

# 1 Overall Response

2 We sincerely thank all the reviewers for their helpful comments and constructive suggestions.

3 **Regarding comments on writing of this paper.** We will carefully revise the manuscript and fix all the missing  
4 citations (R4), missing explanations (R2), confusing sentences (R3 and R4) and organization (R2 and R3).

5 **Regarding the significance of our approach.** The main goal is to propose a framework for dynamic inference with a  
6 common backbone. Although both distillation and attention are well-developed, to the best of our knowledge, it is the  
7 first attempt for combing them successfully with noticeable accuracy and speed gain. These improvements might open  
8 up new research opportunities for exploring the resource efficient ML including mathematical theory and algorithm.  
9 Actually, the motivation of using attention and distillation differs from their origins. More specifically, in proposed  
10 method, attention is utilized to separate features of different classifiers rather than to find the most informative pixels;  
11 and distillation is to transfer knowledge among classifiers in the common backbone, instead of among different models.

## 12 2 Response to Reviewer 2

13 **Regarding the comparison with related work.** Fig. 1 shows the comparison with Feedback Network[C1], Fractal-  
14 Net[C2] and NestedNet[C3] (results taken from their papers). Proposed method shows a large margin improvement of  
15 our proposed method. Compared with [C1,C2,C3], our method is different in two folds: (i) computation of shallow  
16 classifiers in the proposed method can be reused by the deeper classifiers, leading to possibility for dynamic inference;  
17 (ii) the accuracy of the early predictions in these work is lower than their baselines. In contrast, two of three shallow  
18 classifiers in the proposed method achieve higher accuracy than the baseline.

19 **Regarding the specific questions of writing.** (i)  $F_i$  and  $F_c$  denotes features of  $i_{th}$  and the deepest classifier. They are  
20 outputs of the convolutional layers, also named as "activation", which is unlearnable. (ii) The ensemble prediction is  
21 obtained by weighted sum of softmax layers outputs from all the classifiers.

22 **Regarding the sensitivity of hyper-parameters  $\alpha$  and  $\lambda$ .** The proposed method is robust to hyper-parameters. All  
23 the reported accuracy in the paper shares the same hyper-parameters settings:  $\alpha = 0.5$ ,  $\lambda = 5 \times 10^{-7}$ . In addition, we  
24 supplement the accuracy of ResNet18 with various  $\alpha$  and  $\lambda$  on CIFAR100 in Figure 2 and 3. The observed accuracy  
25 varying range (less than 0.5%) is negligible compared with the improvements (3.37%) from baselines (77.09%).

26 **Regarding the results in Table 2.** The main goal of the proposed method is to accelerate the DNN inference with  
27 multi-classifiers in a common backbone. Compared with many related work which achieve acceleration at the expense  
28 of accuracy loss, the 1.26% accuracy increment on ImageNet in Table 2 may not be trivial with acceptable computation  
29 increment (no exceeding 5%, while ResNet50->101 brings another 98% computation for 1.5% accuracy improvement).

## 30 3 Response to Reviewer 3

31 **Regarding the hardware friendliness.** The hardware friendliness of the proposed method lies in the comparison with  
32 scalable neural networks such as SkipNet, as introduced in line69-74. The hardware friendliness of lightweight designs  
33 such as MobileNet, will be destroyed when they are utilized as backbones of SkipNet, because SkipNets added complex  
34 gated control units on every layer. In contrast, the thresholds based control method in SCAN can maintain the hardware  
35 friendliness of backbones for its simpleness.

## 36 4 Response to Reviewer 4

37 **Regarding the ensemble for every level (Q1 and Q2).** We implement your suggestion and the experiments on  
38 ResNet18-CIFAR100 show that every level ensemble leads to 1.64%, 1.48% and 1.67% accuracy increment on the  $2_{nd}$ ,  
39  $3_{rd}$  and  $4_{th}$  depth prediction, respectively. With the improved ensemble classifiers on all the levels, scalable inference  
40 with new thresholds achieves 1.54% accuracy increment, compared with the origin model with same acceleration. We  
41 will further explore this as future works.

42 **Regarding line192-193 (Q3).** Observation (iv): the ensemble accuracy is 1.11% higher than that of the deepest classifier  
43 on average with almost no computation penalty. Observation (v): compared with classifiers trained individually, the  
44 proposed method can bring more accuracy increment on shallow classifiers than the deepest one.

45 **Regarding the missing experiments (Q4).** The details of parameter count and total memory usage will be added in  
46 revision. An average acceleration (2.17 $\times$ ) and compression (3.20 $\times$ ) ratio is given in line190-line191.

47 **Regarding training time (Q5).** Taking a pretrained model, ResNet50 on ImageNet takes 8.75 hours for 25 epoches  
48 trained on two 2080Ti GPU devices.

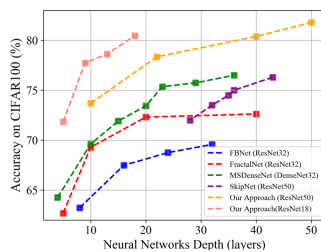


Figure 1: Comparisons.

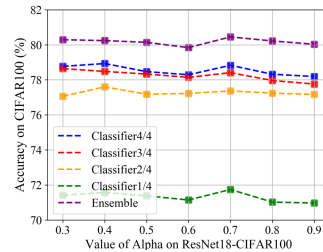


Figure 2: Sensitivity of  $\alpha$ .

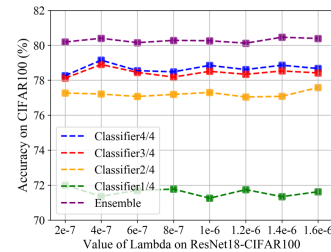


Figure 3: Sensitivity of  $\lambda$ .