

1 We thank all the reviewers for their detailed comments. Please find our response to the major comments below (we will
2 fix all typos/minor comments in the final version).

3 **Response to Reviewer 1:**

4 **Related work:** We will include comparison with [Fiat et al., 2019] in the main section.

5 **Feature Learning:** It is true that NPF learning happens when we train standard DNN with ReLU and it is only
6 natural that it closes the gap. However, NPF learning (interpretation) is different from the standard interpretation of
7 feature learning, where, the hidden features are learnt in the penultimate layer and the final layer learns a linear classifier
8 on these features. To see the difference, consider $S1 = \{\text{FLNPF, DNPFL, and ReLU-DNN}\}$ vs $S2 = \text{FRNPF}$ vs $S3 =$
9 Infinite width CNTK. $S1, S2, S3$, all of them generalise ($S2$ is the least with 67%, yet, on a 10 class task like CIFAR-10,
10 is still way better than random classification accuracy of 10%). Both $S1$ and $S2$ are finite width, so standard feature
11 learning happens in both $S1$ and $S2$, but, $S1$ with NPF learning is better (78% or above in CIFAR-10) than $S2$ (67% in
12 CIFAR-10) with no NPF learning. Thus neither finite width, nor the standard feature learning is useful to explain the
13 difference between $S1$ and $S2$. $S2$ and $S3$, both have no NPF learning, yet, $S3$ generalises better, which can be attributed
14 to the fact that infinite width ensures better averaging and hence results in a well formed kernel. NPF learning also
15 happens in the DNPFL setting, which is different from ReLU-DNN. To conclude, *NPF learning is a measure which*
16 *discriminates/describes the different regimes (i.e., $S1, S2, S3$) better than the standard feature learning explanation.*

17 **Response to Reviewer 2:**

18 **Assumption 5.1 and Goal of the paper:** Most analysis of DNNs with ReLU is on what happens at initialisation.
19 In a DNN with ReLU, NPV and NPF are not statistically independent at initialisation, i.e., Assumption 5.1 does not
20 hold. However, in the current state-of-the-art analysis, in $w \rightarrow \infty$ regime, activations change only at rate of $\sqrt{\frac{1}{w}}$
21 ([Jacot et al., 2019]), i.e., activations/NPFs do not change during training. Hence, though Assumption 5.1 may not hold
22 exactly, it is *not a strong assumption* to fix the NPFs for the purpose of analysis. Thus, statistically decoupling the NPV
23 from NPFs is only natural, and furthermore it adds strength: the $\text{NTK} = \text{const} * \text{NPK}$ is *interpretable* in terms of the
24 active sub-networks (NTK is defined in terms of gradients with no interpretation), which also shows that the active
25 sub-networks are *fundamental entities arising naturally in the NTK framework*.

26 Fundamental role of gates is further accomplished by the FLNPF experiments, where we show that by copying gating
27 information alone, and resetting and training NPV (i.e., Θ^V) from scratch (Assumption 5.1 holds in this case), we can
28 recover the performance of the DNN with ReLU. Further, in the DNPFL setting, (which is not hypothetical) Assumption
29 5.1 holds, and the DNPFL does generalise well in the experiments.

30 **Prior experimental work:** The goal of this paper is not to study the utility of active sub-networks as representations,
31 but to directly look at the generalisation capability. However, we will move relevant work (example [Srivastava et al.,
32 2014]) in the appendix to the main body.

33 **Regularisation:** Similar to [Arora et al., 2019], data-augmentation, batch norm, residual connections, dropout and
34 other forms of regularisations are avoided. However, studying these in our framework is future work.

35 **Response to Reviewer 3:**

36 MNIST and CIFAR-10 are used as standard datasets in most analytical works such as ours, see [Arora et al., 2019] for
37 example.

38 **Response to Reviewer 4:**

39 **Analytical part in supplementary:** The main contribution of the paper is analytical in nature, and aims at
40 providing an understanding into the internal workings of a DNN. We believe that it is critical that the setup and
41 subsequent explanations belong to the main body. While we give details of the experimental setup in the appendix, we
42 will be happy to move the important points in the main body to improve clarity.

43 **Explanation of generalisation:** By generalisation we mean performance on test data. Agreed that on MNIST,
44 all the cases namely FRNPF, DNPFL, FLNPF and ReLU-DNN have marginal performance difference. However, on
45 CIFAR-10 the difference between FRNPF (67%) and FLNPF, DNPFL, ReLU-DNN (all above 78%) is more than 10%.
46 The crucial insight from this work is that mere gating/masking property is enough to give us 67% (on CIFAR-10, this is
47 non-trivial because a random classifier will only have 10% accuracy), and in addition if the gates also *adapt* during
48 training (as in standard ReLU-DNN) gives the rest 10%. Further, once we have the *learnt* gates, we can reset and learn
49 NPV from scratch without loss in performance. Thus, the experiments were designed to test the role of gating and we
50 believe we have extensive experiments to support our claims.

51 **Expanded conclusion section:** We will make the summary of main contributions more clearer in the conclusion.