

1 We sincerely thank all reviewers for their constructive comments and positive support. In particular, we respectively ask  
 2 R3 to reconsider his/her evaluation since we can address all the comments, as detailed below.

3 **Q1 (R1):** Will the estimation error for  $T^{\clubsuit}$  be zero if the anchor points are estimated, if so why? And ablation study.

4 **A1:** There will be an estimation error for  $T^{\clubsuit}$  if there is an estimation error for the anchor points. The same estimation  
 5 error for the anchor points will also go for the  $T$ -estimator. For simplicity, we study the setting where anchor points are  
 6 given. In this case, the error only comes from estimating  $T^{\clubsuit}$ . Thus we do not need to do ablation study for the errors  
 7 of the two matrices  $T^{\clubsuit}$  and  $T^{\spadesuit}$ . Note that even given anchor points, there is an estimation error for the  $T$ -estimator  
 8 because of estimating the noisy class posterior while there is no estimation error for  $T^{\clubsuit}$  because the intermediate class  
 9 posterior is designed to be given. Considering how the estimation error for anchor points affects the two estimators is  
 10 also very interesting. We leave it as a future work.

11 **Q2 (R1):** Will intermediate class labels obtained ... have an estimation error that will affect the error of both matrices?

12 **A2:** The intermediate class labels are defined by the ones that maximize the intermediate class posteriors. Since we  
 13 have given the true intermediate class posterior distribution, i.e.,  $P(Y'|x) = \hat{P}(\bar{Y}|x)$ , there will be no estimation error  
 14 for the intermediate class labels.

15 **Q3 (R1):** On the synthetic dataset, dual T-estimator performs better ... Is there an explanation on why this happens?

16 **A3:** On the real-world datasets, the estimation error of the dual  $T$ -estimator can be larger than that of the  $T$ -estimator  
 17 only with small training sample sizes. It is because the number of images per class is too small to estimate the transition  
 18 matrix  $\hat{T}^{\spadesuit}$  which can be sparse and lead to a large estimation error. We will discuss this in the final version.

19 **Q4 (R3):** The justification and explanation of Eq. 3 is not clear.  $\hat{P}(\bar{Y}|Y')$  (which is computed using the validation set).

20 **A4:** We will add more discussions about Eq.3. It discusses how to make  $T^{\spadesuit}$  to be independent of  $Y$  because  $Y$  is unavail-  
 21 able. Specifically, we have explained a sufficient condition for letting Eq. 3 hold, i.e., let the intermediate class labels be  
 22 identical to noisy labels. We have also discussed that the condition may be hard to be satisfied, therefore, an estimation  
 23 error for fitting the noisy labels to intermediate class labels  $\Delta_3$  is introduced to our estimation. In Appendix 2, we have  
 24 also empirically validated the estimation error. In line 193, we have stated that  $\hat{P}(\bar{Y}|Y')$  is NOT computed on the valida-  
 25 tion set. It is computed on the training set. The validation set functions as the same role for both the proposed method and  
 26 the baselines for model selection. Your concerns on the unfair training/validation split may not stand. However, we agree  
 27 that to see how the size of validation set influences the performance is interesting. We have redone experiments according  
 28 to your suggestions to reduce the size of validation set, i.e., using 10% of the training examples as a validation set for the  
 29 both estimators with five repeated trials (random seeds are from 1 to 5). The estimation errors for transition matrices are  
 30 illustrated in the following table, which show that the dual  $T$ -estimator still significantly outperforms the  $T$ -estimator.

	MNIST			F-MNIST			CIFAR-10			CIFAR-100		
	Sym-20%	Sym-50%	Pair-45%	Sym-20%	Sym-50%	Pair-45%	Sym-20%	Pair-45%	Sym-20%	Pair-45%	Sym-20%	Pair-45%
$T$	0.333±0.002	0.459±0.001	0.644±0.069	0.261±0.012	0.428±0.019	0.687±0.097	0.383±0.004	0.772±0.077	0.861±0.028	0.389±0.001	0.760±0.117	0.894±0.002
Dual $T$	<b>0.106±0.003</b>	<b>0.271±0.007</b>	<b>0.298±0.021</b>	<b>0.220±0.014</b>	<b>0.315±0.012</b>	<b>0.414±0.020</b>	<b>0.215±0.011</b>	<b>0.267±0.075</b>	<b>0.578±0.026</b>	<b>0.255±0.020</b>	<b>0.516±0.187</b>	<b>0.835±0.014</b>

32 **Q5 (R3):** Why use T-Revision models of poorer classification accuracy in the Dual T-estimator paper?

33 **A5:** This could be a common concern if we focus on improving the classification accuracy. However, **our aim is to**  
 34 **reduce the estimation error for the transition matrix. The experiments are set to FAIRLY verify the effectiveness of the**  
 35 **Dual  $T$ -estimator.** As mentioned by the reviewer, the settings of our experiments are different from the original paper,  
 36 thus the reported accuracy is different. Specifically, to boost the classification performance, different tricks has been  
 37 employed in the baselines, e.g., transition information is given (e.g., Co-Teaching); choosing different hyper-parameters  
 38 to estimate anchor points (e.g., using the 97% and 100% largest estimated noisy class posteriors for *CIFAR-10* and  
 39 *CIFAR-100*, respectively, in Forward); training networks with different epochs to estimate noisy class posteriors (e.g.,  
 40 40 epochs in Forward while 20 epochs in T-revision on *mnist*); with or without using the validation set to estimate noisy  
 41 class posteriors (e.g., T-revision uses a validation set while Forward does not). To fairly show the superiority of Dual  
 42  $T$ , all the baselines in our paper use the same tricks, e.g., using a noisy validation set to choose model for estimating  
 43 the noisy class posterior and using the largest noisy class posteriors for estimating anchor points. So the classification  
 44 performance would be worse than that in the original paper. However, the effectiveness of the Dual  $T$ -estimator has  
 45 been successfully illustrated. We will add the discussions in the final version.

46 **Q6 (R4):** There is little discussion about the impact of (1) in the paper. e.g. will different types of the classifiers, which  
 47 have different capabilities to fit the noisy labels, influence the performance of the estimator and how large?

48 **A6:** Thanks for the insightful question. We will add more discussions. For example, the estimation error for  $T$ -estimator  
 49 is from estimating the noisy class posterior in (1); the estimator error for Dual  $T$ -estimator (or  $T^{\spadesuit}$  because  $T^{\clubsuit}$  is  
 50 noise free) is also related to (1) because to push the intermediate class  $Y'$  to be close to the noisy class  $\bar{Y}$ . Studying  
 51 how different classifiers influence the performance of the estimator via impacting (1) is interesting. In our paper, we  
 52 assume that  $\Delta_1$ , which is the estimation error for the noisy class posterior, is bigger than  $\Delta_3$ , which is the estimation  
 53 error for pushing the intermediate class  $Y'$  to be equal to the noisy class  $\bar{Y}$ . Note that if  $\Delta_1 \geq \Delta_2 + \Delta_3$ , the dual  
 54  $T$ -estimator will outperform  $T$ -estimator, where  $\Delta_2 = |P(\bar{Y}|Y') - \hat{P}(\bar{Y}|Y')|$  is the error introduced by counting the  
 55 noisy and intermediate class labels. We discussed when this will hold and empirically validated this assumption in the  
 56 supplementary material. We agree with the reviewer that it is interesting to study when the assumption is invalid and  
 57 how different types of classifiers influence it. We would add some empirical study in the supplementary material.