

1 We thank the reviewers for their feedback. Our paper will be updated to reflect the responses below.

2 **ALL: Novelty?** Our contribution are: (1) Demonstrating training is highly sensitive to sparsifying
3 back-propagated output error gradients; (2) A simple training algorithm, sparse in forward and
4 backward passes, with activation and weight sparsity that fits emerging sparse accelerators. *Why*
5 *SWAT outperforms related work?* We believe this is because SWAT can perform dynamic topology
6 exploration: Backpropagation with sparse weights and activations approximates backpropagation on
7 a network with sparse connectivity and sparsely activated neurons. The gradients generated during
8 back-propagation minimize loss for the current sparse connectivity. However, each iteration SWAT
9 generates a potentially new sparse network using the sparsifying function. Non-active weights are also
10 updated thus capturing fine-grained temporal importance of connectivity during training. Figure 7
11 and 8 show the importance of unmasked gradient updates and dynamic exploration of connectivity.

12 **Reviewer 1: (1)** When sparsifying Equation 2 and 3, we expect combinations sparsifying both
13 gradients and weights and/or activations (i.e., (g,a), (g,w), (g,w,a)) would underperform as sparsifying
14 gradients alone harms convergence. Also, as sparsifying weights or activations alone will not sparsify
15 both equations, we skip these combinations too. **(2)** DSG loses considerable accuracy even at low
16 sparsity. E.g., for ResNet18 on ImageNet at 50% sparsity DSG suffers an accuracy loss of 4.6%.

17 **Reviewer 2: (1)** “*Drastic drop due to sparse activations in forward pass*”: In Figure 1 we isolate the
18 forward and backward pass and examine sensitivity of training to sparsifying only the forward pass.
19 Notably, this means we use the full activation for the backward pass. So, in Figure 1 the backward
20 pass is not optimizing the network for sparse activation. Typical filter pruning methods use sparse
21 activation in the backward pass so the network adapts for activation sparsity. Kurtz et al. (ICML
22 2020) and Georgiadis (CVPR 2019) show how to increase activation sparsity in the forward pass.
23 Combining such techniques with SWAT may reduce FLOPS without increasing error and thus may be
24 a good direction for future work. **(2)** “*Never faced such drastic drops from sparse output gradients*”:
25 We believe there is a misunderstanding: GMP, STR, CS and RigL use a sparse weight gradient during
26 the parameter update stage whereas in Figure 2 we sparsify the back-propagated error gradient. The
27 back-propagated output error gradient is sparsified before performing the convolutions to generate
28 weight and input gradients. The generated weight gradient is dense because the convolution between
29 a sparse activation and dense back-propagated error gradient tensor yields a dense result. This back-
30 propagated error gradient is not dropped in STR, CS, and GMP but rather dense back-propagated
31 error gradients are used to generate weight gradients that are masked during the parameter update
32 stage. Thus, STR, CS, GMP only update the active parameters. **(3)** Yes, sparsity budget proposed
33 by Kusupati et al. (ICML 2020) could be used for more speedup. **(4)** Yes, the different channels
34 will be selected for different filters and there is overhead for introducing structured sparsity since the
35 L1 response of channels is computed. However, the FLOP reduction due to high sparsity is more
36 than the overhead of computing the L1 response. Moreover, the idea of periodically computing the
37 Top-K channel selection is applicable. **(5)** SWAT doesn’t have dense phase. Base model accuracy
38 can be calculated from the accuracy drop. E.g., for ResNet-50 on ImageNet, the Top-1 accuracy is
39 76.8. **(6)** Yes, in general uniform will have lower flop than ERK especially when the input resolution
40 is high. ERK generally applies lower sparsity at initial layer which have significant computation
41 especially true for imagenet. However, in Table 2, ERK is more efficient because the initial layer has
42 small computation due to small input resolution (CIFAR-10) and computationally expensive layers
43 have higher sparsity. Similarly, we observe DST allocates higher sparsity to more computationally
44 expensive layers and quickly reaches a sparsity where overall computation is low. **(7)** We will add
45 discussion on Li et al. 2017 and follow-up work on structured sparsity.

46 **Reviewer 3: (1)** *The speedup gap is small:* The theoretical speedup by the earlier method would be
47 around $\frac{1}{1-0.67} = 3.03\times$ whereas with our method the theoretical speedup would be $\frac{1}{1-0.76} = 4.16\times$.
48 Theoretical speedup does not fully capture the benefit of SWAT since SWAT also reduces memory
49 footprint during training. **(2)** *Methodology for cycle count estimation?* We will add relevant detail
50 from the supplementary material. Simulator counts the cycles taken to spatially map and schedule the
51 computation present in each layer. The memory hierarchy is similar to the DaDianNo architecture.
52 *Table 1:* (i) The column represents whether the input gradient and weight gradient computation is
53 sparse or not. (ii) Second, convolution between sparse input activation and dense back-propagated
54 error gradient tensor generates a dense gradient for weight and the dense gradient is used in parameter
55 update. (iii) Yes, related work can be adapted for structured sparsity.

56 **Reviewer 4:** Selecting the sparsity budget of SWAT (e.g., using ERK) is not our main contribution,
57 but rather demonstrates SWAT can readily leverage such techniques. We have covered a range of
58 configurations in Table 4, 5 6 and 7 in the main paper and Table 1 and 2 of the supplemental. Joint
59 optimization of backward pass activation and weight sparsity is an good direction for future work.