

1 Dear all reviewers: Thank you very much for taking your time to review our paper and providing us valuable and
 2 insightful comments. Below, we answer all of the questions.

3

4 **R1 Q1)** It wasn't totally clear to me how you ensure $T + \Delta T$ remains a valid transition matrix throughout training.
 5 Also, is ΔT initialized to zero matrix, or randomly?

6 A1) We ensure $T + \Delta T$ to be a valid transition matrix by first projecting their negative entries to be zero and then
 7 performing row normalization. About ΔT , we initialize it to be a zero matrix in our experiments.

8 **R1 Q2)** Can this approach take advantage of a small clean set?

9 A2) Yes. If a small clean set is available, it is helpful. It can be used to (1) better initialize the transition matrix, (2)
 10 better validate the slack variable ΔT , and to (3) fine-tune the deep network.

11 **R1 Q3)** **It would be interesting to also include results using a subset of the WebVision dataset to see if the method**
 12 **works there too.**

13 A3) **Due to the limited time, we do experiments on a subset of WebVision dataset.** Specifically, we create the
 14 training data by sampling a hundred classes from the thousand classes and sampling 1,600 images for each class. 10%
 15 of the training data is held out for validation. We use the total of 5,000 images from the sampled 100 classes in the
 16 original validation set as the test set. We compare "Forward", "Reweight", "Forward-R", "Reweight-R" on this subset.
 17 We use a ResNet-50 model pretrained on Imagenet. The results are 80.48% (Forward), 81.08% (Reweight), 85.12%
 18 (Forward-R), 85.42% (Reweight-R). We can see that "Reweight" and "Reweight-R" achieve better results and the
 19 revision technique greatly boost the performance. Due to the limited time, we only compared with "Forward" and
 20 "Reweighting" (in our setting, those two methods consistently work better than other baselines).

21 **R1 Q4)** Discuss the relationship with "Multiclass learning with partially corrupted labels"(Wang et al)

22 A4) Thank you for the valuable feedback. They both employ the importance reweighting technique. However, their
 23 approach requires anchor points to estimate the transition matrix; while the proposed approach is designed to release
 24 the strong requirement of anchor points.

25 **R2 Q5)** More experiments on real data.

26 A5) We perform our experiments on the subset (mentioned in A3) of WebVision using four models. The results of
 27 experiments are 80.48% (Forward), 81.08% (Reweight), 85.12% (Forward-R), 85.42% (Reweight-R). More details can
 28 be found in A3.

29 **R3 Q6)** Please discuss the model performance on MNIST with more label noise.

30 A6) We raise the noise rates to be 60%, 70%, 80%. Other experiment settings are unchanged. The results are presented
 31 in Tables 1 and 2. We can see that the proposed model outperforms the baselines more significantly as the noise rate
 32 grows. Due to the limited time, we only compared with "Forward" and "Reweighting" (in our setting, those two methods
 33 consistently work better than other baselines).

34 **R3 Q7)** Discuss the potential extension of the proposed approach.

35 A7) We can extend our approach to mixture proportion estimation [27] and learning with complementary label.

Table 1: Mean test accuracy (in %, \pm std dev).

| | Sym-60% | Sym-70% | Sym-80% |
|--------------|----------------------------------|----------------------------------|----------------------------------|
| Forward-A | 97.10 \pm 0.08 | 96.06 \pm 0.41 | 91.46 \pm 1.03 |
| Forward-A-R | 97.65 \pm 0.11 | 96.42 \pm 0.35 | 91.77 \pm 0.22 |
| Reweight-A | 97.39 \pm 0.27 | 96.25 \pm 0.26 | 93.79 \pm 0.52 |
| Reweight-A-R | 97.83\pm0.18 | 97.13\pm0.08 | 94.19\pm0.45 |

36 Table 2: Mean test accuracy (in %, \pm std dev), anchor points removed.

| | Sym-60% | Sym-70% | Sym-80% |
|----------------|----------------------------------|----------------------------------|----------------------------------|
| Forward-N/A | 96.82 \pm 0.14 | 94.61 \pm 0.28 | 85.95 \pm 1.01 |
| Forward-N/A-R | 96.99 \pm 0.16 | 95.02 \pm 0.17 | 86.04 \pm 1.03 |
| Reweight-N/A | 97.01 \pm 0.20 | 95.94 \pm 0.14 | 91.59 \pm 0.70 |
| Reweight-N/A-R | 97.81\pm0.12 | 96.59\pm0.15 | 91.91\pm0.65 |

- "-A" means the transition matrix is estimated by using the instance X with highest estimated $P(\bar{Y}|X)$ (which are likely to be anchor points).
- "-N/A" means instances with high estimated $P(Y|X)$ are removed from the dataset.
- "-R" means that the transition matrix used is revised by a revision ΔT .
- The highest accuracy in each column is bold faced.