

1 We thank the reviewers for their helpful comments and we will fix accordingly. Below we address some comments and  
2 questions that have been raised.

### 3 **Reviewers 1 and 3**

4 The reviewers wrote that requiring only polynomial versus not polynomial in the separation definition is a bit limiting.  
5 The benefits of depth can be studied from different angles. We focus on the notion of depth-separation that has been  
6 extensively studied in many prior works, and requires polynomial vs. exponential width. We agree that other notions of  
7 depth-separation should also be explored, and we mention in the “related work” section an example for such a different  
8 notion that has been studied in [23,16,31].

### 9 **Reviewer 2**

10 The reviewer asked: “bounded functions can oscillate exponentially and has an  $\exp(d)$  Lipschitz norm. However,  
11 the first theorem suggests they can be approximated with only  $\text{poly}(d)$  weights. This looks very counterintuitive.  
12 How should we understand it?”. This observation is indeed surprising, however, it has a simple intuitive explanation.  
13 Consider a univariate function  $f : \mathbb{R} \rightarrow \mathbb{R}$  that is expressible by a neural network  $N$  of constant depth and  $\text{poly}(d)$   
14 width (for some integer  $d$ ). Such a function is piecewise-linear with a  $\text{poly}(d)$  number of pieces. Since the weights  
15 in  $N$  might be exponential in  $d$ , then some of the linear pieces might have exponential derivatives. Thus,  $f$  might  
16 oscillate quickly in some intervals. However, since  $f$  is bounded, then an interval where  $f$  has an exponential derivative  
17 must be very small (exponentially small). Hence,  $f$  consists of  $\text{poly}(d)$  linear pieces, but may oscillate quickly only in  
18 exponentially-small intervals. A network of constant depth and  $\text{poly}(d)$  weights cannot approximate such  $f$  in the  $L_\infty$   
19 sense. However, since we assume that the input distribution  $\mu$  is not too concentrated in very small intervals, then it is  
20 possible to approximate  $f$  in the  $L_2(\mu)$  sense. The case where  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is more complicated, but follows the same  
21 intuition.

### 22 **Reviewer 4**

23 Thanks for your feedback and comments, we will incorporate them into the final version. Below we address your  
24 specific questions:

25 (2) Not for  $k > 2$ .

26 (5) Since we focus on constant-depth networks, then it does not matter. We will comment on that in the final version.

27 (7) There is no special reason for this choice of notations. We will change.

28 (9) The results hold also for approximation w.r.t.  $L_1(\mu)$  (it follows easily from our proof), but do not hold for  
29 approximation w.r.t.  $L_\infty$ .