

1 **Tightness of bounds in Theorem 3.1 (All reviewers)** First we note that there is a nontrivial class of networks for  
2 which the upper bounds are tight when considering their linear interpolations. Namely, networks with ReLU activation  
3 function and mean squared error (MSE) loss. Section C.1 (supp. material) details the proof of Theorem 3.1. The bound  
4 depends on the following inequalities; Lipschitz continuity of the activation function, matrix norm inequalities for the  
5 layer weights  $W_i$ , a triangle inequality when expressing the intermediate activations in layer  $i$  as the combination of  
6 those in layer  $i - 1$ , a triangle inequality related to the  $\epsilon$  optimality of the endpoints, and the local Lipschitz continuity  
7 of the loss function. For the class of networks mentioned, the last inequality becomes trivial given we have MSE loss.  
8 Weights can be found for ReLU networks such that the other inequalities are tight.

9 We now address the tightness of the bound for piecewise linear curves between networks, which themselves can  
10 approximate continuous curves. Under the assumption that a piecewise linear curve or limit of such sequence truly does  
11 exist along which the loss is optimal and constant (reasonable as mode connectivity assumes this), for two models on  
12 the same line segment of this curve, we can consider the optimal curve between them. It follows that this optimal curve  
13 is their linear interpolation. Therefore, we have tightness in the upper bound for curve-finding restricted to piecewise  
14 linear curves between a class of networks for which we have tightness in the linear interpolation for that class. Via  
15 an argument by continuity, this tightness extends to continuous curves. Thus, these bounds are nontrivial as we have  
16 tightness for a wide class of networks and curve parameterizations under a reasonable assumption. We will add the  
17 discussion on tightness in the revised version.

18 **Small gap between theory and experiments (R1)** We discuss the use of post-activations in section C.2. From Figure  
19 9 (supp. material), it is clear the curve learned with alignment via pre-activations outperforms the unaligned curve.  
20 Thus, the theorem applies to a successful method. Since we see best performance from alignment using the correlation  
21 of post-activations, we use that method for the majority of the experiments. We discuss the relation between these two  
22 techniques in section C.2. To the reviewer’s other point, we also note that ReLU is Lipschitz continuous with constant 1.

23 **Neuron Alignment algorithm and symmetry (R1)** The symmetry for dense layers, convolutional layers, and residual  
24 blocks do have different properties as the reviewer mentions. By conducting experiments on TinyTen, ResNet32, and  
25 GoogLeNet, we show that mode connectivity with neuron alignment is able to offer performance gains with these  
26 different layers present. We note that Figure 1a, shows that there is still meaningfully symmetry detected by neuron  
27 alignment for all of these mentioned layers. Thus, we believe the algorithm can be widely applied.

28 **Sketch of the proof of Theorem 3.1 (R1)** We can approximate the intermediate activations of a network along the  
29 linear interpolation as an interpolation of the endpoint activations. Then the distance between the activations to those of  
30 an endpoint is bounded by the sum of this approximation error and the distance between the two endpoint activations.  
31 Thus for aligned networks, this second term in that bound is smaller as the permutation minimizes that value. As we  
32 iteratively do this for each layer, we finally bound the distance between the output of the network along the curve to that  
33 of the endpoint models, giving us a bound on the error of the network that is tighter with alignment.

34 **Multi-stage alignment (R1)** We appreciate the suggestion of a multi-stage neuron alignment. We were curious if  
35 models along the curve are aligned to the endpoint models. Figures 1b and 1c show this for the curve midpoint.  
36 Interestingly, for aligned curves, the midpoint on the learned curve is already optimally aligned to the endpoints.

37 **Significance of performance (R3)** We emphasize our contribution by including the minimal accuracy along the curve  
38 in Table 1. We can see that there is a model along the unaligned curves that clearly has subpar performance compared  
39 to independently trained models (i.e. the endpoints). Alignment notably improves this. Additionally, we have included  
40 figures that show the performance of the linear mode connectivity with and without alignment. This can be seen in  
41 Figure 2 (top left) and Figure 4 (supp. material). The performance gain from alignment is very notable, up to 60%  
42 difference in minimal accuracy along the curve, in these figures.

43 **Expression of the bounds (R3)** The bounds can be stated recursively via the equations in section C.1. We chose to  
44 omit it from the main body in the interest of space.  $B_u$  is a nonnegative linear combination of the endpoint activation  
45 distances,  $\|f_i^u(0) - f_i^u(1)\|_2$ . For  $B_a$ , the only difference in the expression is the use of  $\|f_i^u(0) - P_i f_i^u(1)\|_2$  for each  
46 layer. So when the intermediate activations are closer, as in aligned networks, the upper bound becomes smaller.

47 **Comparisons to Garipov et al. (R4)** It seems that there was a misunderstanding. We compare the performance of  
48 aligned curves to unaligned curves in our results. When the curves are unaligned, these are precisely curves trained via  
49 the method of Garipov et al. Thus our paper shows improvement over existing methods.

50 **Choice of subset for Algorithm 1 (R4)** By a subset of the training data, we were referring to a random subset of  
51 the training data that well approximates the underlying data manifold. This is to emphasize that the alignment can  
52 be computed more quickly on a random subset compared to the whole training dataset. As a test, we computed the  
53 alignment for a pair of TinyTen networks trained on CIFAR100, where the subsets contained 2500, 5000, and 10000  
54 images respectively. Despite these being random subsets of differing size, they produced the same alignment.