

1 We thank the reviewers for the helpful feedback and the positive assessment of our submission. We plan to update the
2 text and bibliography following their suggestions.

3 *Reviewer #1, “It is interesting to see if further increase the width of the network (from linear in d to polynomial in d and
4 even exponential in d), does the discontinuous approximation requires a smaller depth.”*

5 In the setting of our paper (minimization of the total network size) a large depth is in some sense unavoidable (as e.g.
6 Theorem 3.2 shows). However, in general there is of course some trade-off between width and depth. Depth is more
7 important for expressiveness (since a parallel computation can be serialized, but not vice versa) – for example, even
8 for 3-layer nets, Eldan & Shamir¹ show existence of small 3-layer nets that require exponential width if expressed
9 as 2-layer nets. As for the dependence on the input dimension d , the key point is how the approximated functional
10 family depends on d . Assuming a sufficiently constrained family (e.g. a ball in the Barron space² or functions with a
11 compositional structure³), one can in a sense avoid the curse of dimensionality and find low-depth approximations with
12 the width determined mostly by the required accuracy rather than d . We remark also that the weight discontinuity is
13 present to a certain extent in all optimized nonlinear models (not necessarily deep or involving coding constructions)⁴.

14 *Reviewer #3, “Concentrating information in a small proportion of high-precision weights in the network according to a
15 discontinuous assignment function seems unstable in the face of noise or imperfect optimization... Can the authors
16 comment on how their results might interact with these other sources of error?”*

17 There is indeed a significant instability, especially for approximations with periodic functions, when the information
18 concentration is highest. In our construction for ReLU networks, the issue is ameliorated by the fact that the information
19 is divided into independent chunks: first on the level of weights corresponding to different patches, and then also on
20 the level of weight digits corresponding to different positions in a patch. This suggests that the errors can be localized
21 and, borrowing again from the coding theory, one can hypothesize that we can improve stability by allowing some
22 redundancy and, for example, using something like error-correcting codes. On the other hand, in the construction with
23 periodic activations, the classifier output is a chain of interdependent computations, so any error will have a stronger
24 negative effect.

25 *Reviewer #4, “Theorem 5.1 extends the approximation results to all piece-wise linear activation functions and not just
26 ReLUs. So in theory, this should also apply to max-outs and other variants of ReLUs such as Leaky ReLUs?”*

27 That’s right, all these functions are easily expressible one via another using just linear operations ($\text{ReLU}(x) =$
28 $\max(0, x)$, $\text{LeakyReLU}_a(x) = \text{ReLU}(x) - a\text{ReLU}(-x)$, $\max(x, y) = \frac{x+y}{2} + \frac{1}{2(1-a)}(\text{LeakyReLU}_a(x - y) +$
29 $\text{LeakyReLU}_a(y - x))$), so any network of one type can be converted into another, exactly equivalent one, at the cost
30 of merely increasing the number of neurons by a constant factor.

31 *Reviewer #4, “I fail to see some intuitions regarding the typical values of r , d , and H for the networks used in practice.
32 While the approximation rate depends on r and d , the power law relationship uses a term $W^{-r/d}$ and d is the dimension
33 of the input and W is the size of the network parameters. In practice, d can be of the order of millions in imaging-related
34 applications, since we have million pixels in the image. In such cases, the approximation guarantees may be very weak
35 and may not provide any insights.”*

36 In imaging, an important difference from our setting is that reasonable images form only a small and complex subset
37 of the whole ambient million-dimensional space, whereas in our setting the approximation is defined for each input
38 vector from $[0, 1]^d$. Accordingly, the number of pixels is not the right value of d here, a more appropriate value would
39 be something like the intrinsic dimension of the image manifold. For example, for MNIST, suitable feature extraction
40 and dimension reduction allows to reparameterize the data set by 9 parameters⁵ while retaining classification accuracy
41 above 98%, suggesting $d \lesssim 10$. As for smoothness r , usual classification problems such as MNIST do not quite fit our
42 setting as the predicted output (the image label) is piecewise constant. One can assume some low value of smoothness
43 (say $r \sim 1$, assuming a Lipschitz continuation), or refer to results on approximation of piecewise smooth functions⁶. Of
44 course, all these estimates are very crude, and there are various other considerations for imaging-related problems that
45 must be taken into account (e.g., our results assume $W \rightarrow \infty$ at a fixed d , but in a practical problem larger networks
46 will have a higher “effective input dimension” in the earlier mentioned sense).

¹R. Eldan, O. Shamir, The power of depth for feedforward neural networks, COLT 2016

²A. R. Barron, Universal Approximation Bounds for Superpositions of a Sigmoidal Function, 1993, DOI: 10.1109/18.256500

³T. Poggio et al., Why and when can deep-but not shallow-networks avoid the curse of dimensionality: A review. Int. J. Autom. Comput. 14, 503–519 (2017)

⁴P. Kainen et al., Approximation by neural networks is not continuous, Neurocomputing 29(1), 47–56 (1999)

⁵A. Das et al, Dimensionality Reduction for Handwritten Digit Recognition, 2018, DOI: 10.4108/eai.12-2-2019.156590

⁶Ph. Peterson, F. Voigtlaender, Optimal approximation of piecewise smooth functions using deep ReLU neural networks, 2018, DOI: 10.1016/j.neunet.2018.08.019