

1 We sincerely thank all reviewers for their meaningful and detailed comments that will help improve our work. In the  
2 following, we address the raised concerns one by one:

3 **Reviewer 1.** We thank Reviewer 1 for his/her constructive feedback. (i) Linear-  
4 and kernel models: We consider rank collapse to be a phenomenon that is specific  
5 to the compositional structure of deep networks and hence  
6 shallow models are naturally robust against it. (ii) Binary classification: we agree that mapping all data points to a line does  
7 not impede this per se but we would like to point out that the direction of the rank collapse after initialization is in fact random  
8 and *independent* from the data, which makes it very unlikely  
9 to coincide with the principal direction in the data. Particularly  
10 so in very high dimensional neural networks. Fig. 1 shows  
11 that rank collapse remains an issue when reducing CIFAR10  
12 to a binary classification task (using class 0 and 1 only). (iii)  
13 Networks wider than the input: Good point, we will make sure  
14 to clarify this in the paper. If  $N < c\sqrt{d}$  for some constant  $c$ , then the rank indeed remains at  $\Omega(N)$ .  
15  
16

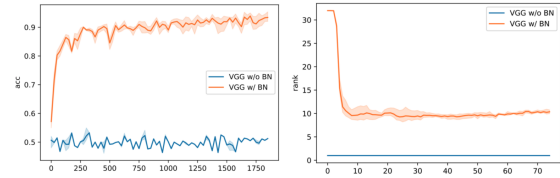


Figure 1: VGG19, two class CIFAR, SGD

17 **Reviewer 2.** We thank Reviewer 2 very much for sharing his/her thoughts on our work. (i) Building new DNNs: We  
18 would like to emphasize that the main goal of the work at hand was to deepen the understanding of batch normalization  
19 and its crucial interplay with random i.i.d. initialization. Furthermore, as we show in Section 4, rank collapse is not  
20 necessarily a problem of the architecture per se but it is closely linked to the way the networks are initialized. As can  
21 be seen in Figure 4, a sophisticated (rank preserving) initialization can even outperform BN in very deep networks.  
22 However, we do agree that developing new architectures with rank collapse in mind is indeed a very interesting direction  
23 to follow up on. As we elaborate next, we think our findings are already of interest to the community, as they offer a  
24 novel view on BN which we have not seen in the literature.

25 (ii) Broadness of impact: We believe that the work at hand makes a significant contribution to a better theoretical  
26 understanding of the delicate interplay between architecture, initialization and optimization algorithm. Such a theoretic-  
27 ally sound perspective can serve as the basis for further practical advancements on all three of these fronts. One such  
28 advancement can be found in Section 4 which includes an accelerating pre-training step that is developed directly from  
29 our theoretical analysis. Furthermore, since batch normalization itself is one of the major architectural developments of  
30 the last decade with successful application in countless settings, we do believe that this result can be an important step  
31 towards a systematic improvement of optimization methods for neural networks based on theoretical intuitions.

32 **Reviewer 3.** We are very grateful for the comments made by Reviewer 3. (i)  $\gamma$ : Indeed, our analysis needs  $\gamma$  to be small  
33 (depending on the width), which comes from the fact that our proof technique relies on low-order Taylor approximations.  
34 However, we note that  $\gamma$  is, importantly, independent of  $L$ ! Furthermore, as can be seen in Fig. 3 this seems to be an  
35 artefact of our proof technique, since the result holds empirically even for networks without skip connections ( $\gamma = \infty$ ).

36 (ii) Overfitting: Well spotted! In our examples (Fig. 11) it seems to us that the overfitting tends to happen whenever  
37 a certain accuracy is achieved on the training set, regardless of the actual method that is used. Please see Fig. 2 to  
38 confirm this intuition: Here, one can see that when reaching  $\sim 75\%$  training accuracy (epoch 100 for pre-training and  
39 200 for pre-training) both methods yield a similar test loss ( $\sim 1.75$ ). This is not necessarily surprising as the networks  
40 considered are MLPs that might have a harder time at finding generalizable patterns in image data compared to e.g.  
41 convnets.

42 (iii) Extensions to CNNs: This is definitely a very interesting  
43 follow up directions. In fact we are currently looking into  
44 extending our theoretical result to convolutional neural nets.  
45 The notion of rank is a bit more subtle in convolutional layers  
46 as the hidden presentation are third order tensors but we still  
47 observe that (after unfolding the tensors) randomly initialized  
48 convnets suffer from rank collapse unless batch normalization  
49 layers are in place (please see Figure 1 above as well as Fig. 11  
50 in the appendix of our paper).

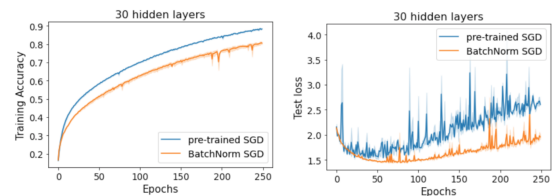


Figure 2: Fig 11 from paper for more epochs