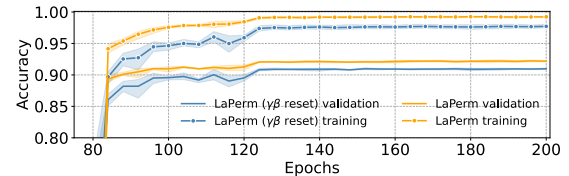


1 We greatly appreciate the reviewers for the time and expertise they have invested in the reviews. This essay focuses on
2 addressing the main concerns raised by each reviewer. Due to the page limit, we will not be able to mention and answer
3 all the questions. However, we would like to thank the reviewers for every observation and suggestions for revision and
4 clarification and will make sure to improve on every aspect of our paper.

5 **1. The meaningfulness of the analysis of “weight profiles”.** (Reviewer 1) In Section 2, our focus is on the statistical
6 similarity between weight vectors within a trained weight matrix and it is crucial for later developing the insight:
7 reordering of neuron connections can result in learning. A majority of the main results in our paper are obtained using
8 *uniform random distribution* as opposed to the gaussian-like distribution seen in well trained neural networks. Moreover,
9 we have conducted experiments (not mentioned in the paper) using different initializations, e.g., uniform, with different
10 scales and observed that the resulting weights to be still roughly gaussian with SoWP.

11 **2. Training BatchNorm (BN) alone yields decent performance [1].** (Reviewer 1) To isolate the contribution of BN’s
12 scaling (γ) and shifting (β) terms from LaPerm reconnected
13 networks, we propose to add experiments using the following
14 “ $\gamma\beta$ reset” training scheme into the appendix. Take ResNet50 as
15 an example: we use BN as usual and update γ and β using the
16 inner optimizer of LaPerm. At each synchronization, we reset γ and β to 1 and 0. We compare its performance with



17 LaPerm (never reset γ and β) both using $k = 800$ (other experiment settings are the same as Section 5.4). We repeat the
18 experiment three times and show their results on the right. We observe only roughly 1% decrease in the final accuracies
19 when γ and β do not hold information. In contrast, [1] focuses on the expressive power of γ and β .

21 **3. Section 5.2 lacks sufficient detail.** (Reviewer 1) We would like to formalize our claim as following: We hypothesize
22 that there are 2+1 dimensions of differences between initial and learned weights: D_0 : Common distribution of weights,
23 D_1 : Locations of weights, and D_2 : Deviations of exact learned weights from the best weights representable by
24 well-chosen distribution (D_0) and locations (D_1). Each weight vector w_j is represented as $w_j = \sigma_j(\theta) + \delta_j$, where θ is
25 a vector drawn from the common distribution (D_0), σ_j is a permutation operation (D_1), and δ_j is the remainder (D_2).
26 This decomposition, however, is not unique, unless we put more restrictions on the choice of θ and σ_j . — The rest of
27 the subsection will be revised accordingly.

28 **4. LaPerm can be understood as a regularizer and is sensitive to k .** (Reviewer 1) We agree that LaPerm has
29 inherent regularization effects. But our claim — DNNs can be trained by only reconnecting random weights — would
30 imply more than just regularization. Sections 5.4 and 5.5 provide partial results supporting our view. As for sensitivity
31 of k , if we define being sensitive as "responding unpredictably to slight changes", then LaPerm is not quite sensitive.
32 For example, we showed in Section 5.2 that the performance of LaPerm responded monotonically w.r.t. k and changing
33 it from 500 to 2000 effects only 2% of the accuracies.

34 **5. The paper does not seem focused enough.** (Reviewer 1) Our focus is dual: (1) most of the information of DNNs
35 is stored in the orderings of weights, (2) we can actually construct algorithms to find good orderings. We think both are
36 equally important and cannot defocus one of them. We would like to improve the first paragraph of sections 2 to 5 in
37 the final version of the paper to further emphasize this structure for clarity.

38 **6. What is the reason to have multiple synchronizations in LaPerm?** (Reviewer 2) LaPerm appears to need
39 repeated synchronizations to find the optimal reordering. We conducted experiments on Conv4 under the same setting
40 as in Section 5.2, but choose to synchronize only once at the end of the training, we obtained on average 13.8% (both
41 validation and training) accuracies. We will merge this experiment into Figure 5 to stress this point.

42 **7. Performance under a higher sparsity.** (Reviewer 3) Results for higher sparsity for Conv4 on line 222 are reported
43 in Figure 8 (2). Instead of only showing their final performance w.r.t. “% of weights” as we did in the paper, we will
44 add their detailed training trend in the appendix. In Figure 7, we will show two more trends for $p = 5\%$ and 1% .

45 **8. The costs of lookahead.** (Reviewer 4) Except when k is extremely small, LaPerm, on average, has few extra
46 computational overheads in addition to the cost of its inner optimizer. For example, running the scripts provided in
47 the supplementary material on a Google Colab GPU runtime, synchronizing Conv13 (14.9M parameters) once takes
48 around 200ms. For Conv13 in Figure 6, we needed to synchronize totally 125 times which only added 25s to the overall
49 training time. Moreover, a larger k is observed to work well (e.g. Figure 5 in the paper) and thus should often be used.
50 We will improve the main text in Section 4 to stress this point.

51 **9. Response to details and clarity.** (All Reviewers) We thank all the reviewers for the careful observations. We
52 will improve the clarity of the figures. We will revise the main text and expand the paper’s references and appendix
53 according to the reviewers’ suggestions in their “Correctness”, “Clarity”, and “Prior work” sections, especially focus on
54 constraining our claims to only classification problems, stressing the architecture choices and their purposes, and fixing
55 errors such as changing L2 weight decay to L2 regularization, fixing ambiguity and mistakes in notations and formulas.
56 Due to limited time and space, we are unable to answer the comments about: trying other inner optimizers, conducting
57 large-scale experiments, analyzing decision boundaries, theoretical justifications. We will pursue them in future works.

58 [1] Frankle et al. "Training BatchNorm and Only BatchNorm..." arXiv 2003.00152, 2020.