

Table 1: Comparison with other KD methods on CIFAR100.

Model	Baseline	VID	AT	FSP	Jacobian	SVD	Heo	Ours
ResNet18	77.09	78.93	78.45	78.75	78.45	78.53	79.21	82.92
ResNet50	77.42	79.21	78.73	79.02	79.03	78.82	79.72	84.74
MobileNetV2	69.04	70.62	70.34	70.48	70.26	69.35	71.02	73.57

Table 2: Comparison experiments on ImageNet.

Model	Baseline	KD	FitNet	Self-Distill	DML	Ours
ResNet18	69.76	70.45	70.26	70.51	70.39	70.92
MobileNetV2	71.52	72.23	71.95	72.37	72.29	72.81
ShuffleNetV2	69.36	70.14	69.86	70.26	70.21	70.74

1 We sincerely appreciate all reviewers for their most detailed and thoughtful reviews.

2 **Overall Response: Comparison with more KD methods (R1’Question-3 and R4’Question-2).** Tab.1 shows the
 3 comparison with additional 6 kinds of KD methods, including the VID (Ahn *et al.*, CVPR2019) requested by R1, and
 4 5 kinds of KD methods (from Tab.1 of the paper) requested by R4. It is observed that our method outperforms the
 5 second-best KD method by 3.71%, 5.02%, 2.55% on ResNet18, ResNet50, MobileNetV2 @ CIFAR100, respectively.

6 **Response to R1: Question-1 and Question-2.** Tab.2 shows the comparison experiments on ImageNet with ResNet18,
 7 MobileNetV2 and ShuffleNetV2. It is observed that our method leads to 1.16%, 1.29%, 1.38% accuracy improvements
 8 on ImageNet with ResNet18, MobileNetV2, ShuffleNetV2 respectively, outperforming the second-best KD method
 9 by 0.41%, 0.44%, 0.48%, respectively. **Question-3.** We have re-conducted experiments on ModelNet10/40 with
 10 ResGCN8 and each experiment is repeated 5 times. On ModelNet10, the accuracy of our method and DML are
 11 $94.35 \pm 0.21\%$ and $93.63 \pm 0.17\%$, respectively. On ModelNet40, the the accuracy of our method and DML are
 12 $91.78 \pm 0.12\%$ and $91.67 \pm 0.09\%$, respectively. Note that DML is a logit-based KD method while the orthogonal loss
 13 can only be utilized in feature-based KD methods so DML can not be improved by the orthogonal loss. **Question-4.**
 14 Please refer to the overall response (the VID Column). **Question-5.** We do not introduce new hyper-parameters for
 15 balancing these two losses because we find that the fixed 1:1 loss ratio is good enough. Additional experiments show
 16 that the performance of our method can be improved by 0.15% by introducing and adjusting a new hyper-parameter
 17 here (ResNet18, CIFAR100). **Question-6.** The proposed method has given strong experimental clue to the intuition
 18 that *task-oriented features are more essential in the knowledge transfer related domains*. These improvements might
 19 open up new research opportunities for exploring the transfer-based ML including mathematical theory and algorithm.
 20 **Question-7.** Thanks for your advice, we will remove this sentence in the final version.

21 **Response to R4: Reply to additional feedback.** The training loss for KD baselines includes both distillation loss
 22 and the supervised loss. Thanks for the valuable question and we will modify our writing to clarify this confusion.
 23 **Question-1.** We apply both supervised loss and distillation loss to the auxiliary classifiers instead of only supervised
 24 loss because the distillation loss can facilitate the training of auxiliary classifiers. A more accurate auxiliary classifier
 25 is able to extract better task-oriented features, which improves the performance of knowledge distillation in turn. *How*
 26 *to Avoid the Case?* The case that auxiliary classifiers have learned to minimize their loss but the backbone model fails
 27 to learn useful features will not happen because (i) Compared with the whole backbone, the auxiliary classifiers have
 28 much fewer parameters so they have weaker learning ability than the backbone model. (ii) The auxiliary classifiers
 29 are trained to minimize both distillation loss and supervised loss, which is harder than only the supervised loss. (iii) In
 30 the overall loss function, we multiply the loss of auxiliary classifiers by hyper-parameters whose value is less than 1
 31 ($\alpha = 0.03, \beta = 0.5$), which forces the whole model to pay more attention to the learning of the backbone and the final
 32 classifier, instead of the auxiliary classifiers. **Question-2.** With the same or fewer parameters in the training period,
 33 our method still outperforms the other KD methods. For instance, ResNet18 with our method has 12.10M and 11.22M
 34 parameters in the training and testing period respectively, achieving 82.92% accuracy in the testing period. In contrast,
 35 ResNet50 with Hinto’s KD has 23.70M parameters in both the training and testing period, achieving 78.58% accuracy
 36 in the testing period. Note that the parameters of auxiliary classifiers are very tiny (only 7% of the whole model, on
 37 ResNet18). **Question-3.** We have added comparison experiments with them. Please refer to the overall response.

38 **Response to R5:** The proposed method has given strong experimental clue to the intuition that *task-oriented features*
 39 *are more essential in the knowledge transfer related domains*. These improvements might open up new research
 40 opportunities for exploring the transfer-based ML including mathematical theory and algorithm.

41 **Response to R6: Question-1.** This question can be considered in two aspects. First, when the feature resizing layer is
 42 used to increase the dimension of student features, due to the energy-preserving property of orthogonality, the energy
 43 of student features will not be amplified by the feature resizing layer. Second, when the feature resizing layer is used
 44 to decrease the dimension of teacher features, the weight of feature resizing layer works as an orthogonal projection
 45 to the teacher features, which leads to less information loss than non-orthogonal projection. **Question-2.** As shown in
 46 the ablation study (Tab.7), compared with naive feature distillation, the task loss brings 3.50% accuracy improvement,
 47 which accounts for 80% of the total accuracy gain, indicating that the task loss is the key of our method.