1 We thank reviewers sincerely for their detailed feedback. It has enabled us to improve and clarify several sections.

2 **1. The results are not tested on ImageNet / ResNet [R1, R2, R3].** We agree it would be great to scale these experiments to deeper models (ResNet) and larger datasets (ImageNet) and see how well they generalize. As found by Frankle et al in "Stabilizing..." (2019), not all Lottery Ticket (LT) results generalized to ImageNet-sized models. We suspect that some of our results may similarly not generalize (such as mask-1 actions) but some will (such as Supermasks and mask-0 actions because our hypothesis on masking as training is independent of the original initialization of the weights). Due to time and computational constraints we leave this for future work, but have added a note in our paper discussing the lack of generalization results. *[See Update a below]*

9 **2. Significance of results are unclear; mask criteria curves unclear [R1].** Some of the results were indeed poorly described; we have updated several descriptions and the plot captions and think the clarity has been improved significantly *[Updates b, e, f]*. Summary: the significance of "magnitude increase" results are two-fold: it shows the LT phenomenon is not exclusive to the large final criterion, and the results are consistent with our hypothesis of masking as training. The mask-1 "significance of the sign" results show, in contrast to the LT paper, that the basin of attraction for a lottery ticket network is actually quite large: anywhere in the correct weight quadrant optimizes well; optimizers encounter difficulty crossing the zero barrier between signs. Finally, "masking as training" proposes an explanation for many of our results (and is predictive of performance of both mask criteria and alternate mask-0 treatments). We have also made the plots easier to interpret by simplifying the confidence bands and highlighting the main takeaways in the caption *[Updates d]*.

19 **3. Missing results on trained performance of additional Supermasks [R2].** This is a great suggestion. We have run these experiments and included them in SI *[Update g]*. We found that large_final_same_sign actually slightly underperforms large_final when trained in the iterative pruning LT paradigm. This is likely because we break 0 ties randomly when a mask criterion does not contain enough 1's, and selecting only weights with the same sign entails the need for much random selection.

24 **4. Questionable use of statistical significance [R2].** Iterative pruning does indeed cause correlation across different pruning percentages. However, we never report p-values aggregated in this way: all statistical tests are performed at a single pruning percentage, aggregating only across five independent runs. We've clarified in the text *[Update h]*.

27 **5. Mask-0 experiments and Fig 5 poorly motivated [R2].** We've removed the unnecessary straw man argument and state the real motivation more simply: to test the hypothesis that the value of frozen weights affects training *[Update c]*.

29 **6. Further analysis of masks themselves + principled understanding of lottery tickets [R2, R3].** We agree that it would be a valuable research direction to better understand and provide a theoretical foundation for what makes good masks, as well as their connection to better initializations. We believe our results represent a significant though not consummate step toward better understanding. Relating to the **question on the performance of signed_reinit [R3]**, our results hint that using a bimodal distribution to initialize weights may be more beneficial than a Gaussian distribution. This may explain why signed_reinit underperforms signed_reshuffle and signed_constant, both of which are initialized using a bimodal distribution.

36 **7. Clarifications on Supermask training [R2, R3].** The Bernoulli is sampled each forward pass, so during training each mini-batch entails a fresh sample of the mask. At test time, we report the average accuracy over 10 random samples of the mask (aggregating accuracy of the single-sample models, not ensembling predictions). As **[R3]** mentions this is akin to learning a per-weight dropout probability, but we should be careful not to interpret this as learned regularization but as the entire learning process itself. **[R3 point 2.3]** We respond to two possible interpretations of your suggestion here: (1) Supermasks are indeed trainable from initial weights that are const magnitude, random sign (Fig 7, right). (2) If signs are randomized *after* training a Supermask, performance degrades to chance (experiment run but not shown).

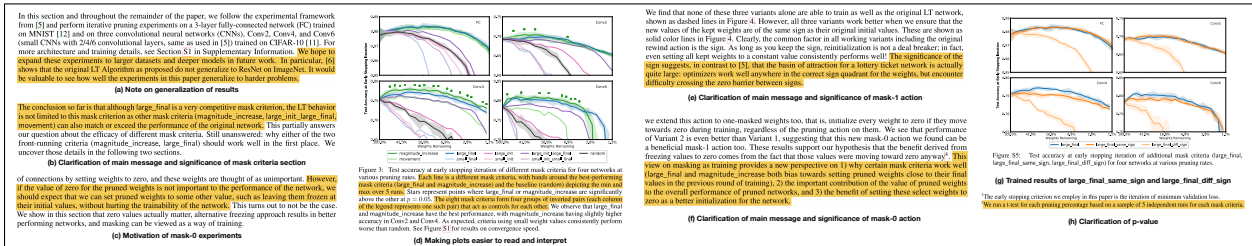43 **8. Initialize the learned Supermasks based on a well-performing heuristic mask criteria [R3].** This is a great idea and would likely work, though for now, because it works and is simple, we have trained only from scratch.

45 **9. Mask using other points on optimization trajectory instead of endpoints [R3].** Great idea. We have not yet explored computing masks using (intermediate, final) weights as you suggest instead of (initial, final) weights. We did try masking using (initial, min-val-err) weights, but performance was not as good.

48


49