

1 We greatly appreciate all the reviewers' comments and feedback. Here are our responses.

2 **Reviewer 1**

3 – *Significance: Unclear. The authors use the adjectives "shallow" and "deep" in several places but really they seem*
4 *to be contrasting 1- vs. 2-hidden-layer networks. It is not clear how/whether the approach can lead to non-trivial*
5 *separation results for *actual* deep networks.*

6 **A:** We agree with the reviewer that the end goal should be understanding the effectiveness of DNNs that have many
7 layers, however, we believe that our results are steps in that direction.

8 In this work, we are using the terms “shallow” versus “deep” as is common in the deep learning literature. (For instance,
9 one of our main references, (Eldan-Shamir 2016), also uses these terms exactly as defined in this paper.)

10 – *Consider including a more comprehensive comparison with existing separation results (e.g Eldan-Shamir 2016).*
11 *For example, even though both show polynomial upper bounds in the size of the hidden layers, it appears that there*
12 *exists a significant reduction in the degree of the polynomial in the current manuscript.*

13 **A:** This is a good point. In the revised version, we will add the following paragraph at the end of Section 7 which is on
14 related work: “For a specialized radial function in R^n , [Eldan-Shamir 2016] shows that while any SNN would require
15 at least exponentially many nodes to approximate the function, there exists a DNN with two hidden layers and $O(n^{19/4})$
16 nodes that well approximates the same function. In this paper, for a general class of widely-used functions viz GMM
17 discriminant functions, we show that while SNNs require at least exponentially many nodes, for any GMM discriminant
18 function there exists a DNN with two hidden layers and only $O(n)$ nodes that approximates it.”

19 **Reviewer 2**

20 – *No numerical results are provided to this paper. For a theoretical work, this can be acceptable. But some*
21 *experimental results should be no harm.*

22 **A:** We agree that some numerical evaluations could be interesting for our achievability results. As a matter of due
23 diligence, we did run numerical cross checks but decided to prioritize discussion of the theoretical results.

24 **Reviewer 3**

25 – *If the authors had managed to relax the Assumptions 1-3 to cover more practical settings, my rating would be higher,*
26 *but I think it is entirely appropriate to postpone such improvements to future work.*

27 **A:** First, please note that Assumptions 1-3 are only required for proving Theorem 1 (our achievability result). As
28 explained in Section 4, Theorem 2 (i.e., our converse result) holds for general activation functions. In fact, it also
29 applies to the cases where each node is allowed to have a different activation function. Second, as explained in the
30 conclusion section, we could relax these assumptions, and allow e.g. ReLU activation function at a minor expense to
31 the size of the network. In fact, we are currently working on various relaxations of these assumptions.

32 – *Consider changing the title: I personally find the word "approximation" misleading, the objective here is not to*
33 *learn or to approximate GMMs, but instead to separate them;*

34 **A:** Even though we are motivated by classification of GMMs, we chose to put ‘approximation’ in the title, as our main
35 technical results focus on approximating the discriminant functions of GMMs.

36 – *Clarify if the conclusion in Section 5 on the sufficiency of exponentially many nodes holds true if the covariance of x*
37 *is *not* proportional to identity, and if its holds for a larger class of activation function sigma;*

38 **A:** We are able to show that the result holds for a larger class of activation functions (recent work) and we expect that it
39 holds for more general covariance matrices as well (planned next step).