

1 We thank the reviewers for their detailed comments and questions. Typos are fixed.

2 **R1 Line 137-138:** In these experiments we change the size of the shuffled blocks all the way to 1 and even try shuffling
3 the channels of the input. Therefore, the only object that is preserved here is the set of all pixel values which can be
4 looked at as a histogram/distribution and be fully characterized by its moments. We have now clarified this in the paper.
5 **Line 179-188:** It has been observed in the literature that *any* two neural network minimizers can be connected via a
6 non-linear low-loss path. In contrast, as mentioned by the reviewer, due to non-linear nature of neural nets one would
7 not expect the model obtained by the linear interpolation of two neural nets in the parameter space to perform well, and
8 that is precisely why we find this phenomenon interesting and investigate it. Note that insights from this observation can
9 be used for improving ensemble methods. **Line 108-109:** One can come up with quantitative measure by comparing
10 domains higher level representations by passing them through a network but here what we meant was looking at data
11 itself, and we have included some samples in the paper. Real includes real images, while Clipart is a cartoon version,
12 both have colors and texture. On the other hand ChexPert is black and white and has texture while Quickdraw has no
13 color or texture. In light of addressing your main concerns, we respectfully ask you to consider accepting the paper.

14 **R2 Back n forth to Appendix:** Thank you for your suggestions. We have now moved the distance in feature space
15 results to the main part of the paper and moved the ‘train’ part of figure 5 to the Appendix. **L151-152:** We mean
16 ambiguous data points in the target domain. We will revise the text to clarify this. **L191-192:** Two P-T models in the
17 same basin are not trivially (approximately) equal, as higher layers compute quite different features (Table 1). We
18 agree that using interpolation coefficient as x-axis does not tell the whole story. But with ReLU you can artificially
19 scale up weight magnitudes and distances. It’s the connectivity rather than the actual distance that is important here.
20 Regarding RI-T having different random seeds, we added experiments of different RI-T models from *the same* random
21 initialization values. They still converge to *different disconnected* basins. We also added extrapolation (with coefficient
22 in [-1, 2]) to our connectivity analysis. **Spectrum section:** Thank you for your question, what we meant is among two
23 models that can classify with certain margin, which translates to having low cross-entropy loss, then we can look at
24 concentration of spectrum and concentration towards low values shows less confidence. More mathematically speaking,
25 the confident model requires a lower-rank to get ϵ -approximation of the function and therefore, has lower capacity.
26 However, the model at init (either random or pre-train) has higher loss and is not fitted to the data.

27 **R3 Examples of negative transfer?** Pre-training with random labels *can* hurt transfer learning [Maennel et al 2020].
28 **Practical usage:** Our findings on basin in the loss landscape can be used to improve ensemble methods. Our observation
29 of higher order statistics improving training speed could lead to better network initialization methods.

30 **R4 Novelty:** The block-shuffling experiments and the linear interpolation of two neural networks experiments are
31 not trivial and it was unclear what to expect beforehand. They were not done before either. Moreover, experiment
32 outcomes matching the intuitions does not necessarily renders the experiments uninteresting or unimportant. **Claims:**
33 We clarify that we mention in the paper “lower layers are in charge of more general features” (comparing to more
34 class-specific features in higher layers) and we do not claim they are in charge of feature reuse. (line 240-242) **Explain**
35 **methods used:** CKA[Kornblith et al] is the latest work on estimating feature similarity with superior performance
36 over earlier works. The algorithm in [Sedghi et al] is the one that provides correct singular values for convolutions
37 and has been used for various applications. We cited the works and referred the reader to original papers for details
38 for space constraints. We have now added details of both methods to the Appendix. **Related work in NLP:** We have
39 now added [Brown+2020, Devlin+2018, Tamkin+2020, Raffel+2019, Liu+2015, Radford+2019, Roberts+2020] to the
40 related works on transfer learning in NLP. **Difference between RI-T & P-T:** The reviewer has brought up this point
41 multiple times that better performance of P-T, P-T’s being in the same basin, and the results of shuffling experiments is
42 *expected* due to the fact that it starts from an optimized point. We need to emphasize that the pre-training model has
43 been optimized to fit a *different* dataset (e.g. ImageNet). It would not necessarily converge faster or perform better in
44 downstream tasks (esp. with shuffled blocks) simply because it is ‘already optimized’. In fact, if the outcome is the
45 opposite, we could perfectly justify it by saying that it is trapped in a local optima because it is ‘already optimized’
46 (to fit a different task). Moreover, optimization speed is affected differently from accuracy. **Two P-T’s in the same**
47 **basin being a result of large number of epochs:** We have performed many experiments and also results from Section
48 4 show that once the models pass the initial finetuning stage, they belong to the same basin and a large number of
49 epochs is not needed. We also added new experiments of two RI-T models from *the same* initializations, and unlike
50 P-T, they do *not* converge to the same basin even after many epochs. **Higher layers being less similar compared to**
51 **lower layers:** The two P-T models are in the same basin but are not identical functions. Stochasticity of SGD leads to
52 non-identical models. Lower layers are the foundation of the function and higher layers depend on lower layers. If
53 the lower layer feature is more robust there are many solutions at higher layers and each of them can perform well.
54 Moreover, in a ResNet model you may not be using the weights in the higher layers and may be using skip connection.
55 Nevertheless, P-T’s are more similar compared to other models (Table 1). In light of addressing your main concerns, we
56 respectfully ask you to consider accepting the paper.