

1 We thank the reviewers for their time and valuable feedback on our paper. We've responded to comments below:

2 **Motivation and Contribution:** In this study, our primary goal was to understand whether winning tickets contain
3 generic inductive biases which apply to multiple problems or rather, whether winning tickets are simply overfit to the
4 particular dataset and optimizer used to generate them. One of the most exciting aspects of the lottery ticket hypothesis
5 is that it suggests that we may be able to generate new initialization schemes which can substantially improve training
6 of neural networks from scratch without requiring any pruning or fine-tuning. However, such initialization schemes are
7 only possible if winning tickets contain generic inductive biases. We therefore view the primary contribution of this
8 study as a scientific one – understanding the generality of winning tickets – rather than a prescriptive one, though we
9 hope that this study will help to lay the groundwork for improved training methods in the near future. Importantly, the
10 absolute performance of winning tickets is independent of the generality of winning tickets.

11 **Comparisons against other pruning methods (R1, R2):** We completely agree with the reviewers that comparisons
12 between winning ticket trained networks and more traditional model compression approaches are interesting and
13 important problems, but we feel this an orthogonal question to the primary question of this study (how generic are
14 winning tickets), especially since these comparisons are difficult to make fairly. Pruning approaches generally require
15 the full model to be trained from scratch, after which the model is generally iteratively pruned and fine-tuned. In contrast,
16 winning tickets (especially transferred winning tickets) can be trained from scratch, and do not require any further
17 fine-tuning. As such, the costs of generating these models are quite different, making fair comparisons challenging. We
18 will add a discussion point clarifying these differences in section 5.1.

19 As Reviewer 1 acknowledged, our study required massive amounts of compute, especially since we performed six
20 replicates of each experiment. For example, we trained over 1500 ImageNet/Places365 models from scratch to generate
21 the results in Figs. 3e,f and 4e,f. We therefore view such comparisons as beyond the scope of the present work.

22 **Analysis of lottery tickets (R1) :** We wholeheartedly agree with the reviewer that this an exciting and important line
23 of inquiry (and one that we are currently pursuing!). However, detailed analysis of the structure of winning tickets is
24 quite challenging since winning tickets and random tickets are often statistically quite similar (at least to the first few
25 moments). We also agree that comparing pruning patterns to those derived from L_0 or L_1 regularization are interesting,
26 but we haven't analyzed these yet. As such, we leave this problem for future work.

27 **Rescaling of weights (R2):** To be clear, after each pruning iteration, a new mask is generated, after which the
28 subsequent model is trained from scratch with the new mask. Because converged weights are not used after pruning (as
29 in typical compression approaches), there is no need to re-scale the weights. Additionally, the relevant comparisons for
30 our primary question are between the transfer ticket performance, same dataset ticket performance, and random ticket
31 performance, all of which have the same scaling of initialization.

32 **Differences in input dimension (R2):** For both the VGG19 model and the ResNet50, a global average pool is applied
33 across all spatial dimensions prior to the final linear output layer. As a result, changes in input dimension do not require
34 changes in model architecture. We will add a paragraph clarifying this point within section 3.3 in the final paper.

35 **Global pruning scaling (R2):** We have not evaluated rescaling the weights based on the pruning fraction. However,
36 the relevant comparison here is between winning tickets and random tickets neither of which is rescaled. We therefore
37 consider it unlikely that rescaling would change our core results since we have no reason to expect that rescaling would
38 preferentially benefit winning or random tickets.

39 **Clarification of target/source in lines 150-155 (R2):** In order to evaluate the generality of winning tickets, winning
40 tickets were generated by iteratively training and pruning a model on one dataset/optimizer ("source" configuration,
41 as defined in lines 144-145), and evaluated on a second dataset/optimizer ("target" configuration). For standard
42 lottery ticket experiments in which the source and target dataset are identical, each iteration of training represents the
43 winning ticket performance for the model at the current pruning fraction. However, because the source and target
44 dataset/optimizer are different for our experiments and because we primarily care about performance on the target
45 dataset for this study, we must re-evaluate each winning ticket's performance on the target dataset, adding an additional
46 training run for each pruning iteration. In the final paper, we will expand this section to better clarify this point.

47 **Inconsistent colors (R1, R3):** Thank you for pointing out the inconsistency in colors between plots! We agree that
48 this makes the plots difficult to read and will correct the colors to be consistent across plots and figures in the final paper.

49 **Inclusion of preserved mask (R2):** While we feel the use of globally random masks is important to properly evaluate
50 the performance of random tickets (because masks can leak substantial amounts of information from the final trained
51 model to the initialization, as discussed in Section 3.1), we will add additional plots with the preserved mask case to the
52 appendix for our primary results in the final paper.

53 **Generality claims (R2):** We agree with the reviewer that this claim is a little too strong since, indeed, some winning
54 tickets generated on small datasets do not generalize to larger datasets. We will soften this claim in the final paper.