

1 We thank all three reviewers (**R#1, R#3, R#4**) for appreciating our novelty, strong results, and broad impact to the  
2 community, and promise to release all codes upon acceptance and to improve writing clarity. Below, we first respond to  
3 general E<sup>2</sup>-Train questions, and then specific questions on **SMD, SLU** and **PSG**, being grouped between lines.

4 **Reducing accuracy loss? (R#1).** We continued striving to reduce the accuracy loss after NeurIPS submission, and  
5 found a stochastic weight averaging (SWA) technique (“SWALP” paper in ICML’19) to be helpful. As requested by R#1,  
6 we report the improved new result : after applying SWA to E<sup>2</sup>-Train on ResNet 74, we obtain top-1 93.01% on CIFAR-10  
7 (0.56% loss) and top-1 71.63% on CIFAR-100 (no accuracy loss) with 83.40% and 81.27% energy saving, respectively.

8 **Experiments on adapting a pre-trained model (R#1, R#3).** We performed a proof-of-concept experiment for CNN  
9 fine-tuning by splitting CIFAR-10 training set into half, where each class was i.i.d. split evenly. We first pre-train  
10 ResNet-74 on the first half, then fine-tune it on the second half. During fine-tuning, we compared two energy-efficient  
11 options: (1) fine-tuning only the last FC layer using standard training; (2) fine-tuning all layers using E<sup>2</sup>-Train. With all  
12 hyperparameters being tuned to best efforts, the two fine-tuning methods improve over the pre-trained model top-1  
13 accuracy by [0.30%, 1.37%] respectively, while (2) saves 61.58% more energy (FPGA-measured) than (1). That shows  
14 that E<sup>2</sup>-Train is the preferred option: higher accuracy, and saving much more energy. We will report it in camera-ready.

15 **On different CNNs (R#3, R#4).** E<sup>2</sup>-Train is effective in even compact networks, e.g., MobileNetV2: 92.06% / 71.61%  
16 top-1 accuracy on CIFAR-10/-100 (baseline: 92.47% / 71.91%), with 87.73% / 88.17% energy savings, respectively.

17 **Training time (R#3).** E<sup>2</sup>-Train takes roughly the same epoch numbers to converge (line 324). Similar to energy, the  
18 per-epoch time is also reduced, thanks to our three-level savings. We will report detailed measurements in camera-ready.

19 **Why SMD works, by theory? (R#3, R#4).** We are happy to confirm that SMD is not simply a heuristic - **we recently**  
20 **proved** that SMD can outperform SGD at certain range of mini-batch sampling ratios and epoch numbers. The proof  
21 is inspired by “*Random Shuffling Beats SGD after Finite Epochs*”, ICML’19, and the core idea relies on the finite  
22 population correction technique. We plan to report our theoretical findings and proof in another upcoming submission.

23 **Is SMD stable/reproducible? (R#3, R#4).** We repeated training ResNet74 on CIFAR-10 using SMD for 10 times (all  
24 end up saving  $\frac{1}{3}$  energy), with different random initializations. The accuracy std. is only 0.132%, showing high stability.

25 **Is SMD generalizable to more models? (R#4).** For ResNet74 on CIFAR-100, [SMD, SMB] with  $\frac{1}{3}$  energy saving  
26 achieve top-1 [71.37%, 71.11%], respectively; and for ResNet-110 on CIFAR-10, [SMD, SMB] with  $\frac{1}{3}$  saving achieve  
27 [93.05%, 92.75%], respectively. Various of such experiment results confirm that SMD is consistently better than SMB.

28 **Clarifying Fig. 3b (R#4).** We tried to find more baselines for solid evaluation of SMD. As SMD outperformed SMB  
29 with standard learning rates (lr), we conjectured increasing lr might accelerate SMB’s convergence (reducing epochs),  
30 thus making another strong baseline. We also repeated SMD/SMB experiments and found their gap to persist (>0.2%).

31 **SLU on non-residual CNNs (R#3, R#4).** We thank both reviewers to point out. Indeed, our current experiments are  
32 based on CNNs with skip connections, partially because they dominate in SOTA CNNs. We will make it clear in  
33 camera-ready. Furthermore, we conjecture SLU can be extended to non-residual CNNs by “appending” skip connections  
34 with gates to plain backbones (i.e., creating “ResNet” versions for their training). We leave it for future work.

35 **Error bar of SLU (R#4).** 20 trials of SLU experiments to ResNet38 on CIFAR-10 show that, with 95% confidence  
36 level, the confidence interval for the mean of the top-1 accuracy and the energy saving are [92.47%, 92.58%] (baseline:  
37 92.50%) and [39.55%, 40.52%], respectively, verifying SLU’s trustworthy effectiveness.

38 **What auxiliary RNNs learn(R#3)/ SLU training dynamics(R#4).** After training, visualizations show RNNs learned  
39 (i) the class-wise discriminative selectivity of layers (ResBlocks); and (ii) input hardness-aware routing, consistent with  
40 SkipNet’s observations. During training, we observe that the auxiliary RNNs converge much faster (<20 epochs) than  
41 the main backbone part, showing “layer selectivity” can be identified at early training stage. Interestingly, this coincides  
42 with observations of “critical learning periods”, ICLR’19. We will include the visualizations/analysis into camera-ready.

43 **Adjusting  $p_L$  (R#4).**  $p_L$  controls SD’s drop ratio, and we always tune it to make SD’s drop ratio the same as SLU (line  
44 291). Similar to SLU, dropping more layers decreases SD’s performance, in exchange for more energy saving.

45 **Clarifying PSG description (R#4).** PSG is motivated by two facts: (i) reducing precision is very effective (expo-  
46 nentially) for reducing hardware energy cost; and (ii) MSB parts contribute exponentially larger than LSB parts for  
47 calculating outputs. Specifically, while SGD uses  $x$  and  $g_y$  to calculate  $g_w$ , PSG uses  $x^{\text{msb}}$  and  $g_y^{\text{msb}}$  to calculate  $g_w^{\text{msb}}$   
48 for predicting the sign of  $g_w$  when  $|g_w|^{\text{msb}} \geq \tau$  (this happens with high probability thanks to fact (ii), e.g., > 60% in  
49 ResNet74 on CIFAR10). Otherwise (i.e.,  $|g_w|^{\text{msb}} < \tau$ ), PSG will proceed to finish the remaining residual computation  
50 of calculating  $g_w$  and its sign. In hardware, the computation to obtain  $g_w^{\text{msb}}$  is embedded within that of  $g_w$ . Therefore,  
51 the PSG’s predictors (i.e., calculating  $g_w^{\text{msb}}$  using  $x^{\text{msb}}$  and  $g_y^{\text{msb}}$ ) do not incur any energy overhead.

52 **More requested details on PSG (R#3).** We use default  $\beta = 0.05$  to control MSBs, while the effectiveness of PSG is  
53 observed to be insensitive to the choice of  $\beta$  when it is in the range of [0.05, 0.1]. If  $\beta$  is too small, it will result in too  
54 frequent coarse gradients and might hurt convergence. The ratio of using merely  $g^{\text{msb}}$  typically remains at least 60% in  
55 the training process. We will include these in camera-ready, along with detailed explanation of  $E_1$  and  $E_2$  in (3).