



HAL
open science

RINO: Robust INner and Outer Approximated Reachability of Neural Networks Controlled Systems

Eric Goubault, Sylvie Putot

► **To cite this version:**

Eric Goubault, Sylvie Putot. RINO: Robust INner and Outer Approximated Reachability of Neural Networks Controlled Systems. Lecture Notes in Computer Science, 2022, Computer Aided Verification 34th International Conference, CAV 2022, Haifa, Israel, August 7–10, 2022, Proceedings, Part I, pp.511 - 523. <10.1007/978-3-031-13185-1_25>. <hal-03805908>

HAL Id: hal-03805908

<https://hal.science/hal-03805908v1>

Submitted on 7 Oct 2022

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



HAL Authorization



RINO: Robust INner and Outer Approximated Reachability of Neural Networks Controlled Systems

Eric Goubault^(✉) and Sylvie Putot

LIX, Ecole Polytechnique, CNRS and Institut
Polytechnique de Paris, 91128 Palaiseau, France
{eric.goubault,sylvie.putot}@polytechnique.edu



Abstract. We present a unified approach, implemented in the RINO tool, for the computation of inner and outer-approximations of reachable sets of discrete-time and continuous-time dynamical systems, possibly controlled by neural networks with differentiable activation functions. RINO combines a zonotopic set representation with generalized mean-value AE extensions to compute under and over-approximations of the robust range of differentiable functions, and applies these techniques to the particular case of learning-enabled dynamical systems. The AE extensions require an efficient and accurate evaluation of the function and its Jacobian with respect to the inputs and initial conditions. For continuous-time systems, possibly controlled by neural networks, the function to evaluate is the solution of the dynamical system. It is over-approximated in RINO using Taylor methods in time coupled with a set-based evaluation with zonotopes. We demonstrate the good performances of RINO compared to state-of-the-art tools Verisig 2.0 and ReachNN* on a set of classical benchmark examples of neural network controlled closed loop systems. For generally comparable precision to Verisig 2.0 and higher precision than ReachNN*, RINO is always at least one order of magnitude faster, while also computing the more involved inner-approximations that the other tools do not compute.

Keywords: Neural networks verification · Reachability analysis · Robustness · Inner-approximation

1 Introduction

Over the last few years, neural networks have emerged as an increasingly classical choice for the control of autonomous systems, in particular due to their properties as universal function approximators. However, their adoption in safety-critical systems, the inherent uncertainties from the dynamic environment, and their sensitivity to adversarial examples make it crucial to establish their safety and robustness. This verification is challenging because of the complex non-linear characteristics of neural networks. Recent works come up with some approaches

© The Author(s) 2022

S. Shoham and Y. Vizel (Eds.): CAV 2022, LNCS 13371, pp. 511–523, 2022.

https://doi.org/10.1007/978-3-031-13185-1_25

and tools to bound the output uncertainty of neural networks with respect to input perturbations. However, many of them are restricted to the analysis of networks with ReLU activation functions. Moreover, the approaches considering general differentiable activation functions and systems with general non linear dynamics provide over-approximations, which conservatism is difficult to estimate. RINO proposes a scalable and adaptive approach to compute both inner (or under) and outer (or over) approximations for the closed loop reachability problem of neural network controlled systems, with differentiable activation functions. The outer-approximation allows for property verification, while the inner-approximation allows for property refutation. Combined, the inner and outer-approximations allow to assess the conservatism of the approximations.

As the behavior of a neural network controlled closed-loop system relies on the interaction between the continuous dynamics and the neural network controller, a good precision requires to not only compute the output range but also describe the input-output mapping for the controller. In this work, we propose to use a zonotope-based abstraction to compute in a unified way both the reachable sets of neural networks and dynamical systems. This seamless integration of the reachability of neural networks and dynamical systems presents the advantage of a natural propagation of useful correlations through the different components of the closed loop system, resulting in an efficient and precise approach compared to many existing works which rely on external reachability tools.

Contributions

- RINO implements all ideas presented in [8–11] for the joint computation of inner and outer approximations of robustly reachable sets of differentiable nonlinear discrete-time or continuous-time systems (without neural networks in the loop), possibly with constant delays. These previous works demonstrated the good scaling properties of our approach on different examples including a full nonlinear quadcopter flight model but the tool was never presented as such.
- Additionally, we demonstrate here that an application of these ideas to the case of neural networks enabled dynamical systems provides very competitive results for the over-approximation compared to the state of the art (at least similar precision and one order of magnitude faster) while also providing the first approach for inner-approximation of the reachable sets of such systems, which we use to falsify some safety properties.
- Finally, RINO also computes approximations of output ranges that are reachable robustly or adversarially with respect to a subset of inputs: while these robust ranges are mostly used in this work to compute inner-approximations of joint ranges of state variables instead of projections, we believe this sensitivity information can be a useful tool in the future in particular to assess global robustness properties of neural networks.

Related Work. The safety verification for DNNs has received considerable attention recently, with several threads of work being developed. We draw below a non exhaustive panorama focusing on available tools for reachability analysis of neural network controlled systems with smooth activation functions.

Different approaches have been proposed to the reachability analysis closed-loop systems with neural network controllers, often by a transformation to a continuous or hybrid system reachability. Sherlock [6] targets both the open-loop and closed-loop problems with ReLU activation functions, in particular using the regressive polynomial rule inference approach [5] for the closed-loop, and Flow* [3] for the reachability of the dynamical system. NNV [24] also targets both the open loop and closed loop verification problems, with various activation functions and set representations such as polyhedra or star sets [23], and different reachability algorithms for dynamical systems relying on CORA [1] and the MPT toolbox [18]. ReachNN [13] and its successor ReachNN* [7] propose a reachability analysis based on Bernstein polynomials for closed-loop systems with general activation functions, also relying on Flow* [3] for the reachability of the dynamical system. Verisig [14] handles NNCS with nonlinear plants controlled by sigmoid-based networks, exploiting the fact that the sigmoid is the solution to a differential equation to transform the neural network into an equivalent hybrid system, which is then fed to Flow*. Verisig 2.0 [15] uses preconditioned Taylor Models to propagate reachable sets in neural networks, and also relies on Flow* for reachability of the hybrid system component.

The very recent works [21] and [12] implemented respectively over JuliaReach and in POLAR are also closely related to our work. In [21], the authors implement a bridge between zonotope abstractions and Taylor model abstractions in order to combine tools analyzing controllers (e.g. using zonotopes like deepZ [22]) with tools analyzing ordinary differential equations (e.g. Flow* [3]). In [12], the authors use a polynomial arithmetic made up of a combination of Bernstein polynomials and Taylor models to iteratively overapproximate networks layers, according to whether the activation function is differentiable or not.

2 Problem Statement and Background

2.1 Robust Reachability of Closed-Loop Dynamical Systems

We consider in this work a closed-loop system consisting of a plant with states x , modeled as a discrete-time or continuous-time system with time-varying disturbances w and inputs u , where some components of the control inputs can be the output a neural network h taking x as input. For notation's simplicity, we focus on continuous-time systems and define:

$$\begin{cases} \dot{x}(t) = f(x(t), u(t), w(t)) & \text{if } t \geq 0 \\ x(t) = x_0 & \text{if } t = 0 \end{cases} \quad (1)$$

where f is a sufficiently smooth function and at least C^1 , and controls u and disturbances w are also supposed to be sufficiently smooth C^k for some $k \geq 0$ stepwise. This allows discontinuous controls and disturbances, where the discontinuities can only appear at discrete times t_j .

The neural network h is a fully-connected feedforward NN with differentiable activation functions, defined as the composition $h(x) = h_L \circ h_{L-1} \circ \dots \circ h_1(x)$ of

L layers where each layer $h_i(x) = \sigma(W_i x + b_i)$ performs a linear transform followed by a sigmoid or hyperbolic tangent activation σ . We assume the control is decomposed as $u(t) = (u_1(t), u_2(t))$ where $u_2(t)$ is a control input defined in \mathbb{U}_2 and $u_1(t)$ is the output of the neural network controller. This controller is executed in a time-triggered fashion with control step T , so that $u_1(t) = h(x(t_k))$, for $t \in [t_k, t_k + T)$, where $t_k = kT$ for positive integers k . System (1) can then be rewritten as

$$\begin{cases} \dot{x}(t) = f(x(t), h(x(t_k)), u_2(t), w(t)) & \text{if } t \in [t_k, t_k + T), t_k = kT, k \geq 0 \\ x(t) = x_0 & \text{if } t = 0 \end{cases} \quad (2)$$

Let $\varphi^f(t; x_0, u_2, w)$ for time $t \in \mathbb{T}$ denote the *time trajectory* of (2) with initial state $x(0) = x_0$, for input signal u_2 and disturbance w .

We consider the problem of computing inner and outer-approximations of robust reachable sets as introduced in [9], defined here as

$$R_{\mathcal{AE}}^f(t; \mathbb{X}_0, \mathbb{U}_2, \mathbb{W}) = \{x \mid \forall w \in \mathbb{W}, \exists u_2 \in \mathbb{U}_2, \exists x_0 \in \mathbb{X}_0, x = \varphi^f(t; x_0, u, w)\}$$

Note that this notion of robust reachability extends the classical notions of minimal and maximal reachability [20]. We use the subscript notation \mathcal{AE} to indicate that the reachable set is minimal with respect to the disturbances w (universal quantification \mathcal{A}) and maximal with respect to the input u_2 (existential quantification indicated by \mathcal{E}), and that the universal quantification always precedes the existential quantification.

2.2 Mean-Value Inner and Outer-Approximating Robust Extensions

A classical but often overly conservative way to overapproximate the image of a set by a real-valued function $f : \mathbb{R}^m \rightarrow \mathbb{R}$ is the natural interval extension $\mathcal{F} : \mathbb{IR}^m \rightarrow \mathbb{IR}$, \mathbb{IR} being the set of intervals with real bounds, which consists in replacing real operations by their interval counterparts in the expression of the function.

A generally more accurate extension relies on a linearization by the mean-value theorem. Mean-value extensions can be generalized to compute ranges that are robust to disturbances, identified as a subset of the input components. Let f be a continuously differentiable function from \mathbb{R}^m to \mathbb{R} with input decomposed as $x = (u, w) \in (\mathcal{U}, \mathcal{W}) \subseteq \mathbb{IR}^m$. We define the robust range of function f on \mathbf{x} , robust with respect to component $w \in \mathcal{W}$, as $R_{\mathcal{AE}}^f(\mathcal{U}, \mathcal{W}) = \{z \mid \forall w \in \mathcal{W}, \exists u \in \mathcal{U}, z = f(u, w)\}$.

For a continuously differentiable function $f : \mathbb{R}^m \rightarrow \mathbb{R}^n$, we note $\nabla f = (\nabla_j f_i)_{ij} = (\frac{\partial f_i}{\partial x_j})_{1 \leq i \leq n, 1 \leq j \leq m}$ its Jacobian matrix. We note $\langle x, y \rangle$ the scalar product of vectors x and y , and $|x|$ the absolute value extended componentwise. For a vector of intervals $\mathcal{X} = [\underline{\mathcal{X}}, \overline{\mathcal{X}}]$, we note $c(\mathcal{X}) = (\overline{\mathcal{X}} + \underline{\mathcal{X}})/2.0$ and $r(\mathcal{X}) = (\overline{\mathcal{X}} - \underline{\mathcal{X}})/2.0$ its center and radius defined componentwise.

Theorem 1. ([8], slightly simplified version of Thm. 2). *Let f be a continuously differentiable function from \mathbb{R}^m to \mathbb{R} and $\mathcal{X} = \mathcal{U} \times \mathcal{W} \subseteq \mathbb{IR}^m$. Let \mathcal{F}^0 ,*

$\nabla_w^{\mathcal{X}}$ and $\nabla_u^{\mathcal{X}}$ be vectors of intervals such that $c(\mathcal{X}) \subseteq \mathcal{F}^0$, $\{|\nabla_w f(u, w)|, (u, w) \in \mathcal{X}\} \subseteq \nabla_w^{\mathcal{X}}$ and $\{|\nabla_u f(u, w)|, (u, w) \in \mathcal{X}\} \subseteq \nabla_u^{\mathcal{X}}$. We have:

$$\begin{aligned} & [\overline{\mathcal{F}^0} - \langle \underline{\nabla}_u^{\mathcal{X}}, r(\mathcal{U}) \rangle + \langle \overline{\nabla}_w^{\mathcal{X}}, r(\mathcal{W}) \rangle, \underline{\mathcal{F}^0} + \langle \underline{\nabla}_u^{\mathcal{X}}, r(\mathcal{U}) \rangle - \langle \overline{\nabla}_w^{\mathcal{X}}, r(\mathcal{W}) \rangle] \subseteq R_{\mathcal{AE}}^f(\mathcal{U}, \mathcal{W}) \\ R_{\mathcal{AE}}^f(\mathcal{U}, \mathcal{W}) & \subseteq [\underline{\mathcal{F}^0} - \langle \overline{\nabla}_u^{\mathcal{X}}, r(\mathcal{U}) \rangle + \langle \underline{\nabla}_w^{\mathcal{X}}, r(\mathcal{W}) \rangle, \overline{\mathcal{F}^0} + \langle \overline{\nabla}_u^{\mathcal{X}}, r(\mathcal{U}) \rangle - \langle \underline{\nabla}_w^{\mathcal{X}}, r(\mathcal{W}) \rangle] \end{aligned}$$

Theorem 1 provides inner and outer-approximations of the robust range (or of the classical range when there is no disturbance component w) of scalar-valued functions, or of the projections on each component of vector-valued functions, using bounds on the slopes on the input set. The result is useful to compute a projected range that is robustly reachable with respect to the disturbances w , or as a brick in computing an under-approximation of the image of a vector-valued function, as stated in Theorem 3 in [8].

Note that the accuracy of the mean-value AE extension can be improved with an evaluation by a quadrature formula ([10], Sect. 4.2). Alternatively, an order 2 Taylor-based extension ([10], Sect. 3) can be used.

2.3 Reachability of Neural Network Controlled Closed-Loop Systems

The inner and outer approximations defined in Sect. 2.2 can be computed for f being a simple function, possibly involving a neural network evaluation, or f being the function defined by the iterated values of a discrete systems, or finally f being the solution flow of closed-loop system (2).

In both discrete-time and the continuous-time cases, and whether some neural network controller is present or not, the evaluation of an outer-approximation of the image of the solution and its Jacobian with respect to inputs and disturbances over sets is needed in order to apply Theorem 1.

In our work and implementation, we advocate the use of a unique abstraction by affine forms (or zonotopes for the geometric view of a tuple of variables represented by affine forms) for these sets and these evaluations, including performing reachability of the neural network controller. This abstraction is very convenient and versatile to over-approximate any smooth function, providing a good trade-off between efficiency and precision in most cases (and for more precision, one can consider extensions with e.g. polynomial zonotopes [2]).

For continuous-time systems, we use Taylor expansions in time of the solution on a time grid. To build these Taylor expansions, we evaluate function f and its (Lie) derivatives over affine forms by a combination of automatic differentiation and numerical evaluation by affine arithmetic, as described in e.g. [9]. The neural network is seen as a nonlinear function h , composed with f to build function g for which we compute the solution flow. Theorem 1 is applied to this solution flow. We build the abstraction of h and thus g by a simple propagation of affine forms by affine arithmetic in the network: linear transformers are exact, and we propagate affine forms through the activation functions seen as standard nonlinear functions relying on the elementary exponential function, $\tanh(x) = 2/(1 + e^{-2x}) - 1$ and $\text{sig}(x) = 1/(1 + e^{-x})$. For differentiating the activation functions, we use $\tanh'(x) = 1.0 - \tanh(x)^2$ and $\text{sig}'(x) = \text{sig}(x)(1 - \text{sig}(x))$.

3 Implementation

As mentioned in the introduction, RINO implements all ideas presented in [8–11] for the joint computation of inner and outer approximations of robustly reachable sets of differentiable nonlinear discrete-time [8, 10] or continuous-time systems [8, 9], possibly with constant delays [11]. For experiments with systems without neural networks, we refer to the results presented in these works, obtained with a previous version of RINO.

RINO is written in C++. Intervals and zonotopes are used for set representation: the tool relies on the FILIB++ library [19] for interval computations and the aafib library¹ for affine arithmetic [4]. Ole Stauning’s FADBAD++ library² is used for automatic differentiation: its implementation with template enables us to easily evaluate the differentiation in the set representation of our choice (affine forms or zonotopes mostly). The tool takes as inputs:

- an open-loop or closed loop system, either discrete time or continuous-time, which for now is hard-coded in C++,
- an optional neural network, provided to the tool in a format directly inspired from the format analyzed by Sherlock [6], which can be used as some inputs of the closed-loop system,
- an optional configuration file to set initial values, input and disturbances ranges, and some parameter of the analysis (such as time step, order of Taylor expansion in time)

It computes inner and outer-approximations of the projection on each component of ranges, as well as joint 2D and 3D inner-approximations (provided as yaml file and Jupyter/python-produced figures). Additionally to the classical ranges, RINO computes approximations of output ranges that are reachable robustly or adversarially with respect to disturbances, specified as a subset of inputs. In the experiments presented hereafter, we consider examples only of classical reachability, for which comparisons with existing work are available, but the extension to robust reachability based on our previous work is straightforward.

4 Experiments

For space reasons, we focus here on the main novelty which is the extension of this previous work to compute under and over-approximations of (robust) reachable sets of neural network controlled systems (2).

Choice of Tools and Benchmark Examples. We compare RINO against ReachNN* and Verisig 2.0 that are the most recent fully-fledged reachability analyzers for neural network based control systems, and for which comparisons with other tools on classical benchmarks are well documented in e.g. [15]. They both improve on previous versions, Verisig and ReachNN, and on state of the art

¹ <http://aafib.sourceforge.net>.

² <http://www.fadbad.com>.

tools Sherlock, also based on Flow*, and NNV. As noted in e.g. [15]: “Firstly, note that Verisig takes significantly more time to compute reachable sets (21 times slower in the case of the B5 benchmark). Furthermore, Verisig is unable to verify some properties due to increasing error. Note that NNV is unable to verify any of the properties considered in this paper due to high approximation error.”. Remark though that there has been some amelioration to the internal solvers used in NNV which should qualify the latter statement (see e.g. [16]). We do not compare with the implementation in JuliaReach [21] since, first, timings are difficult to compare with an interpreted framework, and second, because it would require mixing several tools together, with many potential combinations. We try to provide elements of comparison with POLAR [12], but in many ways the latter addresses a different problem, with the emphasis on being able to interpret e.g. ReLU activation functions.

Table 1. List of benchmarks (see [15])

Name	Dynamics	Initial set	Horizon	Control step
Mountain Car	$\dot{x}_1 = x_2$ $\dot{x}_2 = 0.0015u - 0.0025 \cos(3x_1)$	$[-0.5, -0.48]$ $[0, 0.001]$	$T = 75$	1
discrete MC (stepsize 1)	$x_1^{n+1} = x_1^n + x_2^n$ $x_2^{n+1} = x_2^n + 0.0015u^n$ $-0.0025 \cos(3x_1^n)$	$[-0.5, -0.48]$ $[0, 0.001]$	$T = 75$	1
TORA	$\dot{x}_1 = x_2$ $\dot{x}_2 = -x_1 + 0.1 * \sin(x_3)$ $\dot{x}_3 = x_4$ $\dot{x}_4 = u$	$[-0.77, -0.75]$ $[-0.45, -0.43]$ $[0.51, 0.54]$ $[-0.3, -0.28]$	$T = 5$	0.1
ACC	$\dot{x}_1 = x_2, \dot{x}_4 = x_5$ $\dot{x}_2 = x_3, \dot{x}_5 = x_6$ $\dot{x}_3 = -4 - 0.0001x_2^2 - 2x_3$ $\dot{x}_6 = 2u - 0.0001x_5^2 - 2x_6$	$x_1 = [90, 91]$ $x_2 = [32, 32.05]$ $x_4 = [10, 11]$ $x_5 = [30, 30.05]$	$T = 5$	0.1
B1 (Ex 1 in [7])	$\dot{x}_1 = x_2$ $\dot{x}_2 = ux_2^2 - x_1$	$[0.8, 0.9]$ $[0.5, 0.6]$	$T = 7$	0.2
B2 (Ex 2 in [7])	$\dot{x}_1 = x_2 - x_1^3$ $\dot{x}_2 = u$	$[0.7, 0.9]$ $[0.7, 0.9]$	$T = 1.8$	0.2
B3 (Ex 3 in [7])	$\dot{x}_1 = -x_1(0.1 + (x_1 + x_2)^2)$ $\dot{x}_2 = (u + x_1)(0.1 + (x_1 + x_2)^2)$	$[0.8, 0.9]$ $[0.4, 0.5]$	$T = 6$	0.1
B4 (Ex 4 in [7])	$\dot{x}_1 = -x_1 + x_2 - x_3$ $\dot{x}_2 = -x_1(x_3 + 1) - x_2$ $\dot{x}_3 = -x_1 + u$	$[0.25, 0.27]$ $[0.08, 0.1]$ $[0.25, 0.27]$	$T = 1$	0.1
B5 (Ex 5 in [7])	$\dot{x}_1 = x_3^3 - x_2$ $\dot{x}_2 = x_3$ $\dot{x}_3 = u$	$[0.38, 0.4]$ $[0.45, 0.47]$ $[0.25, 0.27]$	$T = 2$	0.2

We use a large subset (7/10) of the examples from Verisig 2.0 [15], which are benchmarks used by most of the tools in the field, through e.g. the ARCH competition [17]. We also consider the same settings in terms of initial sets and the same time horizon. These are recalled in Table 1.

We indicate some of RINO’s reachability results on these benchmarks in Table 2, before comparing the tightness and computing times with other tools.

Table 2. RINO’s results for time step 0.05 (except Mountain Car, step 1.)

Name	over-approx	under-approx	t (s)	t docker
Mountain Car sigmoid (2×200)	$[-0.78197, -0.64704]$ $[-0.019387, -0.0093975]$	\perp	31.	40.41
Discrete MC sigmoid (2×200)	$[-0.8711, -0.68326]$ $[-0.026888, -0.01411]$	$[-0.82466, -0.7297]$ $[-0.023716, -0.017282]$	35	19.85
TORA tanh (3×20)	$[0.022471, 0.04829]$ $[-0.80790, -0.78039]$ $[-0.37201, -0.3433]$ $[0.30682, 0.33235]$	$[0.029133, 0.041776]$ $[-0.8037, -0.78452]$ \emptyset \emptyset	1.6	2.54
ACC tanh	$[229.05, 230.29]$ $[22.819, 22.868]$ $[-2.0285, -2.0284]$ $[159.88, 161.02]$ $[29.893, 30.006]$ $[-0.30836, 0.01398]$	$[229.05, 230.29]$ $[22.819, 22.868]$ $[-2.0285, -2.0284]$ $[160.03, 160.87]$ \emptyset \emptyset	6.	7.65
B1 tanh (3×20)	$[0.012957, 0.1349]$ $[0.18089, 0.23235]$	\emptyset \emptyset	0.7	0,92
B1 sigmoid (3×20)	$[0.10155, 0.15331]$ $[0.17188, 0.20041]$	$[0.12092, 0.13398]$ \emptyset	0.6	0.77
B2 sigmoid (3×20)	$[-0.12356, -0.0811]$ $[0.16682, 0.26396]$	\perp	0.2	0.21
B3 tanh	$[0.2256, 0.25296]$ $[-0.17777, -0.16092]$	$[0.23507, 0.24352]$ \emptyset	1.3	1.67
B4 tanh	$[-0.0017942, 0.010039]$ $[-0.03494, -0.02305]$ $[0.064524, 0.070953]$	\emptyset $[-0.032405, -0.02557]$ \emptyset	0.1	0,098
B5 tanh	$[-0.42399, -0.38098]$ $[0.16388, 0.17547]$ $[-0.24869, -0.23363]$	\perp	2.7	3.8

Settings. All tools, Verisig 2.0 and ReachNN* and RINO, were run without GPU support, under Ubuntu 18.04 docker, on a Mac running Mac OS Big Sur 11.2.3 on a 2.3 GHz Intel Core i9 processor with 16Gb of memory. Verisig 2.0 and ReachNN* were run with the Reproducibility Package of Verisig 2.0 [15].

For fairness of timing results, we also run RINO with docker, and the running ratios given in Table 3 are those using these docker versions. RINO was also run natively on the same Mac. The performance degradation between the two versions of RINO can be estimated from the full data given in Table 2 from none to a 40% increase (with one exception at 80%), and most between 20 and 30%. This is higher than generally observed with docker, but due to the fact that docker on Macintosh is known to perform badly when it comes to IOs, using the underlying file system. Therefore, the performance degrades more when the system is of higher dimension and have more time steps to evaluate, since RINO logs all estimated ranges for all variables in separate files.

Comparisons Results. We compare in Table 3 the running times of Verisig 2.0, ReachNN* and RINO, and volumes of their final over-approximations, more precisely the widths of the projections of each component at final time horizon.

The three tools depend on some parameters, in particular integration time steps and order of approximation. RINO does not require tuning the integration time steps and order of Taylor models so much, so we use one fixed time step of 0.05 for all examples. We use for Verisig 2.0 and ReachNN* the settings of the CAV Reproducibility package, that we suppose give good results. Verisig 2.0 and ReachNN* actually perform poorly on the same examples with a fixed time steps of 0.05 s.

We experimented RINO with different time steps. The precision is relatively stable and does not necessarily improve when decreasing the time step. Indeed, as already noted [25], the improvement in approximation by Taylor models on smaller time steps is balanced by the loss of precision due to set-based abstraction being performed more often. Note also that the analysis time does not depend linearly on the time step: the control step, which rules the frequency at which the analysis of the neural net controller has to be performed, is fixed (see Table 1) and does not depend on the integration time step.

Column 2 in Table 3 describes the relative width of the intervals given by Verisig 2.0 for each variable at the final time and for each system, with respect to the one given by RINO. Column 4 is the same, but for ReachNN*. Columns 3 and 5 give the ratio of the analysis time of Verisig 2.0 (respectively ReachNN*), with respect to the analysis time of RINO.

In all cases, RINO is much faster than both Verisig 2.0 and ReachNN*, by factors ranging from 13 to 638.5. Moreover, this includes for RINO the time to compute the inner-approximations that Verisig 2.0 and ReachNN* do not compute. ReachNN* could not analyze TORA because of lack of memory on our platform, and timed out on ACC. Finally, interpolating the timings given in Table 1 of [12], e.g. for B1 (sig), Verisig 2.0 is reported to take 47s whereas POLAR is reported to take 20s on their platform. As Verisig 2.0 took 81.33s on our platform, we can infer that RINO is most certainly much faster, with e.g. 3.62s for B1, than POLAR.

RINO's precision is of the same order as Verisig 2.0, and always better than ReachNN* by a factor of about 2 to 10. RINO is in fact even substantially more precise than Verisig 2.0 in some cases (B1 and B2 in particular).

Inner-Approximations. Let us take example B1 (with sigmoid-based controller), and suppose we have a safety property that the value of x_1 should never be bigger than 1. Figure 1a represents in filled blue region the inner-approximation, as plain black lines the bounds of the outer-approximation, and as purple dots values actually reached, obtained by trajectories for sample initial conditions. The over-approximation alone does raise a potential alarm with respect to the unsafe zone (in red), only the inner-approximation actually proves that the safety

Table 3. Precision and running time comparisons RINO [timestep=0.05] vs Verisig 2.0 [time steps of [15]] vs ReachNN* [time steps of [15]]

Example	% width Verisig2 over RINO	Ratio time Verisig2/RINO	% width ReachNN* over RINO	Ratio time ReachNN*/RINO
TORA (tanh)	117,6%	38,6	Mem full	Mem full
	98,4%			
	106,7%			
	128,0%			
TORA (sig)	115,7%	43,4	Mem full	Mem full
	68,0%			
	110,1%			
	133,3%			
ACC (tanh)	101,9%	500,8	Time out	Time out
	105,6%			
	103,3%			
	110,1%			
	105,1%			
B1 (tanh)	84,9%	88,8	96,7%	85,1
	287,8%		245,0%	
B1 (sig)	112,1%	105,4	227,8%	86,8
	140,6%		441,9%	
B2 (sig)	263,2%	77,6	408,8%	121,9
	60,4%		513,7%	
B3 (tanh)	99,5%	57,5	103,9%	81,9
	98,5%		287,3%	
B3 (sig)	99,1%	55,2	176,8%	76,4
	98,0%		1043,9%	
B4 (tanh)	105,0%	187,9	224,2%	214,6
	101,6%		130,6%	
	108,7%		896,0%	
B4 (sig)	105,4%	154,4	226,9%	173,5
	101,9%		132,3%	
	107,7%		908,9%	
B5 (tanh)	100,2%	365,3	192,6%	8,9
	99,2%		826,4%	
	100,4%		635,6%	
B5 (sig)	100,2%	360,2	192,5%	9,0
	99,1%		851,6%	
	100,4%		1437,4%	

property is falsified. We also note on this picture that the over-approximation is very tight, given that samples give almost indistinguishable ranges. Figure 1b represents the inner and outer approximations of joint range (x_1, x_2) as well as estimation by sampling. As shown by the samples, (x_1, x_2) becomes almost a 1D curve after some time, making inner approximation extremely difficult to estimate. Indeed our inner-approximation in orange is fairly precise for the first time steps, and the corresponding inner skewed boxes are rotated to match the curvy, 1D, shape of the samples. The green boxes printed on the picture are the box enclosure of the actually computed outer-approximation. Note that the inner-approximation of the projections on each component can be non-empty while having an empty joint inner range, as some approximation is committed in the joint inner range computation (as a skewed box) from the projected ranges.

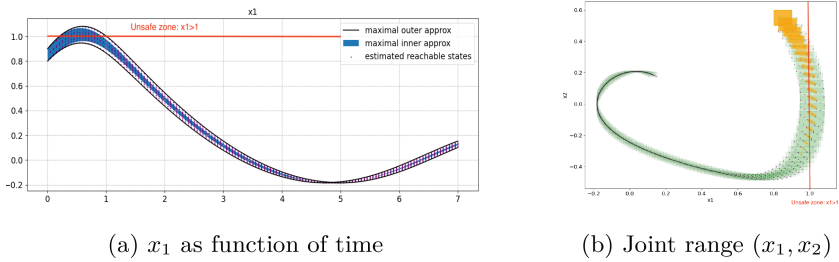


Fig. 1. B1: inner-approximation, outer-approximation and sampling (purple dots) (Color figure online)

5 Conclusion and Future Work

We presented the RINO tool, dedicated to the reachability analysis of dynamical systems, possibly controlled by neural networks. While providing accurate results, RINO is significantly faster than other state-of-the-art tools, which is key in view to address real-life reachability problems, where the systems and neural networks can be of high dimension. Moreover, as far as we are aware, it is the only existing tool to propose inner-approximations of the reachable sets of such systems. We currently handle only differentiable activation functions. We are thinking of some abstractions to handle ReLU activations as well, even though the approach is less natural in that case as it will introduce conservatism. We also plan to improve the accuracy of our current results by further specializing this work to exploit the structure of neural network, such as monotonicity of activation functions. Finally, robustness is a crucial property for neural networks enabled systems, and we plan to explore the possibilities offered by the computation of robust reachable sets.

References

1. Althoff, M.: An introduction to CORA 2015. In: Proceedings of the Workshop on Applied Verification for Continuous and Hybrid Systems (2015)
2. Althoff, M.: Reachability analysis of nonlinear systems using conservative polynomialization and non-convex sets. In: HSCC 2013, pp. 173–182. Association for Computing Machinery, New York (2013). <https://doi.org/10.1145/2461328.2461358>
3. Chen, X., Ábrahám, E., Sankaranarayanan, S.: Taylor model flowpipe construction for non-linear hybrid systems. In: 2012 IEEE 33rd Real-Time Systems Symposium, pp. 183–192 (2012)
4. Comba, J., Stolfi, J.: Affine arithmetic and its applications to computer graphics. In: SIBGRAPI (1993)
5. Dutta, S., Chen, X., Sankaranarayanan, S.: Reachability analysis for neural feedback systems using regressive rule inference (2019)
6. Dutta, S., Kushner, T., Jha, S., Sankaranarayanan, S., Shankar, N., Tiwari, A.: Sherlock: a tool for verification of deep neural networks (2019)
7. Fan, J., Huang, C., Chen, X., Li, W., Zhu, Q.: ReachNN*: a tool for reachability analysis of neural-network controlled systems. In: Hung, D.V., Sokolsky, O. (eds.) ATVA 2020. LNCS, vol. 12302, pp. 537–542. Springer, Cham (2020). https://doi.org/10.1007/978-3-030-59152-6_30
8. Goubault, E., Putot, S.: Robust under-approximations and application to reachability of non-linear control systems with disturbances. *IEEE Control Syst. Lett.* **4**(4), 928–933 (2020)
9. Goubault, E., Putot, S.: Inner and outer reachability for the verification of control systems. In: Proceedings of the 22nd ACM International Conference on Hybrid Systems: Computation and Control, HSCC 2019, Montreal, QC, Canada, 16–18 April 2019, pp. 11–22. ACM (2019)
10. Goubault, E., Putot, S.: Tractable higher-order under-approximating AE extensions for non-linear systems. In: Jungers, R.M., Ozay, N., Abate, A. (eds.) ADHS 2021. IFAC-PapersOnLine, vol. 54, pp. 235–240. Elsevier (2021)
11. Goubault, E., Putot, S., Sahlmann, L.: Inner and outer approximating flowpipes for delay differential equations. In: Chockler, H., Weissenbacher, G. (eds.) CAV 2018. LNCS, vol. 10982, pp. 523–541. Springer, Cham (2018). https://doi.org/10.1007/978-3-319-96142-2_31
12. Huang, C., Fan, J., Chen, X., Li, W., Zhu, Q.: POLAR: a polynomial arithmetic framework for verifying neural-network controlled systems (2021)
13. Huang, C., Fan, J., Li, W., Chen, X., Zhu, Q.: ReachNN: reachability analysis of neural-network controlled systems. *ACM Trans. Embed. Comput. Syst.* **18**, 1–22 (2019)
14. Ivanov, R., Weimer, J., Alur, R., Pappas, G.J., Lee, I.: Verisig: verifying safety properties of hybrid systems with neural network controllers (2019)
15. Ivanov, R., Carpenter, T., Weimer, J., Alur, R., Pappas, G., Lee, I.: Verisig 2.0: verification of neural network controllers using Taylor model preconditioning. In: Silva, A., Leino, K.R.M. (eds.) CAV 2021. LNCS, vol. 12759, pp. 249–262. Springer, Cham (2021). https://doi.org/10.1007/978-3-030-81685-8_11
16. Johnson, T.T., et al.: ARCH-COMP21 category report: artificial intelligence and neural network control systems (AINNCS) for continuous and hybrid systems plants. In: Frehse, G., Althoff, M. (eds.) 8th International Workshop on Applied Verification of Continuous and Hybrid Systems (ARCH 2021). EPiC Series in Computing, vol. 80, pp. 90–119. EasyChair (2021)

17. Johnson, T.T., et al.: ARCH-COMP20 category report: artificial intelligence and neural network control systems (AINNCS) for continuous and hybrid systems plants. In: ARCH20. 7th International Workshop on Applied Verification of Continuous and Hybrid Systems (ARCH 2020). EPiC Series in Computing, vol. 74, pp. 107–139. EasyChair (2020). <https://doi.org/10.29007/9xgv>
18. Kvasnica, M., Grieder, P., Baotić, M., Morari, M.: Multi-parametric toolbox (MPT). In: Alur, R., Pappas, G.J. (eds.) HSCC 2004. LNCS, vol. 2993, pp. 448–462. Springer, Heidelberg (2004). https://doi.org/10.1007/978-3-540-24743-2_30
19. Lerch, M., Tischler, G., von Gudenberg, J.W., Hofschuster, W., Kramer, W.: FILIB++, a fast interval library supporting containment computations. *ACM Trans. Math. Soft* **32**, 299–324 (2006)
20. Mitchell, I.M.: Comparing forward and backward reachability as tools for safety analysis. In: Bemporad, A., Bicchi, A., Buttazzo, G. (eds.) HSCC 2007. LNCS, vol. 4416, pp. 428–443. Springer, Heidelberg (2007). https://doi.org/10.1007/978-3-540-71493-4_34
21. Schilling, C., Forets, M., Guadalupe, S.: Verification of neural-network control systems by integrating Taylor models and zonotopes (2021)
22. Singh, G., Gehr, T., Mirman, M., Püschel, M., Vechev, M.T.: Fast and effective robustness certification. In: *Advances in Neural Information Processing Systems*, NeurIPS, pp. 10825–10836 (2018)
23. Tran, H.-D., et al.: Star-based reachability analysis of deep neural networks. In: ter Beek, M.H., McIver, A., Oliveira, J.N. (eds.) FM 2019. LNCS, vol. 11800, pp. 670–686. Springer, Cham (2019). https://doi.org/10.1007/978-3-030-30942-8_39
24. Tran, H.-D., et al.: NNV: the neural network verification tool for deep neural networks and learning-enabled cyber-physical systems. In: Lahiri, S.K., Wang, C. (eds.) CAV 2020. LNCS, vol. 12224, pp. 3–17. Springer, Cham (2020). https://doi.org/10.1007/978-3-030-53288-8_1
25. Wetzlinger, M., Kulmburg, A., Althoff, M.: Adaptive parameter tuning for reachability analysis of nonlinear systems. In: *Proceedings of the 24th International Conference on Hybrid Systems: Computation and Control*. HSCC 2021. Association for Computing Machinery (2021)

Open Access This chapter is licensed under the terms of the Creative Commons Attribution 4.0 International License (<http://creativecommons.org/licenses/by/4.0/>), which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license and indicate if changes were made.

The images or other third party material in this chapter are included in the chapter’s Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the chapter’s Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder.

