1  We thank the reviewers for their constructive feedback. Our study explores how the intrinsic dimension (ID) of data
2  manifolds varies across the layers of state-of-the-art deep neural networks (DNNs). Our findings can be summarized as
3  follows: **1)** the ID follows a "hunchback" shape across the layers of a trained DNN (Fig. 3); **2)** data representations
4  live on low-dimensional, but curved manifolds (Fig. 5); **3)** in an untrained network, the ID in all layers is very similar
5  to the ID of the input (Fig. 5C), which can be made arbitrarily low by adding low-level features (Fig. 6); and **4)** the
6  classification accuracy of a network depends on how low the ID of data manifolds is in its last hidden layer (Fig. 4).
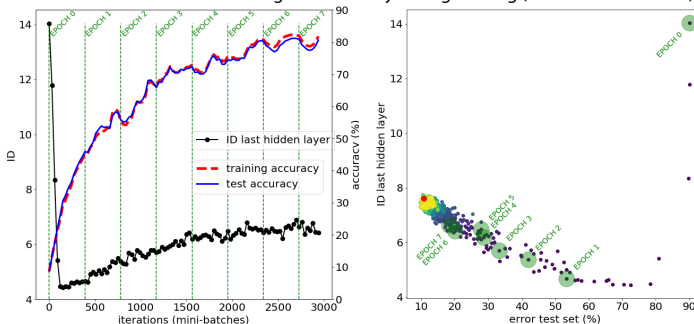
7  As recognized by reviewers R2 and R3, these findings are new. Therefore, we are surprised that reviewer R1 found our
8  conclusions not to be novel enough. We agree that we unfortunately overlooked one very pertinent reference [Ma et
9  al, 2018] and we thank the reviewer for pointing it out – we will cite and discuss this paper in our revision. However,
10 **our results are complementary with [Ma], rather than overlapping.** In fact, [Ma] studied how the ID of a specific
11 layer (the second last) in a DNN varies during optimization, and not how the ID varies across the layers of fully trained
12 DNNs. Only one analysis in our study is partially redundant with [Ma] (and we will acknowledge this accordingly)
13 – the increase of ID produced by random labels reported in Sec. 3.5 (red curve in Fig. 6B). But **none of the major**
14 **findings summarized above can be found in, or easily inferred from, [Ma] or any other study we are aware of**.

15 We are particularly surprised by R1 conclusion that "the findings of 3.2 ... have limited novelty", since "[Ma] has already
16 observed that the intrinsic dimension decreases during training (and the accuracy increases ...)". First, [Ma] analysis is
17 limited to the last hidden layer. Our study shows that, for early and middle layers, the trend is exactly the opposite
18 - as a result of training, the ID increases (Fig. 5C). Second, [Ma] does not show that the ID in the last hidden layer
19 predicts the accuracy achieved by different state-of-the-art DNNs, as demonstrated in our Fig. 4. Finally, **following R2**
20 **suggestion, we run new analyses to monitor the ID evolution during training**. We observed a monotonic decrease
21 of the ID only in some layers of a NN trained with MNIST, but not in the last hidden layer of VGG-16 trained with
22 CIFAR-10 (see New Fig A). Here, after an initial drop, the ID slowly increased, but, differently from [Ma], without
23 producing any overfitting. This shows that our results not only are largely novel with respect to those of [Ma], but also
24 challenge their generality. Specifically, this new analysis suggests that the evolution of the ID during training may
25 depend on the specific architecture and dataset. This, of course, calls for further investigation.

26 Further expanding on **R2 suggestion**, we also estimated the variability of the ID in the last hidden layer of a VGG-16
27 adapted for CIFAR-10 across 50 different trainings, finding no correlation with accuracy (r=-0.003), likely because of
28 the little variation in accuracy produced by different random weight initializations. This suggests that differences in
29 accuracy across well-trained networks (Fig. 4) are mostly due to differences in the architecture. In **response to R2**, we
30 have also verified that the ID variation across layers is generally consistent across object classes (see New Fig. B).

31 In our revision, following **R1**, we will underline that our ID estimator is a global, not a local one. We will also address
32 **R3 questions** about the reliability of our ID estimates. **First**, we will expand the discussion on the robustness of the ID
33 estimator with respect to the dimension of the embedding space. We already performed a test on artificial data of known
34 ID, embedded in a 100,000 dimensional space. The test did not reveal any significant degradation of accuracy. **Second**,
35 we will comment on the relevance of the independence assumption we made before eq 1. Indeed, the ID can also be
36 estimated by restricting the product in eq 1 to non-intersecting triplets of points, for which independence is strictly
37 satisfied, but in practice this does not significantly affect the estimate. **Third**, we will stress that our ID estimate is not
38 severely affected by the presence of "hubs". Indeed, the ID is not strongly affected by aggressive decimation of the
39 dataset (Fig. 2B), a procedure which would kill the hubs. **Fourth**, we will provide an interpretation of the value of the
40 ID in the output layer. This value cannot be 1, since all the networks we considered perform classification among $N_c$
41 categories. Its value should be of the order of $\log(N_c)/\log(2)$. We will also extend the discussion on the implications
42 of the reported low ID. Indeed, we believe that the coordinates on the manifolds should capture the invariance of the
43 representation of objects, and are therefore "meaningful", even if their explicit representation cannot be extracted easily.
44 We will also cite and discuss ref [Achille et al], which we agree is relevant.



**A** Evolution of ID and training/test accuracy during training (VGG-16 on CIFAR-10)

**B** Single manifold IDs in AlexNet (ImageNet)