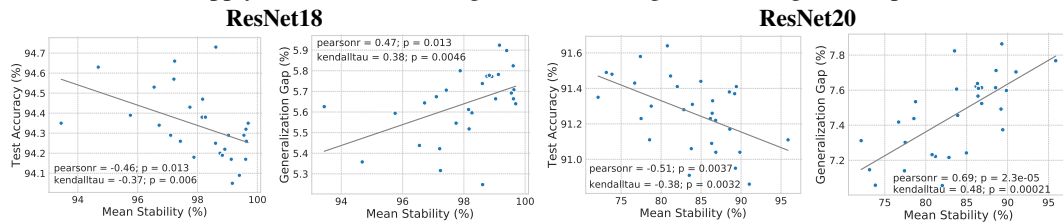


1 We thank the reviewers for their detailed, valuable reviews. We are glad that the reviewers saw that our work was:
 2 “novel and insightful” with “a very intriguing and intuitive explanation for the mechanics of pruning” and “well
 3 designed” experiments (R1); possessing an “interesting” core observation about pruning stability and generalization
 4 and a “[comprehensive] set of research questions”, including a “worthwhile study” of “pruning the largest magnitude
 5 filters” and “characterizations of flatness” (R2); and addressing “an interesting phenomenon” with a “proposed metric”
 6 and approach that “could lead towards more insight in the training dynamics of sparse neural networks” (R3). We
 7 address reviewer comments as much as possible below but will incorporate all feedback in the final version.

8 **R1,R2,R3 Concern about generality of results due to experimental hyperparameters and overparameterized**
 9 **networks:** We agree with the reviewers’ concern and will ensure the final version includes the following data that
 10 show the generalization-stability tradeoff is not merely an artifact of the particular setting we studied. Specifically,
 11 the following graphs go “beyond the two CIFAR-10 networks” (R3) to include the significantly less-parameterized
 12 “ResNet20” (R3). **To show that our results hold without tuned hyperparameters for training and pruning**, we
 13 trained ResNet18 and ResNet20 with the exact training configuration described in the CIFAR section of He et al.
 14 (2015)—**SGD, data augmentation, weight decay**, etc.; i.e., all of the factors requested by R1 except for dropout.
 15 These studies use a “more robust pruning regime” (R2) by pruning a **constant fraction of all layers in all blocks** with
 16 the common filter ℓ_1 norm (Li et al., 2017) “for calculating filter importance” (R2). We emphasize that, even when
 17 using this “simplest, possible prune setup” and “standard training pipelines” (R1), the generalization-stability tradeoff
 18 is clearly present, confirming that our observations do indeed apply to “more general settings” (R1)—and based on the
 19 ResNet20 results, seem to apply to less “modern” regimes (R1), though modern regimes inspired our central question.



20 **R2,R3 Generalization gap vs. test accuracy; train accuracy not reported:** Correct, we neglected to state that all
 21 models had 100% training accuracy. With constant training accuracy, higher generalization (test accuracy) implies
 22 a smaller generalization gap. Thus, lower stability improves generalization *and* reduces the generalization gaps
 23 (overfitting)! We will update our manuscript to clearly discuss training accuracies and plot the generalization gaps.

24 **R3 “Pearson correlation and slope do not give an accurate characterization”:** Correct, the graphed relationships
 25 are not always linear. The manuscript will add a Kendall rank coefficient (Kendall, 1938; Jiang, 2019), τ , for each
 26 Pearson coefficient. Of the 20 statistically significant Pearson tests, 19 of the Kendall tests are statistically significant
 27 (one p-value went from 0.04 to 0.07) and all tests had the same sign, further supporting the hypothesized tradeoff.

28 **R2,R3 Methodology in “main body” and its “clarity”:** We will move methodological details to the “main body”
 29 (R2, R3), improve their “clarity” (R3), and accent that our method lacks “deal-breakers” (R1) as shown above.

30 **R1, R3 Hyperparameter choices (“These networks reach much lower accuracy than expected... L1/L2 regular-**
 31 **ization is disabled”):** Section 2 justified our exposition’s focus on less-regularized models, which is not unprecedented:
 32 the phenomena explored in the main text of a double descent paper referenced by R1 (Nakkiran et al., 2019) were
 33 produced without weight decay—in fact, our CIFAR-100 experiment in Table F.1 explicitly mimicked their approach,
 34 including use of “data augmentation” and networks with “residual connections” (R1). As for our “per-layer” (R1) and
 35 other pruning settings, these are not “deal-breakers” (again, see above), though our final draft will explain their origin
 36 was our initial, unpublished work (excluded to preserve anonymity) that pruned only the last dense layer of a small
 37 network. It led to our exploring pruning of the last convolutional layers of VGG11/ResNet18. The present study added
 38 pruning of other layers until we saw breakdown points—illustrated in the Section 3.2 “limit studies” appreciated by R2.

39 **R2,R3 “[DSD] is worth a comparison” and “the claim... is hard to extract”:** We thank the reviewer for pointing
 40 out this reference and we will include the following discussion in the related work. DSD (Han et al., 2017) shows that
 41 pruning based generalization improvements are retained and improved after restoring pruned connections. Relative to
 42 DSD, we show that the parameters can reenter at zero *or* their original values (Figure D.2) while achieving the full
 43 benefit of pruning and temporary Gaussian noise can replace pruning events to achieve a benefit.

44 **R3 “[15] is not found to improve generalization”:** [15] (LeCun et al., 1990) says OBD improved test error in the last
 45 sentence of page 603 and improved “recognition accuracy” in the third sentence of the conclusion.

46 **R1 Most experiments use VGG; VGG lacks batch norm:** We emphasize that many of our experiments were
 47 performed with ResNet18: Figures 2 and B.2, and Tables C.1 and F.1. We will clarify in the text that our VGG11 model
 48 **includes batch normalization** and was used in part because it allowed us to create replicates with limited resources.