

1 **General comments.** We thank the Reviewers for their detailed reviews and the feedback regarding the **theoretical**
 2 **study** of rational neural networks (NNs) and their **promising applications**. Reviewer 4 noted that the introduction of
 3 rational NNs is motivated by theory, while most exotic activation functions are only empirically supported. We address
 4 the referees' remarks on the theoretical results and numerical experiments here. The paper will be revised accordingly.

5 **Theoretical results.** Reviewer 3 highlights the comparison between this paper and Telgarsky's work. We wish to em-
 6 phasize that our key contribution is to employ a **composition of low-degree rationals** $r(x) = r_{\#layers} \cdots r_2(r_1(x))$,
 7 naturally realized by a NN, to approximate functions efficiently, while Telgarsky's work approximates a ReLU NN by a
 8 high-degree rational function not in the form of a NN. The Newman polynomials used by Telgarsky do not preserve
 9 (minimax) optimality under composition and result in an exponentially larger number of trainable parameters (see Fig. 1
 10 of the paper for a comparison). Finally, we present **optimal lower and upper approximation bounds**. The notion of
 11 size and Zolotarev functions will be explicitly defined in the revised version, as requested by Reviewers 3 and 4.

12 In response to Reviewer 2, concerning the choice of the degree (3, 2), we emphasize that **this degree appears naturally**
 13 **in the technical analysis** due to the composition property of the Zolotarev functions (see Section 3.1): the degree of the
 14 overall rational function r is a whopping $3^{\#layers}$, while the **number of trainable parameters only grows linearly**.
 15 A superdiagonal degree (3, 2) allows r to behave like a **nonconstant function at $\pm\infty$** , unlike a diagonal degree e.g.
 16 (2, 2). The theory ensures that low-degree rational activation functions minimize the number of trainable parameters
 17 given a fixed overall degree. The choice is also motivated empirically, and we do not claim that the degree (3, 2) is the
 18 best choice for all situations as the best configuration may well depend on the application; see below Fig. 1 (right).

19 Reviewers 1 and 3 discussed the apparently marginal difference between $\log(1/\epsilon)$ and $\log \log(1/\epsilon)$ in Thm. 4. To
 20 clarify, the bound for rational NNs is **close to optimal**, given by $\epsilon^{-d/n}$ (DeVore et al., 1989). Most importantly, a
 21 rational NN can achieve this approximation power with a depth of only $\log \log(1/\epsilon)$, which is **exponentially smaller**
 22 than the $\log(1/\epsilon)$ layers needed by a ReLU network to approximate a smooth function to within ϵ . This improvement is
 23 obtained from the composition of low-degree rational functions and is not hidden in multiplicative constants, which do
 24 not depend on ϵ . This improved approximation power has **practical consequences** for larger NNs given that a deep
 25 NN is computationally expensive to train due to expensive gradient evaluation and slower convergence. The constants
 26 inside the \mathcal{O} notations are computed whenever possible (see Lem. 1 and Cor. 3 of the Supplementary Material) but the
 27 main theorems of the paper treat a general setting, for which few (if any) papers in the literature give explicit constants.

28 **Experimental results.** We thank the Reviewers for their comments on the experiment in Section 3.1. While Fig. 2
 29 (right) of the paper showed that rational NNs are easier to train than ReLU and Sinusoid NNs, it did not show the
 30 accuracy loss. We have performed **new experiments with a validation set** independent of the training set and display
 31 the validation loss throughout training (see Fig. 1 (left)). We find that **rational NNs outperform the other NNs** during
 32 the training phase and on the testing set. Reviewer 3 was concerned that this result was due to the difference in number
 33 of trainable parameters but the ReLU NN had 8000 and the rational had $7 \times \#layers + 8000 = 8035$. **Rational NN**
 34 **are never much bigger in terms of trainable parameters** than ReLU NNs since the increase is only linear with
 35 respect to the number of layers. Reviewer 2 suggested a **comparison with polynomial approximation**, which we have
 36 performed (see Fig. 1 (left)). Here, we train a NN with degree 3 polynomial activation functions. We observe that this
 37 NN is **harder to train than rational NNs** as shown by the non-smooth validation loss. Polynomials perform poorly on
 38 non-smooth functions such as ReLU, with an algebraic convergence of $\mathcal{O}(1/\text{degree})$ (Trefethen, 2013) rather than the
 39 (root-)exponential convergence with rationals. Finally, following Reviewer 2's suggestion **we compare rational NNs**
 40 **with different degree activation functions** and find that they **all perform better than ReLU NNs** (see Fig. 1 (right)).

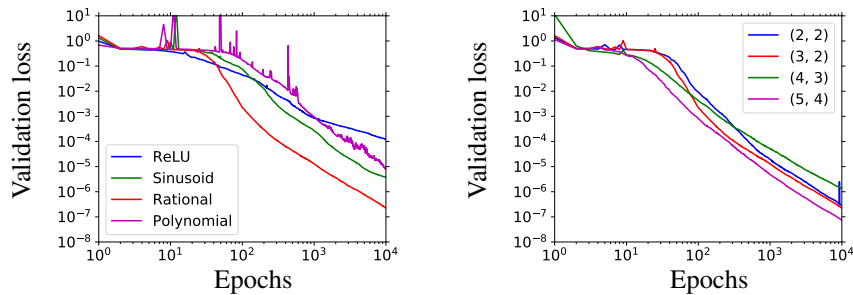


Figure 1: Left: Validation loss of a ReLU (blue), sine (green), polynomial (purple) of degree 3, and rational (red) NNs of type (3, 2) with respect to the number of optimization steps. Right: Comparison between the validation losses of rational neural networks of types (2, 2), (3, 2), (4, 3), and (5, 4).