

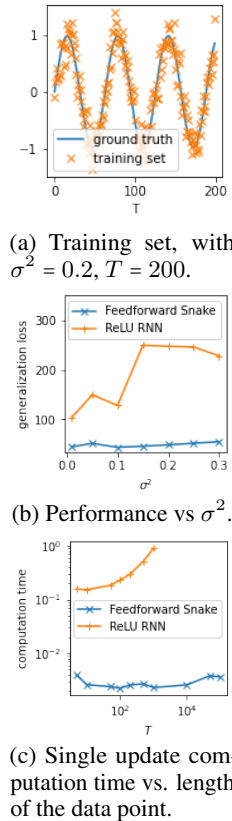
1 General Comment: We are very grateful to all the reviewers for carefully reading our paper.
 2 Their feedback has helped us to improve the paper accordingly. The two messages of this paper
 3 are that (1) learning a periodic or semi-periodic function with neural networks is yet unresolved,
 4 and we argue that the key to solve this is to focus on the extrapolation properties of neural
 5 networks; (2) we proposed to solve this problem with a simple alternative activation function.

6 **“A significant shortcoming is the lack of comparison with recurrent neural networks”**
 7 (R2, R4): We respectfully yet strongly disagree with this. It is indeed interesting to see how
 8 RNN would perform for these tasks. However, the problem with RNNs is that they implicitly
 9 parametrize the data point x by time: $x = x(t)$. It is hence limited to model periodic functions
 10 of at most 1d and cannot generalize to a periodic function of arbitrary dimension; e.g., it is
 11 not clear how one could define RNN to learn the function $f(x, z) = \sin(x) + \sin(z)$, which
 12 is an easy task for feedforward networks with Snake. For this reason, we do not believe our
 13 method (with Snake + feedforward) needs to be a competitor to RNN. This being said, we
 14 perform a comparison of RNN with Snake with feedforward on a 1d problem. See Figure 1.a
 15 for the training set of this task. The simple function we try to model is $y = \sin(0.1x)$, we add
 16 a white noise with variance σ^2 to each y , and the model sees a time series of length T . See
 17 1.b for the performance of both models, when $T = 100$, and validated on a noise-free hold-out
 18 section from $T = 101$ to 300. We see that the proposed method outperforms RNN significantly.
 19 On this task, One major advantage of our method is that it does not need to back-propagate
 20 through time (BPTT), which both causes vanishing gradient and prohibitively high computation
 21 time during training. In Figure 1.c we plot the average computation time of a single gradient
 22 update vs. the length of the time series, we see that, even at smallest $T = 5$, the RNN requires
 23 more than 10 times of computation time to update (when both models have a similar number of
 24 parameters, about 3000). This is a significant advantage of our method over RNN even for
 25 1d periodic problems. These results will be added to the final version.

26 **Stopping Criterion in Experiments** (R4): Our stopping criterion is chosen fairly and reasonably:
 27 all experiments are stopped at the time when the training loss of all the methods stop to
 28 decrease and becomes a constant, i.e., when the model has converged. The performances of
 29 this converged model is not visibly different from the early stopping point for the experiments
 30 we considered. For example, the goal of Figure 7 (for the atmospheric experiment) and Figure
 31 14.b (EUR-USD experiment in the appendix) are plotted in order to show that we stopped at the
 32 point when the model converged, moreover, neither of our model or the baselines seem to suffer significantly
 33 from overfitting, judging from these two figures. Therefore, the comparison is indeed fair and reasonable. In our final version,
 34 we will add similar plots for other experiments to clarify the stopping criterion for the experiments.

35 **Comparison to Swish and Leaky-ReLU Seems Necessary** (R5): The
 36 proposed method indeed outperforms Swish and Leaky-ReLU significantly for tasks in section 6.2 and 6.3. This is because Swish and
 37 Leaky-ReLU suffer from the same problem as ReLU, which is guaranteed by our theorem and by the discussion above. Therefore, we did
 38 not include them for visual clarity. Their performance on the tasks is now shown in Figure 2. We see that, for Swish and Leaky-ReLU, the
 39 learning is hard to mismatched inductive bias, and this leads to their inferior performance. We will add this plot to the appendix in the final
 40 version to avoid confusion. We will add this plot to the appendix in the final
 41 version to avoid confusion. We will add this plot to the appendix in the final
 42 version to avoid confusion. We will add this plot to the appendix in the final
 43 version to avoid confusion. We will add this plot to the appendix in the final
 44 version to avoid confusion.

45 **Strength of the claim** (R5): 1. **Applicability of the theorem.** We
 46 agree that some qualifier is needed in this claim. On the one hand, the
 47 specific statement of the theorem applies to tanh and ReLU, but it is general as discussed in line 64-66. Since the proof
 48 is only based on the asymptotic property of activation function when $x \rightarrow \pm\infty$, one can prove the same theorem for any
 49 continuous activation function that asymptotically converges to a tanh or ReLU; for example, this would include Swish
 50 and Leaky-ReLU (and almost all the other ReLU-based variants), which converge to ReLU; one can follow exactly
 51 the same proving procedure to prove a similar theorem for each of these activation functions. On the other hand, we
 52 agree that this statement might be too strong, and we will state the above condition clearly for the applicability of our
 53 theory in the final version to avoid confusion. 2. **Generality of Snake.** Here, we say that Snake is more “general” in
 54 the sense that it is not only capable of approximating a function in the bounded region, but also capable of extrapolating
 55 beyond a bounded region (in a periodic way). This claim does need some qualification to clarify, and we will modify
 56 this claim in the final version. **Some details. Some Problems in the Appendix** (R1): We will add more details and
 57 correct grammatical errors to the Appendix sections.

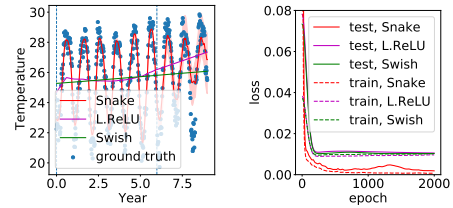


(a) Training set, with $\sigma^2 = 0.2, T = 200$.

(b) Performance vs σ^2 .

(c) Single update computation time vs. length of the data point.

Figure 1: Comparison of RNN with feedforward neural network with Snake.



(a) Same as Fig. 6 in the (b) Training loss and testing loss.

Figure 2: Comparison with Swish etc..