

1 **To Reviewer 1:** Thank you for your positive comments. We have investigated the effect of the disagreement parameter
2 C in Section 4.3. Actually, we also did similar experiments on other teacher-student networks and ensemble sizes (20
3 and 25). The observations are also similar. We will include these results to improve the comprehensiveness of our
4 ablation studies. Also, we will cover more references to enrich the related work.

5 **To Reviewer 2:** Thank you for your detailed comments and suggestions.

6 **Re: motivation.** You may misunderstand our motivation. Our AE-KD actually **does not** treat every teacher in the
7 ensemble equally. On the contrary, we introduce a parameter C to allow disagreement among teachers, so that the
8 obtained update direction of parameters is not necessarily a strict descent direction for all teachers. As you suggested,
9 the gradient directions from weak or noisy teachers are not as reliable as those from good teachers (also see line
10 161-166). With the help of parameter C , the final direction will not accommodate these weak teachers necessarily, and
11 a better weight over teachers can be automatically determined by solving Problem (9) or (11).

12 **Re: comparison with OKDDip.** The main differences lie in three ways. First, the learning paradigms are different.
13 OKDDip is essentially an online KD paradigm which adapts two-level distillation. While AE-KD follows traditional
14 teacher-student paradigm and the teachers won't be updated during the training. Next, we use different strategies to
15 learn the weights. In the second-level distillation of OKDDip, group knowledge is transferred to the group leader,
16 where all diverse peers serve as a group of teachers. The weights for diverse peers are computed using self-attention
17 mechanism. In AE-KD, the weights for teachers are computed based on multi-objective optimization in gradient space.
18 We take disagreement among teachers into consideration and prevent the student suffering from adverse guidance. Last,
19 AE-KD has the tunable parameter of disagreement to reconcile all teachers while OKDDip doesn't.

20 **Re: optimization of weights α_m in Eq. (11).** Eq. (11) is solved for every minibatch with the calculated gradients over
21 all teacher losses. It is a typical One-class SVM problem with constraint $0 \leq \alpha_m \leq C$, and can be easily solved by
22 LIBSVM or other off-the-shelf solvers.

23 **Re: 350 epochs for resnet20.** In experiments, we train all the teachers for standard 240 epochs, following the setting
24 in [1]. And for student resnet20, we train for 350 epochs for better performance.

25 **Re: ensemble networks with various architectures.** Our AE-KD performs in the same way, *i.e.*, every teacher
26 provides a gradient and the final direction is computed according to Eq.(9). To validate this case, we take resnet20
27 as student, and use resnet10, wide_resnet_40_2 and vgg13 as teacher networks with accuracy 74.31%, 75.61% and
28 74.64%, respectively. Results show that our AE-KD can achieve 70.16% accuracy while our baseline method AVER
29 only has 69.40%, showing our superiority of dealing ensemble KD even with various architectures. We will include this
30 setting in our final version.

31 **To Reviewer 4:** Thank you for your detailed comments and constructive suggestions for our experiments.

32 **Re: Section 2.** It intended to formally illustrate related work and baseline methods. We will consider your advice.

33 **Re: issues about experiments.**

34 **1: multiple runs.** Thanks for your suggestion. We experimented with five resnet56 teacher networks and the resnet20
35 student network on CIFAR10, and run our AE-KD for 10 times with different random seeds. The mean and standard
36 variance are 92.49% and 0.02%, respectively. We can see the performance of student network tends to be steady (low
37 variance) due to the distillation from teacher networks. We will cover this in our final version.

38 **2: performance gap between AE-KD* and AE-KD on CIFAR10/100.** The performance gap comes from the weight
39 β in Eq.(7) of the feature-based loss in AE-KD*. In real implementation, for CIFAR10 we determine the optimal β
40 by cross-validation in $\{10^{-1}, 1, 10, 100, 1000\}$, and for simplicity, we solely adopt the same β for CIFAR100, thus the
41 performance of AE-KD* on CIFAR100 can drop a bit than AE-KD since we do not tune β specially for CIFAR100.
42 However, what we emphasize here is that given a KD method (logits or feature based), when it comes to the ensemble
43 setting, our method can develop a better way to distill from their ensemble by adaptively assigning weights. Of course,
44 suitable parameters bring in better results. For example, if we tune β specially for CIFAR100, our AE-KD* can have
45 71.95% accuracy on resnet20 student network with five resnet56 teacher networks, surpassing AE-KD (71.37%).

46 **3: experiments with more teachers.** With more teachers involved, the proportion of weak teachers will tend to be
47 steady. Thus for AVER, their influence on the performance will gradually reduce to a certain level, and the accuracy
48 will increase dramatically with more teachers at first, but stabilize eventually as Figure 2. However, our method remains
49 steady and robust for most time. We experiment with 30, 40, 50 teachers on CIFAR10. The performances of AVER and
50 AE-KD are 92.59%, 92.67%, 92.69% and 92.77%, 92.82%, 92.86%, respectively. AE-KD still outperforms AVER with
51 a comparable and stable gap for large ensemble size.

52 **4: weak teacher models.** Our teacher models follow those in [1] with the same or similar accuracy. Admittedly,
53 stronger teacher models do produce better students. We choose network from [2] on CIFAR10 as teacher models with
54 97.37% accuracy and the performance of AE-KD on resnet20 student network reaches 95.14% compared to AVER
55 (93.66%). The results are even better and AE-KD retains its superiority.

56 **5: line 244.** Yes. It's a typo. Thank you for pointing it out and we will fix it in our final version.

57 [1] Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive representation distillation. ICLR 2020.

58 [2] Hanxiao Liu, Karen Simonyan, and Yiming Yang. Darts: Differentiable architecture search. ICLR 2019.