In Fig. 1, we compare $T_0 = 5$ and $T_0 = 3$: they produce similar test
 48 accuracy. We also compare $\gamma_k = 0.85$ and $\gamma_k = 0.8$: reducing the subset size too fast (γ_k is too small) will degrade the
 49 accuracy.

50 *Reviewer 4:* • **T_0 and hyper-parameter tuning.** Please also see our reply to Reviewer 3. DIHCL is not sensitive to the
 51 choice of T_0 since the earlier epochs after warm starting select almost all samples (exploration), which keeps the DIH’s
 52 estimation accurate. We do not have computation for a full grid search of all hyperparameters so there could exist better
 53 choices. Line 307-312 details how we selected the hyperparameters.

54 • **Method sensitive to noise.** The noisy-label experiments in the Appendix use a 100% noise setting, i.e., all the labels
 55 are randomized and wrong. The purpose is to show that the pattern of DIH is very different on clean and noisy data, so
 56 we can use DIH as an indicator of label noise. In Fig. 1, we report the results under 40% symmetric label noise. It
 57 shows (1) DIH contains critical information to identify the noisy-labeled data, and (2) DIHCL is robust to label noise
 58 and achieves 90.34% test-set accuracy under 40% symmetric label noise on CIFAR10.

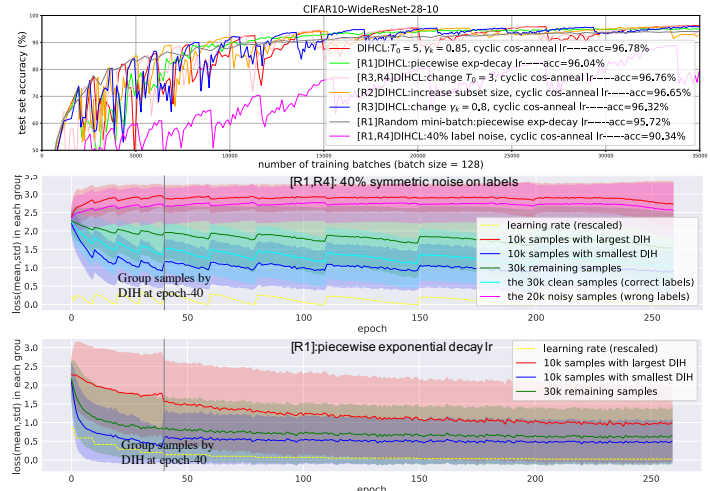


Figure 1: Top: variants of DIHCL-Rand-dLoss; Middle & Bottom: dynamics of DIH-grouped samples under label noise and another lr schedule.

Table 1: [R2] Test accuracy (%) of WideResNet-28-10 on CIFAR100.

Method	Test Accuracy (%)
DIHCL	82.23
[2]	75.04
[3]	71.95
[4]	76.28

Table 2: [R3] Test Acc (mean±variance).

Curriculum	CIFAR10	CIFAR100
DIHCL-Rand, Loss	96.74 ± 0.04	80.80 ± 0.16
DIHCL-Rand, dLoss	96.75 ± 0.06	80.73 ± 0.21
DIHCL-Exp, Loss	97.07 ± 0.11	82.31 ± 0.24
DIHCL-Exp, dLoss	96.44 ± 0.10	81.35 ± 0.27
DIHCL-Beta, Flip	96.48 ± 0.04	81.13 ± 0.18