# Controlling Neural Network Smoothness for Algorithmic Neural Reasoning

Anonymous Author(s) Affiliation Address email

## Abstract

The modelling framework of *neural algorithmic reasoning* [1] postulates that a 1 2 continuous neural network may learn to emulate the discrete reasoning steps of 3 a symbolic algorithm. The purpose of this study is to investigate the underlying hypothesis in the most simple conceivable scenario – the addition of real numbers. 4 5 We find that two layer neural networks fail to learn the structure of this task and that growing the network's width leads to a complex division of input space. This 6 behaviour can be emulated with Gaussian processes using radial basis function ker-7 nels of decreasing length scale. Classical results establish an equivalence between 8 9 Gaussian processes and infinitely wide neural networks. We demonstrate a tight link between the scaling of a network weights' standard deviation and its effective 10 length scale on a sinusoidal regression problem, suggesting simple modifications to 11 control the smoothness of the function learned by a neural network. This provides 12 a partial remedy to the brittleness of neural network predictions. We validate this 13 further in the setting of adversarial examples where we demonstrate the gains in 14 robustness that our modification achieves on a standard classification problem of 15 handwritten digit recognition. In conclusion, we show inherent problems of neural 16 networks emulating even simple algorithmic tasks which, however, may be partially 17 improved with smoothness priors inspired by Gaussian processes. 18

## 19 **1** Introduction

The two most prominent paradigms in artificial intelligence research are discrete, symbolic algorithms 20 on the one side, and continuous, neural information processing systems on the other [2]. While 21 the latter have caused a revolutionary transformation of the field, they are often plagued by hard 22 challenges, such as robustness to changes in the input distributions, for which algorithmic approaches 23 can provide worst-case performance guarantees. Crucially, we know that both approaches are 24 deployed by humans, akin to Kahneman's 1 and 2 reasoning systems [3]. Thus, algorithms must be 25 implemented in biological neural networks in the human brain. Observing a scene in the world, we 26 know that it is represented and processed in the distributed representation of neural activity in visual 27 cortex. However, the same scene is also represented when we describe it with the use of symbols and 28 the syntax of our language. Thus, one of the most mysterious questions in neuroscience as well as in 29 artificial intelligence research is: where and how do these two representation systems interact? 30

The concept of neural algorithmic reasoning [1] is a recent proposal for a modeling framework at the intersection between symbol processing algorithms and continuous distributed information processing systems [see also 4, 5]. The obvious question is how such a hybrid architecture may be trained, since we usually require differentiability of the whole system for end-to-end training. To solve this, the authors [1] propose training a neural network to approximate the output of the algorithm in the middle of the model and allows end-to-end training.



Figure 1: Unit Disc Loss Surfaces. These plots show the network error in input space (remember, we are trying to model  $f : \mathbb{R}^2 \to \mathbb{R}$ ,  $f(x) = x_1 + x_2$ ,  $x \in \mathcal{D} \subseteq \mathbb{R}^2$ ), where brighter regions indicate lower error (red dots — training data points). Top row shows the learned function and its incurred loss for a ReLU neural network with increasing number of hidden units (N, left to right). Bottom row shows the loss surface for a GP with a RBF kernel of increasing length scale ( $\lambda$ ).

The purpose of our work is to investigate the feasibility of neural algorithmic reasoning in one of 37 the most simple conceivable settings: the addition of real numbers. Integer calculus and similarly 38 floating-point arithmetic in binary (symbolic) representations have previously received more attention 39 [6-9]. By contrast, we want to know if a simple (2-layer) neural network can learn to add real-valued 40 numbers on a compact domain such as the unit disc. This is an interesting setting because we know 41 the hypothesis class of neural networks with rectified linear unit (ReLU) activation functions  $\sigma$ 42 trivially contains the correct solution, i.e., given inputs  $x \in \mathbb{R}^2$ , the following function is a perfect 43 representation of the desired output  $y = x_1 + x_2$ 44

$$f(x) := W_2 \sigma(W_1 x + b_1) + b_2, \quad W_1 := \begin{pmatrix} 1 & 0 \\ -1 & 0 \\ 0 & 1 \\ 0 & -1 \end{pmatrix}, \quad W_2 := \begin{pmatrix} 1 \\ -1 \\ 1 \\ -1 \end{pmatrix}, \quad b_1, b_2 := \mathbf{0}.$$

Given input and output training pairs, a neural network can learn to approximate this function. However, it is an open question how accurate the approximation will be within the domain of the training data; possibly in the limit of infinite data and a model that is a universal function approximator [10]. Secondly, it is unclear how the model will generalise outside the domain of the training data.

In this paper, we find that artificial neural networks are unable to learn this simple function, even if abundant training data is available, leaving regions of high error within the training domain and struggling to extrapolate beyond. A comparison to Gaussian processes suggests a simple partial remedy, exploiting classic results about the equivalence between these two model classes, based on a correct adjustment of the smoothness of the learned function. We show that these modifications also translate to increased adversarial robustness on handwritten character image recognition.

## 55 2 Neural Networks and Gaussian Processes Learning Addition

We first investigate learning addition of real numbers in two dimensions. For this, we randomly draw 56 128 points on the unit disc and train a two layer neural network with ReLU nonlinearities after the first 57 layer to solve the addition task using a simple squared loss function  $\mathcal{L}(f, x, y) = (f(x, y) - (x+y))^2$ . 58 We train for 50,000 steps with the Adam optimizer [11] (held fix after initial experiments indicating 59 that this sufficed for convergence), an initial learning rate of 0.001 and a learning rate decay of 0.9. 60 Across the top row of Fig. 1 we change the number of hidden units (i.e., the width) of the network 61 and observe the effect on the learned solution. Few units exhibit the recently proposed polytope 62 structure of neural network approximated functions [12]. While more units slightly improve the 63 performance, we can see that the model uses the additional capacity to cut up the input space into 64



Figure 2: Annulus Loss Surfaces and Generalisation Performance. The first (third) column shows the loss surfaces for a NN (GP) with N hidden units, and an unlimited amount of training data (with optimal length scale and D training data points — red dots). The second (fourth) column shows the negative log loss as a function of eccentricity (dots coloured by angle, black dotted lines radially averaged) in the three regimes of OOD (within convex hull), IID and OOD (outside of convex hull) generalisation.

<sup>65</sup> increasingly refined regions. Interestingly, the network learns intricate ridges of good performance on

which the training data lies. These ridges appear to be connected on continuous paths — an intriguing
 observation for future foundational NN research.

Importantly, these patterns suggest that larger ReLU networks learn sieve-like solutions of increasing 68 resolution, but they do not enter a qualitatively different regime that would resemble the simple 69 algorithmic solution to the task. We can produce a similar sequence of model behaviours with 70 Gaussian process (GP) regressors with varying radial basis function (RBF) kernel scales. We use the 71 standard GP implementation in sklearn [13]. Usually, this would include maximum likelihood length 72 scale selection (optimum near  $\lambda \sim 150$ ), however, setting this by hand lets us visualise the different 73 solutions for *suboptimal* length scales. Specifically, we see that setting the length scale too low (Fig. 74 1 bottom left) forces the model to learn ridges of good solutions through the training data — similar 75 to the NN model. Note that the anti-diagonal line passing through the origin indicates the null space 76 of the target function, i.e., where a 0 output is the correct answer. 77

The similarity between the solution patterns for neural networks with many hidden units and GPs with short length scales is not surprising given classical results [14]. Briefly, Neal established that for NNs with a hyperbolic tangent (TanH) nonlinearity and an infinite number of hidden units (with appropriate scaling of their initialisation variances) the distribution over learned functions (*ab initio*) becomes equivalent to that of a GP. This opens an interesting path forward in understanding and improving the neural network solution exposed in this section (see below).

# **3** Out of Domain Generalization

A crucial difference between algorithms and neural network solutions is the way they generalise to 85 different inputs [2]. Addition is defined on all numbers in  $\mathbb{R}$ , but the approximation learned by a NN 86 can only observe a subset of those inputs in its training data (i.e., *IID* — independent identically 87 distributed). Recent discussions have investigated this from the point of view of interpolation versus 88 extrapolation [15, 16]. Overparameterized NN exhibit a *double descent* phenomenon, which is 89 thought to improve their generalization performance by interpolating between training points [17]. 90 91 although see [18]. Going beyond the convex hull of the training data would, by contrast, require the ability to extrapolate to a new domain (OOD — out of distribution). We adjust our setting slightly to 92 study both aspects of generalisation. 93

Specifically, the training data is now randomly sampled from the annulus  $\mathcal{A} := \{x \in \mathbb{R}^2 \mid 0.5 \leq 1 \}$ 94  $||x|| \leq 1$ . Moreover, to assess that the findings from Fig. 1 do not depend on limited training data or 95 finite network size, we set the number of hidden units to 10,000 and generate a new random batch for 96 every gradient step (totalling  $256 \times 50,000 = 12,800,000$  training examples). The learned solution 97 by the NN is shown in Fig. 2 top left. Again, we see an intricate pattern emerging within the training 98 data domain (IID) with ridges of good performance but valleys of bad predictions. By contrast, the 99 predictions are bad both inside the convex hull of the training data  $||x|| \le 0.5$  as well as outside 100  $||x|| \ge 1.0$ . This is quantified in Fig. 2 top middle as a function of the input norm (we can also see 101 the different ridges outside the training domain distinguished by their angles). For the GP, even with 102 limited training data (D = 256), we see higher performance levels (Fig. 2 bottom middle) both on 103 training data (green), in its convex hull (yellow) as well as outside (red). 104



Figure 3: **Sinusoidal Regression**. Top row shows GP models with different length scales ( $\lambda$ ) fitted to a sinusoidal regression problem with: true function (yellow IID, purple OOD), training points (red), model prediction (light blue), uncertainty (shaded, two standard deviations). The top right plot shows the different losses (mean squared error — MSE, and log likelihood) as a function of GP RBF length scale. The bottom row gives the same regression plots for the NN with TanH activation function and varying length scales ( $\sigma^{-1}$ ) and (right) the training and test error as a function thereof.

### <sup>105</sup> **4** Setting the Length Scale of Neural Networks

An important step in fitting a GP with a RBF kernel 106 to data is finding the best length scale for the kernel. 107 We explore this process, and its equivalent in NNs, 108 in this section. To easily visualise the learned input-109 output mapping, the dataset is now just a noisy sine 110 function on  $x \in [-2\pi, 2\pi]$  (for OOD, we extrapolate 111 to  $[-3\pi, 3\pi]$ ). Precisely, we are trying to model the 112 function  $y = \sin(x) + \epsilon$  with  $\epsilon \sim \mathcal{N}(0, 0.25)$ . In 113 Fig. 3 we can see that there exists a sweet spot for 114 the kernel's length scale (top middle), and that the 115 model overfits (top left) for smaller length scales and 116 underfits (top right) for larger length scales. This 117 is confirmed quantitatively (top right) by looking at 118 the likelihood (blue) as well as the training and test 119 error as a function of the length scale (dotted lines 120 indicate length scales in plots to the right). 121



Figure 4: Adversarial Robustness. Accuracy on MNIST as a function of  $(L_2)$  adversarial perturbation size for different models (TanH NNs, length scale  $\sigma^{-1}$ ).

To get a NN to behave like a GP, we closely follow the construction in [14]. That is, we use TanH 122 and restrict the standard deviation of the weights to  $\sqrt{N}$  where N is the number of inputs to a layer. 123 Importantly, we enforce this standard deviation throughout training by using the scaled weights 124  $\tilde{w} = w(s.d.(w)\sqrt{N})^{-1}\sigma$  with sd(w) the standard deviation of w and  $\sigma$  a scaling factor. Intuitively, 125  $\sigma^{-1}$  behaves like the length scale in GPs (Fig. 3 bottom), i.e., a smaller  $\sigma^{-1}$  means larger weights and 126 more overfitting to the training data whereas too large  $\sigma^{-1}$  means very small weights and underfitting 127 of the training data. Again, we establish the existence of an optimum length scale  $\sigma^{-1}$  (Fig. 3, bottom 128 right) that produces the smallest test error. 129

## 130 5 Controlling Neural Network Smoothness

Returning to the initial two examples of adding real numbers from a disc or an annulus, we can now observe the effect of varying the length scale ( $\sigma^{-1}$ ) of a 2-layer NN with a large number of hidden units (approaching the GP regime) and TanH activation functions. On both datasets (Fig. 5) it is apparent that increasing the length scale makes the learned output function more smooth. Thus, the NN becomes qualitatively more similar to a GP with a well adjusted learning scale for its RBF kernel.

As a further verification for the claim that this approach effectively controls the smoothness of the learned NN input-output mapping, we turn towards a more complicated problem of performing image



Figure 5: **Controlling the Smoothness of Neural Network Functions**. The four columns show learned TanH 2-layer NN loss surfaces for the two datasets (top disc, bottom annulus — same settings as above) for increasing length scales (left to right) producing an increasingly smooth mapping.

recognition on MNIST under worst case (adversarial) distribution shifts. Thus, we search for minimal perturbations that maximally change the output, which is a proxy for the model's smoothness [19].

Fig. 4 shows that controlling the length scale of a neural network on handwritten character recognition does indeed increase the robustness to  $L_2$  adversarial perturbations above a simple ReLU baseline model (both 2-Layer MLPs as above). We can see a trade-off between the clean and the robust accuracy with large length scales increasing the robustness while decreasing the clean accuracy. These performance levels are far from the SOTA robust models on this task, however, they eschew the need for (expensive) adversarial training [20]. Thus, this is a proof-of-principle that the smoothness of a NN function can indeed be controlled with the construction proposed in this paper.

## 147 6 Discussion

The promise of neural algorithmic reasoning to combine distributed and discrete reasoning systems 148 via differentiable NN approximations is intriguing. Here, we show that the functions learned by 149 neural networks, even in one of the most simple conceivable examples of symbolic manipulation 150 (real addition), is prone to learning a highly complex and varying output mapping that falls short 151 of learning the proper algorithmic target. There are many different perspectives onto this problem. 152 153 Many transformers are better at manipulating arithmetic expressions [9]. However, looking at the loss 154 surface of the NN solution for a simple 2-layer MLP, revealed an intriguing structure and analogy to GPs that we decided to pursue in this work. Moreover, the intricate structure of the loss surface is an 155 intriguing pointer for future directions in basic NN theory. 156

We have demonstrated how the NN mapping can be made more similar to that of a well-calibrated GP, specifically, by making it more smooth which also improved adversarial robustness. Note that designing Lipschitz neural networks with bounded smoothness is an open research area [21]. In our case, it is open for future work to investigate how the composite function of more than two layers contributes to global smoothness of a NN, potentially building on recent work extending the NN GP equivalence to multilayer networks [22].

Finally, while this study focused on a technical detail about NN approximations to simple algorithms,
the larger question remains still open how NNs can make the inferential step (see Hume's problem
of induction) from any finite amount of data to an *infinite look-up table* (see MLST episode 061).
Surely, humans are a proof-of-concept that noisy and distributed processing systems (i.e., brains) can
implement discrete symbolic algorithms [23]. We hope that future research in this area will benefit
more from interdisciplinary approaches that take inspiration across fields.

## **169** References

- [1] Petar Veličković and Charles Blundell. Neural algorithmic reasoning. *Patterns*, 2(7):100273,
   2021.
- [2] Gary F Marcus. *The algebraic mind: Integrating connectionism and cognitive science*. MIT press, 2003.
- [3] Daniel Kahneman. *Thinking, fast and slow*. Macmillan, 2011.
- [4] Daniel Bear, Chaofei Fan, Damian Mrowca, Yunzhu Li, Seth Alter, Aran Nayebi, Jeremy
   Schwartz, Li F Fei-Fei, Jiajun Wu, Josh Tenenbaum, et al. Learning physical graph representa tions from visual scenes. Advances in Neural Information Processing Systems, 33:6027–6039,
   2020.
- [5] Sara Sabour, Nicholas Frosst, and Geoffrey E Hinton. Dynamic routing between capsules.
   *Advances in neural information processing systems*, 30, 2017.
- [6] Rodrigo Nogueira, Zhiying Jiang, and Jimmy Lin. Investigating the limitations of transformers with simple arithmetic tasks. *arXiv preprint arXiv:2102.13019*, 2021.
- [7] Alon Talmor, Yanai Elazar, Yoav Goldberg, and Jonathan Berant. olmpics-on what language
   model pre-training captures. *Transactions of the Association for Computational Linguistics*,
   8:743–758, 2020.
- [8] Chengyue Jiang, Zhonglin Nian, Kaihao Guo, Shanbo Chu, Yinggong Zhao, Libin Shen, and
   Kewei Tu. Learning numeral embeddings. *arXiv preprint arXiv:2001.00003*, 2019.
- [9] Avijit Thawani, Jay Pujara, Pedro A Szekely, and Filip Ilievski. Representing numbers in nlp: a
   survey and a vision. *arXiv preprint arXiv:2103.13136*, 2021.
- [10] Kurt Hornik, Maxwell Stinchcombe, and Halbert White. Multilayer feedforward networks are
   universal approximators. *Neural networks*, 2(5):359–366, 1989.
- [11] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [12] Randall Balestriero et al. A spline theory of deep learning. In *International Conference on Machine Learning*, pages 374–383. PMLR, 2018.
- [13] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion,
   Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. Scikit learn: Machine learning in python. *the Journal of machine Learning research*, 12:2825–2830,
   2011.
- [14] Radford M Neal. Priors for infinite networks. In *Bayesian Learning for Neural Networks*, pages 29–53. Springer, 1996.
- [15] Preetum Nakkiran, Gal Kaplun, Yamini Bansal, Tristan Yang, Boaz Barak, and Ilya Sutskever.
   Deep double descent: Where bigger models and more data hurt. *Journal of Statistical Mechanics: Theory and Experiment*, 2021(12):124003, 2021.
- [16] Lukas Schott, Julius Von Kügelgen, Frederik Träuble, Peter Gehler, Chris Russell, Matthias
   Bethge, Bernhard Schölkopf, Francesco Locatello, and Wieland Brendel. Visual representation
   learning does not generalize strongly within the same domain. *arXiv preprint arXiv:2107.08221*, 2021.
- [17] Niladri S Chatterji, Philip M Long, and Peter L Bartlett. When does gradient descent with
   logistic loss find interpolating two-layer networks? J. Mach. Learn. Res., 22:159–1, 2021.
- [18] Randall Balestriero, Jerome Pesenti, and Yann LeCun. Learning in high dimension always
   amounts to extrapolation. *arXiv preprint arXiv:2110.09485*, 2021.
- [19] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversar ial examples. *arXiv preprint arXiv:1412.6572*, 2014.

- [20] Eric Wong, Leslie Rice, and J Zico Kolter. Fast is better than free: Revisiting adversarial
   training. *arXiv preprint arXiv:2001.03994*, 2020.
- [21] Mahyar Fazlyab, Alexander Robey, Hamed Hassani, Manfred Morari, and George Pappas.
   Efficient and accurate estimation of lipschitz constants for deep neural networks. *Advances in Neural Information Processing Systems*, 32, 2019.
- [22] Jaehoon Lee, Yasaman Bahri, Roman Novak, Samuel S Schoenholz, Jeffrey Pennington,
   and Jascha Sohl-Dickstein. Deep neural networks as gaussian processes. *arXiv preprint arXiv:1711.00165*, 2017.
- [23] Wim Fias, Muhammet Ikbal Sahan, Daniel Ansari, and Ian M Lyons. From counting to
   retrieving: Neural networks underlying alphabet arithmetic learning. *Journal of Cognitive Neuroscience*, 34(1):16–33, 2021.