

1 Thank you very much for your detailed reviews and comments. In the rebuttal we will focus on the main issue raised:
2 lack of clarity in the description of our theoretical model. At the end of the rebuttal we will address the remaining
3 comments.

4 **Confusion about our landscape model toy task and the definition of n -wedges.** A common point you brought up
5 is the difficulty in understanding our loss landscape toy task construction, especially what exactly we mean by the
6 n -wedges. We found that in order to be able to verify whether a particular landscape model matched the behaviour
7 observed in real nets, we needed to implement an explicit simulation.

8 The simplest version of our toy landscape is constructed as follows. We populate the D -dim weight space with n -dim
9 low-loss attractors we call n -wedges. Each of these n -wedges has n infinitely extended *long* dimensions, and $D - n$
10 infinitely thin *short* directions. We take each n -tuple of axes, and position a single n -wedge such that its long directions
11 are aligned with them. We then define a surrogate loss $\mathcal{L}_{\text{toy}}(\vec{P} \in \mathbb{R}^D)$ for a network configuration P in this weight
12 space, which we choose to depend monotonically on the L_2 distance to the nearest n -wedge. Luckily for us, this
13 distance is simply $d(P) = \sqrt{\text{sum}(\text{sorted}(P)[D - n]^2)}$ – an easy to understand explicit expression

14 While this construction is very specific, we find that it is the dimensions D and n that influence our results, rather than
15 the specific angles between the n -sheets or their axis-alignment. As such, our toy model serves us well, albeit it doesn't
16 capture many other features of the loss landscape. Nonetheless, on this landscape, we are able to perform connectivity
17 experiments, as well as experiments with optimizing on random hyperplanes, and empirically verify the similarity to
18 real network experiments.

19 In real nets, we find a large number of weight-space directions in which we can move very far, while the loss doesn't
20 change – those would be the n long directions of the wedge; we also find a small number of extremely sensitive
21 directions in which a small motion incurs a high loss cost – those are the $D - n$ short directions. Together, these
22 define locally an n -dimensional hyperplane of finite thickness in the remaining $D - n$ thin direction, i.e. a *cuboid*.
23 Experimentally we notice a strong effect of radius $r^2 = \sum_i w_i^2$, the sum of squares of all weights. While locally a
24 cuboid, we find that individual parts of the manifold of low loss points radiate from the origin at a well-defined range of
25 angles, like a *wedge*. We find the full low-loss manifold to be a union of those in different directions and orientations.
26 We will include this extended discussion in the paper. We will also include an Appendix with a detailed description of
27 the toy landscape + the code that we use to experiment with it + we will publish a demo Jupyter Notebook / Colab.

28 **R1: More experiments, larger networks, and harder datasets.** To strengthen the case for our landscape model,
29 we extended the experiments in our paper to include fully-connected as well as convolutional networks of various
30 sizes (width, depth, non-linearity) including large models such as the **ResNet20v1** (>90% test on CIFAR-10), trained
31 on MNIST, Fashion MNIST and CIFAR-10 & 100. To go beyond classification, we also looked at CNN-based
32 autoencoders. In all cases the results supported our landscape model and we will include them in the final version.
33 This also demonstrates that our landscape model did not overfit to a small CNN on F-MNIST, as it holds for other
34 architectures and datasets.

35 **R5: Overfitting the landscape model to a particular task? New predictions and their empirical observation to
36 the rescue.** We constructed a model for the loss landscape of neural networks based on existing observations in
37 literature and our own verification of them. While we were very happy that our model incorporates them all (people in
38 general had trouble reconciling them together), what gave us confidence were new effects that we predicted based on
39 the model, that we only *later* observed in real networks. Those were 1) the existence of $(N - 1)$ -dimensional low-loss
40 connectors between N -tuples of independent optima, and the scaling of the number of short (=high curvature) directions
41 in their middle with N , 2) the changing of the predicted labels in the middle of a low-loss connector between two
42 optima, 3) stochastic weight ensembling (SWA) not working when checkpoints are too far from each other (belonging
43 to different wedges). We were aware of none of those at the time of building our model, and only later we predict they
44 should happen, and verified them in real networks.

45 **R2: Getting better at visualizing high-dimensional intuitions in 2D.** During the time between the submission and
46 now, we developed a better set of figures and explanations to convey the high-dimensional intuitions in 2D and 3D. For
47 example, we have a better version of Figure 1, where we do not make the wedges circular and smooth, as this was a
48 confusing illustration for some of our readers.

49 **R5: Radial tunnels = what low-dimensional cuts would show.** We noticed a confusion about the two types of
50 "tunnels" we discuss: we use the low-loss connectors between two independent optima as an observation to reconcile
51 with our model + a diagnostic tool. The other type of a tunnel – the radial tunnel – is what we would see on 2D cuts
52 through the landscape. At any point in training, making a random, 2D visualization of the loss around our current
53 point, we would (very likely) see a convex depression. As training progresses mainly radially, and at each point there is
54 convex depression around us, we can visualize this as a radial tunnel going out. We will be clearer with the distinction
55 in the final version.