

University of Dayton

eCommons

---

Computer Science Faculty Publications

Department of Computer Science

---

8-2021

## Verification of Piecewise Deep Neural Networks: A Star Set Approach with Zonotope Pre-filter

Hoang-Dung Tran

Neelanjana Pal

Diego Manzananas Lopez

Patrick Musau

Xiaodong Yang

*See next page for additional authors*

Follow this and additional works at: [https://ecommons.udayton.edu/cps\\_fac\\_pub](https://ecommons.udayton.edu/cps_fac_pub)



Part of the [Graphics and Human Computer Interfaces Commons](#), and the [Other Computer Sciences Commons](#)

---

---

**Author(s)**

Hoang-Dung Tran, Neelanjana Pal, Diego Manzanas Lopez, Patrick Musau, Xiaodong Yang, Luan Viet Nguyen, Weiming Xiang, Stanley Bak, and Taylor T. Johnson (0000-0001-8021-9923)



# Verification of piecewise deep neural networks: a star set approach with zonotope pre-filter

Hoang-Dung Tran<sup>1</sup>, Neelanjana Pal<sup>2</sup>, Diego Manzananas Lopez<sup>2</sup>, Patrick Musau<sup>2</sup>,  
Xiaodong Yang<sup>2</sup>, Luan Viet Nguyen<sup>3</sup>, Weiming Xiang<sup>4</sup>, Stanley Bak<sup>5</sup> and  
Taylor T. Johnson<sup>2</sup>

<sup>1</sup>University of Nebraska Lincoln, NE, USA

<sup>2</sup>Vanderbilt University, TN, USA

<sup>3</sup>University of Dayton, OH, USA

<sup>4</sup>Augusta University, GA, USA

<sup>5</sup>Stony Brook University, NY, USA

**Abstract.** Verification has emerged as a means to provide formal guarantees on learning-based systems incorporating neural network before using them in safety-critical applications. This paper proposes a new verification approach for deep neural networks (DNNs) with piecewise linear activation functions using reachability analysis. The core of our approach is a collection of reachability algorithms using star sets (or shortly, stars), an effective symbolic representation of high-dimensional polytopes. The star-based reachability algorithms compute the output reachable sets of a network with a given input set before using them for verification. For a neural network with piecewise linear activation functions, our approach can construct both exact and over-approximate reachable sets of the neural network. To enhance the scalability of our approach, a star set is equipped with an outer-zonotope (a zonotope over-approximation of the star set) to quickly estimate the lower and upper bounds of an input set at a specific neuron to determine if splitting occurs at that neuron. This zonotope pre-filtering step reduces significantly the number of linear programming optimization problems that must be solved in the analysis, and leads to a reduction in computation time, which enhances the scalability of the star set approach. Our reachability algorithms are implemented in a software prototype called the neural network verification tool, and can be applied to problems analyzing the robustness of machine learning methods, such as safety and robustness verification of DNNs. Our experiments show that our approach can achieve runtimes twenty to 1400 times faster than Reluplex, a satisfiability modulo theory-based approach. Our star set approach is also less conservative than other recent zonotope and abstract domain approaches.

**Keywords:** Formal verification · Safety · Robustness · Neural network verification · Neural networks

## 1. Introduction

Learning-based systems using deep neural networks (DNNs) have become popular choices for solving many complicated problems in practice such as image processing, natural language translation, market prediction [LKB<sup>+</sup>17, HDY<sup>+</sup>12, LWL<sup>+</sup>17]. Recently, these systems have appeared in safety critical applications such as autonomous vehicles [BDTD<sup>+</sup>16] and air traffic collision avoidance systems [JKO18]. Utilizing NNs is promising

but may put the system at a high risk which is unacceptable for safety critical applications. Notably, it has been proved that well-trained NNs are vulnerable to adversarial attacks where a slight change in the input can lead to a huge change at the output [MDFF16]. Recent incidents in autonomous driving (e.g., Tesla [Hei20] and Uber [Muoi17]) raises an urgent need for techniques and tools that can formally verify the safety and robustness of NNs before utilizing them in safety-critical applications.

Safety verification and robustness certification of DNNs have attracted a huge attention from different communities such as machine learning [LM17, KL18, ALMP18, SGM<sup>+</sup>18, AKLP19, WPW<sup>+</sup>18b, WZC<sup>+</sup>18, ZWC<sup>+</sup>18], formal methods [PT10, KBD<sup>+</sup>17, TML<sup>+</sup>19a, XTJ19, XTJ17, HKWW17, DJST17, XTJ18], and security [GMDC<sup>+</sup>18, WPW<sup>+</sup>18b, WPW<sup>+</sup>18a], and a recent survey of the area is available [XMW<sup>+</sup>18]. Analyzing the behavior of a DNN can broadly be categorized into exact and over-approximate analyses. For the exact analysis, the SMT-based [KBD<sup>+</sup>17] and polyhedron-based approaches [TML<sup>+</sup>19a, XTJ17] are notable representatives. For the over-approximate analysis, the mixed-integer linear program (MILP) [DJST17], interval arithmetic-[WPW<sup>+</sup>18b, WPW<sup>+</sup>18a], zonotope-[SGM<sup>+</sup>18], input partition-[XTJ19], linearization-[WZC<sup>+</sup>18], and abstract-domain-[SGPV19] based are fast and efficient approaches. While the over-approximate analysis is usually faster and more scalable than the exact analysis, it guarantees only the soundness of the result. In contrast, the exact analysis is usually more time-consuming and less scalable. However, it guarantees both the soundness and completeness of the result [KBD<sup>+</sup>17]. Although the over-approximate analysis is fast and scalable, it is unclear how good the over-approximation is in term of conservativeness since the exact result is not available for comparison. Importantly, if an over-approximation approach is too conservative for neural networks with small or medium sizes, it will potentially produce huge conservative results for DNNs with a large number of layers and thousands of neurons since the over-approximation error is accumulated quickly over layers. Therefore, a scalable, exact reachability analysis is crucial not only for formal verification of DNNs, but also for estimating the conservativeness of current and up-coming over-approximation approaches.

Recently, a novel approach for verification of DNNs using star set has been proposed in [TML<sup>+</sup>19b]. Star fits perfectly for the reachability analysis of DNNs due to its following essential characteristics: (1) an efficient (exact) representation of large input sets; (2) fast and cheap affine mapping operations; (3) inexpensive intersections with half-spaces and checking empty. *By using star sets, we avoid the expensive affine mapping operation in polyhedron-based approach [TML<sup>+</sup>19a]* and thus, reduce the verification time significantly. The star set approach performs reachability analysis for feedforward DNNs layer-by-layer. For DNNs with piecewise linear activation functions, we can perform both exact analysis and over-approximate analysis. In the case of exact analysis, the output reachable set of each layer is a union of a set of stars. Based on this observation, *the star-based exact reachability algorithm naturally can be designed for efficient execution on multi-core platforms* where each layer can handle multiple input sets at the same time. In the case of over-approximate analysis, the output reachable set of each layer is a single star which can be constructed by doing *point-wise* over-approximation of the reachable set at all neurons of the layer. Notably, *the star set approach is the first approach that can visualize the exact behavior of a DNN with piecewise linear activation functions*. Additionally, by exploiting the power of parallel computing, *the star set approach is applicable for real-world applications*.

In this paper, we enrich the original star set method [TML<sup>+</sup>19b] in following directions. Firstly, we improve the computation time and scalability of the star set approach using zonotope pre-filtering. Secondly, we extend the star set approach to new classes of piecewise activation functions such as saturating linear activation function (or shortly, *satlin*), symmetric saturating linear activation functions (or shortly, *satlins*), and leaky ReLUs. Lastly, we re-evaluate the improved star set approach on safety verification of all 45 ACAS XU networks [JKO18] and robustness verification of a collection of image classification DNNs.

The improvement is based on a zonotope pre-filtering method [BTHJ20] in which a star set is equipped with an *outer-zonotope* to estimate quickly the lower and upper bounds of the star input set at specific neuron to determine if splitting occurs at that neuron without needing to solve LP optimizations. In the original star set method, these bounds are determined by solving two LP problems. Thus, if the number of neurons is large, the number of LP problems increases quickly. Moreover, in the exact analysis, the number of LP problems increases exponentially as the number of stars in the reachable set grows exponentially [TML<sup>+</sup>19b]. Therefore, minimizing the number of LP problems is crucial for scalability of the star-set approach.

We evaluate the improved reachability algorithms compared to the original star set approach [TML<sup>+</sup>19b], Reluplex [KBD<sup>+</sup>17], the zonotope approach [SGM<sup>+</sup>18], and the abstract domain approach [SGPV19]. The experimental results show that the improved algorithms is several times faster the original algorithms. It can achieve 20× to 1408× faster than Reluplex while can visualize the precise behavior of the ACAS Xu networks and can construct the complete set of counter example inputs in the case that a safety property is violated. Our over-approximate reachability algorithm is also much less conservative than the zonotope and new abstract domain approaches. It successfully verifies many safety properties of ACAS Xu networks while the zonotope and abstract domain approaches fail due to their large over-approximation errors. Our over-approximate reachability algorithm also provides a better robustness certification for image classification DNN in comparison with the zonotope and abstract domain approaches.

In summary, the main contributions of this paper are: 1) propose novel, fast and scalable methods for the exact and over-approximate reachability analysis of DNNs with popular classes of piecewise linear activation functions; 2) implement the proposed methods in NNV toolbox that is available online for evaluation and comparison; 3) provide a thorough evaluation of the new methods via real-world case studies.

## 2. Preliminaries

### 2.1. Machine learning models and symbolic verification problem

A feed-forward neural network (FNN) consists of an input layer, an output layer, and multiple hidden layers in which each layer comprises of neurons that are connected to the neurons of preceding layer labeled using weights. Given an input vector, the output of an FNN is determined by three components: the weight matrices  $W_k$ , representing the weighted connection between neurons of two consecutive layers  $k - 1$  and  $k$ , the bias vectors  $b_k$  of each layer, and the activation function  $f$  applied at each layer. Mathematically, the output of a neuron  $i$  is defined by:

$$y_i = f(\sum_{j=1}^n \omega_{ij} x_j + b_i),$$

where  $x_j$  is the  $j^{th}$  input of the  $i^{th}$  neuron,  $\omega_{ij}$  is the weight from the  $j^{th}$  input to the  $i^{th}$  neuron,  $b_i$  is the bias of the  $i^{th}$  neuron. In this paper, we are interested in FNN with ReLU activation functions defined by  $ReLU(x) = \max(0, x)$ .

**Definition 2.1** (*Reachable Set of FNN*) Given a bounded convex polyhedron input set defined as  $\mathcal{I} \triangleq \{x \mid Ax \leq b, x \in \mathbb{R}^n\}$ , and an  $k$ -layers feed-forward neural network  $F \triangleq \{L_1, \dots, L_k\}$ , the reachable set  $F(\mathcal{I}) = \mathcal{R}_{L_k}$  of the neural network  $F$  corresponding to the input set  $\mathcal{I}$  is defined incrementally by:

$$\begin{aligned} \mathcal{R}_{L_1} &\triangleq \{y_1 \mid y_1 = f_1(W_1x + b_1), x \in \mathcal{I}\}, \\ \mathcal{R}_{L_2} &\triangleq \{y_2 \mid y_2 = f_2(W_2y_1 + b_2), y_1 \in \mathcal{R}_{L_1}\}, \\ &\vdots \\ \mathcal{R}_{L_k} &\triangleq \{y_k \mid y_k = f_k(W_ky_{k-1} + b_k), y_{k-1} \in \mathcal{R}_{L_{k-1}}\}, \end{aligned}$$

where  $W_k$ ,  $b_k$  and  $f_k$  are the weight matrix, bias vector and activation function of the  $k^{\text{th}}$  layer  $L_k$ , respectively. The reachable set  $\mathcal{R}_{L_k}$  contains all outputs of the neural network corresponding to all input vectors  $x$  in the input set  $\mathcal{I}$ .

**Definition 2.2** (*Safety Verification of FNN*) Given a  $k$ -layers feed-forward neural network  $F$ , and a safety specification  $\mathcal{S}$  defined as a set of linear constraints on the neural network outputs  $\mathcal{S} \triangleq \{y_k \mid Cy_k \leq d\}$ , the neural network  $F$  is called to be safe corresponding to the input set  $\mathcal{I}$ , we write  $F(\mathcal{I}) \models \mathcal{S}$ , if and only if  $\mathcal{R}_{L_k} \cap \neg\mathcal{S} = \emptyset$ , where  $\mathcal{R}_{L_k}$  is the reachable set of the neural network with the input set  $\mathcal{I}$ , and  $\neg\mathcal{S}$  is the complement of the set  $\mathcal{S}$ . Otherwise, the neural network is called to be unsafe  $F(\mathcal{I}) \not\models \mathcal{S}$ .

## 2.2. Generalized star sets

**Definition 2.3** (*Generalized Star Set [BD17]*) A generalized star set (or simply star)  $\Theta$  is a tuple  $\langle c, V, P \rangle$  where  $c \in \mathbb{R}^n$  is the center,  $V = \{v_1, v_2, \dots, v_m\}$  is a set of  $m$  vectors in  $\mathbb{R}^n$  called basis vectors, and  $P : \mathbb{R}^m \rightarrow \{\top, \perp\}$  is a predicate. The basis vectors are arranged to form the star's  $n \times m$  basis matrix. The set of states represented by the star is given as:

$$\llbracket \Theta \rrbracket = \{x \mid x = c + \sum_{i=1}^m (\alpha_i v_i) \text{ such that } P(\alpha_1, \dots, \alpha_m) = \top\}. \quad (1)$$

Sometimes we will refer to both the tuple  $\Theta$  and the set of states  $\llbracket \Theta \rrbracket$  as  $\Theta$ . In this work, we restrict the predicates to be a conjunction of linear constraints,  $P(\alpha) \triangleq C\alpha \leq d$  where, for  $p$  linear constraints,  $C \in \mathbb{R}^{p \times m}$ ,  $\alpha$  is the vector of  $m$ -variables, i.e.,  $\alpha = [\alpha_1, \dots, \alpha_m]^T$ , and  $d \in \mathbb{R}^{p \times 1}$ . A star is an empty set, i.e.,  $\Theta = \emptyset$  if and only if the predicate  $P(\alpha)$  is infeasible. In other words, we can say the predicate polyhedron  $P(\alpha)$  is an empty set, i.e.,  $P(\alpha) = \emptyset$ .

**Proposition 2.1** Any (bounded/unbounded) convex polyhedron  $\mathcal{P} \triangleq \{x \mid Cx \leq d, x \in \mathbb{R}^n\}$  can be represented as a star.

*Proof.* The polyhedron  $\mathcal{P}$  is equivalent to the star set  $\Theta$  with the center  $c = [0, 0, \dots, 0]^T$ , the basic vectors  $V = \{e_1, e_2, \dots, e_n\}$  in which  $e_i$  is the  $i^{\text{th}}$  basic vector of  $\mathbb{R}^n$ , and the predicate  $P(\alpha) \triangleq C\alpha \leq d$ .  $\square$

**Proposition 2.2** [Affine Mapping of a Star] Given a star set  $\Theta = \langle c, V, P \rangle$ , an affine mapping of the star  $\Theta$  with the affine mapping matrix  $W$  and offset vector  $b$  defined by  $\tilde{\Theta} = \{y \mid y = Wx + b, x \in \Theta\}$  is another star with the following characteristics.

$$\tilde{\Theta} = \langle \tilde{c}, \tilde{V}, \tilde{P} \rangle, \quad \tilde{c} = Wc + b, \quad \tilde{v} = \{Wv_1, Wv_2, \dots, Wv_m\}, \quad \tilde{P} \equiv P.$$

*Proof.* From the definition of a star, we have  $\tilde{\Theta} = \{y \mid y = Wc + b + \sum_{i=1}^m (\alpha_i Wv_i)\}$ , such that  $P(\alpha_1, \dots, \alpha_m) = \top$  which implies that  $\tilde{\Theta}$  is another star with the center  $\tilde{c} = Wc + b$ , basic vectors  $\tilde{V} = \{Wv_1, Wv_2, \dots, Wv_m\}$  and the same predicate  $P$  as the original star  $\Theta$ .  $\square$

**Proposition 2.3** (Star and Half-space Intersection) The intersection of a star  $\Theta \triangleq \langle c, V, P \rangle$  and a half-space  $\mathcal{H} \triangleq \{x \mid Hx \leq g\}$  is another star with following characteristics.

$$\begin{aligned} \bar{\Theta} &= \Theta \cap \mathcal{H} = \langle \bar{c}, \bar{V}, \bar{P} \rangle, \quad \bar{c} = c, \quad \bar{V} = V, \quad \bar{P} = P \wedge P', \\ P'(\alpha) &\triangleq (H \times V_m)\alpha \leq g - H \times c, \quad V_m = [v_1 \ v_2 \ \cdots \ v_m]. \end{aligned}$$

*Proof.* For any  $x \in \Theta \cap \mathcal{H}$ , we have  $x = c + \sum_{i=1}^m (\alpha_i v_i) \wedge Hx \leq g$ , or equivalently,  $x = c + \sum_{i=1}^m (\alpha_i v_i) \wedge H \times V_m \times \alpha \leq g - H \times c$  which implies that the intersection is another star with the same center  $c$  and basic vectors  $V$  as  $\Theta$ , and an updated predicate  $\bar{P} = P \wedge P'$ ,  $P'(\alpha) \triangleq H \times V_m \times \alpha \leq g - H \times c$ .  $\square$

**Remark 1** We note that in our implementation, a generalized star set  $\Theta$  is equipped with a zonotope  $Z$  defined below. This zonotope  $Z$  is an over-approximation of the star set. It is used to estimate quickly the ranges of a state  $x[j]$  in a star set in the analysis. Mathematically, we write  $\Theta \triangleq \langle c, V, P, Z \rangle$ .

**Definition 2.4** (Zonotope) A zonotope  $Z$  is a tuple  $\langle l, G \rangle$  where  $l \in \mathbb{R}^n$  is the center,  $G = \{g_1, g_2, \dots, g_m\}$  is a set of  $m$  generators in  $\mathbb{R}^n$ . The set of states represented by a zonotope is given as:

$$Z = \{x \mid x = l + \sum_{i=1}^m (\alpha_i g_i) \text{ such that } -1 \leq \alpha_i \leq 1\}. \quad (2)$$

A zonotope is basically a star set in which all predicate variables in the ranges of  $[-1, 1]$ . The affine mapping of a zonotope is another zonotope. However, the intersection of a zonotope and a half-space generally is not a zonotope. One advantage of a zonotope compared with a star set is that we can compute quickly the ranges of a state in a zonotope without solving LP optimization. For example, the range of the state  $x(j)$  in a zonotope is:

$$l(j) - \sum_i^m |g_i(j)| \leq x(j) \leq l(j) + \sum_i^m |g_i(j)|.$$

### 3. Reachability of FNNs with ReLU activation functions

#### 3.1. Exact and complete analysis

In this section, we investigate the exact and complete analysis of FNNs with ReLU activation functions. Since any bounded convex polyhedron can be represented as a star (Proposition 2.1), we assume the input set  $\mathcal{I}$  of an FNN is a star set. From Definition 2.1, one can see that the reachable set of an FNN is derived layer-by-layer. Since the affine mapping of a star is also a star (Proposition 2.2), the core step in computing the exact reachable set of a layer with a star input set is applying the ReLU activation function on the star input set, i.e., compute  $ReLU(\Theta)$ ,  $\Theta = \langle c, V, P \rangle$ . For a layer  $L$  with  $n$  neurons, the reachable set of the layer can be computed by executing a sequence of  $n$  stepReLU operations as follows  $\mathcal{R}_L = ReLU_n(ReLU_{n-1}(\dots ReLU_1(\Theta)))$ , where  $ReLU_i(\cdot)$  denotes the stepReLU operation at  $i^{th}$  neuron.

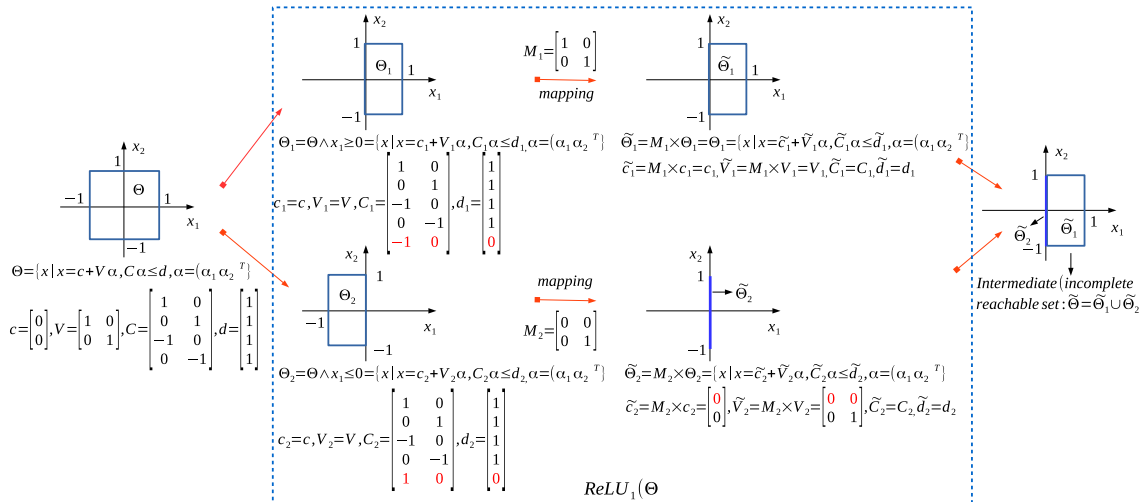


Fig. 1. An example of a stepReLU operation on a layer with two neurons

The stepReLU operation, i.e.,  $ReLU_i(\cdot)$ , works as follows. First, the input star set  $\Theta$  is decomposed into two subsets  $\Theta_1 = \Theta \wedge x_i \geq 0$  and  $\Theta_2 = \Theta \wedge x_i < 0$ . Note that from Proposition 2.3,  $\Theta_1$  and  $\Theta_2$  are also stars. Let assume that  $\Theta_1 = \langle c, V, P_1 \rangle$  and  $\Theta_2 = \langle c, V, P_2 \rangle$ . Since the later set has  $x_i < 0$ , applying the ReLU activation function on the element  $x_i$  of the vector  $x = [x_1 \cdots x_i \ x_{i+1} \cdots x_n]^T \in \Theta_2$  will lead to the new vector  $x' = [x_1 \ x_2 \cdots 0 \ x_{i+1} \cdots x_n]^T$ . This procedure is equivalent to mapping  $\Theta_2$  by the mapping matrix  $M = [e_1 \ e_2 \cdots e_{i-1} \ 0 \ e_{i+1} \cdots e_n]$ . Also, applying the ReLU activation function on the element  $x_i$  of the vector  $x \in \Theta_1$  does not change the set since we have  $x_i \geq 0$ . Consequently, the result of the stepReLU operation on input set  $\Theta$  at the  $i^{th}$  neuron is a union of two star sets  $ReLU_i(\Theta) = \langle c, V, P_1 \rangle \cup \langle Mc, MV, P_2 \rangle$ . A concrete example of the first stepReLU operation on a layer with two neurons is depicted in Fig. 1. We note that in our approach, we do not explicitly define the union operation of two stars. In reachability analysis, if a star set is split into two new stars, we say the reachable set is the union of two stars, which are simply stored in an array. If the reachable set is the union of  $N$  star sets, it is simply an array containing  $N$  star sets. By doing that, we can explore the power of parallel computing for reachability analysis of a neural network.

The number of stepReLU operation can be reduced if we know beforehand the ranges of all states in the input set. For example, if we know that  $x_i$  is always larger than zero, then we have  $ReLU_i(\Theta) = \Theta$ , or in other words, we do not need to execute the stepReLU operation on the  $i^{th}$  neuron. Therefore, to minimize the number of stepReLU operations and the computation time, we first *determine the ranges of all states in the input set by solving  $n$ -linear programming problems.*

**Lemma 3.1** *The worst-case complexity of the number of stars in the reachable set of an  $N$ -neurons FNN is  $\mathcal{O}(2^N)$ .*

*Proof.* Given a star input set, each stepReLU operation produces at most two more stars which leads to the total number of stars in the worst case of one layer is  $2^{n_L}$  where  $n_L$  is the number of neurons in the layer. For an FNN, the output reachable sets of one layer is the inputs of the next layer. Therefore, in the worst-case, the total number of stars in the reachable set of an  $k$ -layers and  $N$ -neurons FNN is  $2^{n_{L_1}} \times \cdots \times 2^{n_{L_k}} = 2^{n_{L_1} + \cdots + n_{L_k}} = 2^N$ .

**Lemma 3.2** *The worst-case complexity of the number of constraints of a star in the reachable set of an  $N$ -neuron FNN is  $\mathcal{O}(N)$ .*

*Proof.* From the stepReLU sub-procedure, we can see that given a star input set  $\Theta$ , each stepReLU operation produces one or two stars that have at most one more constraint than the star input set. Therefore, with a layer of  $n$  neurons, at most  $n$ - stepReLU operations are executed which result star reachable sets in which each one has at most  $n$  constraints more than the star input set. Consequently, the number of constraints in a star input set increases linearly over layers, and thus, the worst-case complexity of the number of constraints of a star in the reachable set of an  $N$ -neurons FNN is  $\mathcal{O}(N)$ .

**Theorem 3.1** (Verification complexity) *Let  $F$  be an  $N$ -neuron FNN,  $\Theta$  be a star set with  $p$  linear constraints and  $m$ -variables in the predicate,  $S$  be a safety specification with  $s$  linear constraints. In the worst case, the safety verification or falsification of the neural network  $F(\Theta) \models S$ ? is equivalent to solving  $2^N$  feasibility problems in which each has  $N + p + s$  linear constraints and  $m$ -variables.*

*Proof.* From Lemma 3.1, there are at most  $2^N$  stars in the reachable set of the neuron network. Also, from Lemma 3.2, each star has at most  $N + p$  constraints. To verify or falsify the safety of the neural network, we need to check if each star in the reachable set intersects with the negation (complement) of the safety specification, i.e., the unsafe region, which is a union of half-spaces. Assume that the safety specification  $S$  is a set of  $k$  linear constraints written as  $S = H_1 \cap H_2 \cdots \cap H_k$ , where  $H_j$  is the  $j^{th}$  half-space corresponding to the  $j^{th}$  constraint. Using the De Morgan's law, the unsafe region can be derived as,  $\neg S = \neg H_1 \cup \neg H_2 \cup \dots \cup \neg H_k$ . From Proposition 2.3, the intersection of each star in the reachable set with the unsafe region creates  $k$  new stars with at most  $N + p + s$  constraints. Note that the number of variables  $m$  in the predicate of a star does not change over stepReLU operations or in the intersection operation with the half-space. Therefore, the new star has  $m$ -variables and at most  $N + p + s$  linear constraints, and checking the intersection is equivalent to checking if the new star is an empty set which is a feasibility linear programming problem which can be solved efficiently in polynomial time.



**Remark 2** Although in the worst-case, the number of stars in the reachable set of an FNN is  $2^N$ , in practice, the actual number of stars is usually much smaller than the worst-case result, which enhances the applicability of the star-based exact reachability analysis for practical DNNs. For example, in the ACAS Xu case study used in the evaluation section, each network has 300 neurons. Therefore, in the worst case, the number of stars in the reachable set is  $2^{300}$ . However, depending on the size of the input set, the number of stars may vary from tens to millions instead of  $2^{300}$ .

**Theorem 3.2** (Safety and complete counter input set) *Let  $F$  be an FNN,  $\Theta = \langle c, V, P \rangle$  be a star input set,  $F(\Theta) = \bigcup_{i=1}^k \Theta_i$ ,  $\Theta_i = \langle c_i, V_i, P_i \rangle$  be the reachable set of the neural network, and  $S$  be a safety specification. Denote  $\bar{\Theta}_i = \Theta_i \cap \neg S = \langle c_i, V_i, \bar{P}_i \rangle$ ,  $i = 1, \dots, k$ . The neural network is safe if and only if  $\bar{P}_i = \emptyset$  for all  $i$ . If the neural network violates its safety property, then the complete counter input set containing all possible inputs in the input set that lead the neural network to unsafe states is  $\mathcal{C}_\Theta = \bigcup_{i=1}^k \langle c, V, \bar{P}_i \rangle$ ,  $\bar{P}_i \neq \emptyset$ .*

*Proof. Safety* The exact reachable set is a union of stars. It is trivial that the neural network is safe if and only if all stars in the reachable set do not intersect with the unsafe region, i.e.,  $\bar{\Theta}_i$  is an empty set for all  $i$ , or equivalently, the predicate  $\bar{P}_i$  is empty for all  $i$  (Definition 2.3).

**Complete counter input set** Note that all star sets in computation process are defined on the same predicate variable  $\alpha = [\alpha_1, \dots, \alpha_m]^T$  which is unchanged in the computation (only the number of constraints on  $\alpha$  changes). Therefore, when  $\bar{P}_i \neq \emptyset$ , it contains values of  $\alpha$  that makes the neural network unsafe. It is worth noticing that from the basic predicate  $P$ , new constraints are added over stepReLU operations, thus,  $\bar{P}_i$  contains all constraints of the basic predicate  $P$ . Consequently, the complete counter input set containing all possible inputs that make the neural network unsafe is defined by  $\mathcal{C}_\Theta = \bigcup_{i=1}^k \langle c, V, \bar{P}_i \rangle$ ,  $\bar{P}_i \neq \emptyset$ .

### 3.2. Over-approximate analysis

Although the exact and complete analysis can compute the exact reachable sets of a ReLU FNN, the number of stars grows exponentially with the number of layers and leads to an increase in computation cost that limits scalability. In this section, we investigate an over-approximation reachability algorithm for ReLU FNNs in which at each layer, only a single star is constructed by using the following approximation rule.

**Lemma 3.3** *For any input  $x \in [l, u]$ , the output set  $Y = \{y \mid y = \text{ReLU}(x)\}$  satisfies:*

- If  $l \geq 0$ , then  $y = x$ .
- If  $u \leq 0$ , then  $y = 0$ .
- If  $l < 0$  and  $u > 0$ , then  $Y \subset \bar{Y} = \{y \mid y \geq 0, y \leq \frac{u(x-l)}{u-l}, y \geq x\}$ .

The over-approximation rules for the ReLU activation function of different approaches are depicted in Fig. 2 which shows that our approximation rule is less conservative than the zonotope’s [SGM<sup>+</sup>18] and new abstract domain’s rules [SGPV19]. The zonotope-based approach [SGM<sup>+</sup>18] over-approximates the ReLU activation function by a minimal parallelogram while the abstract-domain approach [SGPV19] over-approximates the ReLU activation function by a triangle. Our star-based approach also over-approximates the ReLU activation function with a triangle as in the abstract-domain approach. However, the new abstract-domain approach only uses lower bound and upper bound constraints for the output  $y_i = ReLU(x_i)$  to avoid the state space explosion [SGPV19], for example, in Fig. 2, these constraints are  $y_i \geq 0$ ,  $y_i \leq u_i(x_i - l_i)/(u_i - l_i)$ . Notably, ***the zonotope and the new abstract domain approaches construct the reachable set based on estimated ranges which is usually very conservative.*** Consequently, these approaches obtain coarse a over-approximation of the actual reachable set which will be shown in the later section. ***To obtain a tighter over-approximation, our star set approach uses three constraints for the output  $y_i$  instead. Additionally, it constructs the reachable set using the ranges computed from solving LP problems.***

Similar to the exact approach, the over-approximate reachable set of a Layer with  $n$  neurons can be computed by executing a sequence of  $n$  *approximate-stepReLU* operations that work as follows. First, we compute the lower bound and upper bound of the input at the  $i^{th}$  neuron. This can be done by solving the following LP optimization problem:

$$\min(\max)(x[i]) : x[i] = c[i] + \sum_{j=1}^m \alpha_j v_j[i], \text{ subject to } C\alpha \leq d.$$

If the lower bound is not negative, the approximate-stepReLU operation returns a new intermediate reachable set which is exactly the same as its input set. If the upper bound is not positive, the approximate-stepReLU operation returns a new intermediate reachable set which is the same as its input set except the  $i^{th}$  state variable is zero. If the lower bound is negative and the upper bound is positive, the approximate-stepReLU operation introduces a new variable  $\alpha_{m+1}$  to capture the over-approximation of ReLU function at the  $i^{th}$  neuron. We remind that  $m$  is the number of predicate variables of the current star input set. As a result, the obtained intermediate reachable set has one more variable and three more linear constraints in the predicate in comparison with the corresponding input set. Therefore, in the worst case, the over-approximate reachability algorithm will obtain a reachable set with  $N + m_0$  variables and  $3N + n_0$  constraints in the predicate, where  $m_0$ ,  $n_0$  respectively are the number of variables and linear constraints of the predicate of the input set and  $N$  is the total number of neurons of the FNN.

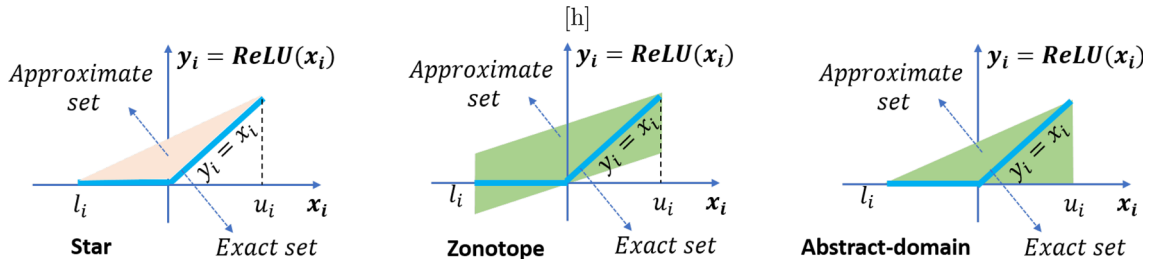


Fig. 2. The star set approach is less conservative than the zonotope [SGM<sup>+</sup>18] and new abstract-domain approaches [SGPV19]

Similar to the exact method, to verify the network’s safety, we check the intersection between the obtained over-approximate reachable set and the unsafe region (specified in the half-space form). If the over-approximate reachable set does not intersect with the unsafe region, the network is safe. If it intersects with the unsafe region, the safety of the network is unknown due to the over-approximation error. One can see that, in the over-approximate method, only soundness is guaranteed.

### 3.3. Zonotope pre-filter

One can see that, computing the ranges of all states in a star set is an important step in our star set approach that requires solving a set of LP optimization problems. In the exact analysis, the number of LP problems increase exponentially since the number of star sets grows exponentially as proved in Lemma 3.1. Therefore, to optimize the computation time and enhance the scalability of the star set approach, we need to minimize the number of LP problems solved in the analysis. Fortunately, *for exact analysis, we do not need to know the exact range of a state of the input to a specific neuron to compute the exact reachable set*. The only information we need to know is whether the range contains the zero point. If this is the case, then the star set is split into two new stars which can constructed efficiently without using the range information. If the range does not contain the zero point, the new star set is constructed even more easily. If the zero point relies on the left hand side of the range, i.e., the input is larger than zero, the output star set at the neuron is equal to the input star set. If the zero point relies on the right hand side of the range, i.e., the input is smaller than zero, the output star set is a projection of the input star set in which the neuron output is projected to zero. *In the over-approximate analysis, the range information is needed to construct an over-approximate reachable set at a specific neuron if and only if it contains the zero point*.

Based on above important observation, *we propose a zonotope pre-filtering step [BTHJ20] in which a star set is equipped with an outer-zonotope which is an over-approximation of the star set*. This outer-zonotope helps to estimate quickly the range of a state in a star set when doing reachability analysis as a zonotope is efficient for this task. Using this estimated range, we neglect all neurons in the layer that do not affect the analysis. If the estimated range contains the zero point, we solve two LP problems to get the actual range of this specific input. We reexamine whether the actual range contain the zero point to perform an appropriate operation, i.e., splitting the input set in the exact analysis or constructing an over-approximation of the reachable set using the range information in the over-approximate analysis. We note that, *the new constructed star set inherits the outer-zonotope from its preceptor. Importantly, this outer-zonotope is also updated through the analysis*. We note that the zonotope pre-filtering technique has been introduced in [BTHJ20] for only the exact reachability analysis. In this paper, we extend this technique to the approximate reachability analysis.

### 3.4. Reachability algorithms (code)

The improved star-based exact reachability algorithm using zonotope pre-filtering given in Algorithm 3.1 works as follows. The layer takes the star output sets of the preceding layer as input sets  $I = [\Theta_1, \dots, \Theta_N]$ . The main procedure in the algorithm is *layerReach* which processes the input sets  $I$  in parallel. On each input element  $\Theta_i = \langle c_i, V_i, P_i, Z_i \rangle$ , the main procedure maps the element with the layer weight matrix  $W$  and bias vector  $b$  which results a new star  $I_1 = \langle Wc_i + b, WV_i, P_i, WZ_i \rangle$ , where  $Z_i$  is the outer-zonotope of the star input set. The reachable set of the layer corresponding to the element  $\Theta_i$  is computed by *reachReLU* sub-procedure which executes a *minimized sequence* of *stepReLU/approxStepReLU* operations on the new star  $I_1$ , i.e., iteratively calls *stepReLU/approxStepReLU* sub-procedure. Note that that the *stepReLU* sub-procedure is designed to handle multiple star input sets since the number of star sets may increase after each *stepReLU* operation.

**Remark 3** The star-based reachability analysis algorithm is much faster and more reliable than the polyhedron-based algorithm [TML<sup>+</sup>19a, Ehl17] because the affine mapping step in reachable set computation can be done efficiently by matrix-vector multiplications while in the polyhedron-based approach, this step is very expensive especially for a layer with a large number of neurons since it may need to compute all vertices of the polyhedron input set [KGBM04].

## 4. Dealing with other piecewise activation functions

The star set method can be extended to FNNs with other classes of piecewise activation functions such as *satlin*, *satlins* and *leaky ReLU*. In this section, we present the extension of the star set method for these type of activation functions. Similar to a *ReLU* layer, a reachable set of a layer with *satlin*, *satlins*, and *leaky ReLU* can be constructed by performing a sequence of step reachability operations. These step operations can produce an exact or over-approximate reachable set at specific neurons.

### 4.1. Reachability of a *satlin* layer

The *satlin* activation function is defined by:

$$f(x) = \begin{cases} 0 & \text{if } x \leq 0 \\ x & \text{if } 0 \leq x \leq 1 \\ 1 & \text{if } x \geq 1 \end{cases}$$

The star-based reachability algorithm with zonotope pre-filtering for a *satlin* layer is given in Algorithm 4.2. Similar to a *ReLU* layer, we use a zonotope pre-filter to determine all neurons where splitting cannot happen. We update quickly the reachable set at these neurons, i.e., project the output to zero or one. Then, we consider the neurons where the splitting may occur. For those neurons, we solve LP optimization to determine their actual ranges to split the input set (in the exact analysis) or to construct an over-approximate reachable set (in the over approximate analysis).

### 4.2. Reachability of a *satlins* layer

The *satlins* activation function is defined by:

$$f(x) = \begin{cases} -1 & \text{if } x \leq -1 \\ x & \text{if } -1 \leq x \leq 1 \\ 1 & \text{if } x \geq 1 \end{cases}$$

The reachability algorithms for a *satlins* layer is given in Algorithm 4.3.

**Algorithm 3.1** Improved star-based reachability with zonotope pre-filtering.

---

**Input:**  $I = [\Theta_1 \dots \Theta_N]$ ,  $W$ ,  $b$   $\triangleright$  star input sets, weight matrix, bias vector  
**Output:**  $R$   $\triangleright$  exact reachable set

```

1: procedure  $R = \text{LAYERREACH}(I, W, b, \text{method})$ 
2:    $R = \emptyset$ 
3:   parfor  $i = 1 : N$  do  $\triangleright$  parallel for loop
4:      $I_1 = W * \Theta_i + b = \langle Wc_i + b, WV_i, P_i, WZ_i \rangle$ 
5:      $R_1 = \text{reachReLU}(I_1, \text{method})$ ,  $R = R \cup R_1$ 
6:   end parfor
7: procedure  $R_1 = \text{REACHReLU}(I_1, \text{method})$ 
8:    $In = I_1$ 
9:    $[lb, ub] = In.Z.\text{getRanges}$   $\triangleright$  estimate ranges of all input variables
10:   $\text{map} = \text{find}(ub < 0)$   $\triangleright$  list of neglected neurons
11:   $In.c(\text{map}, 1) = 0$ ,  $In.V(\text{map}, :) = 0$   $\triangleright$  update the input star set
12:   $In.Z.l(\text{map}, 1) = 0$ ,  $In.Z.G(\text{map}, :) = 0$   $\triangleright$  update the outer-zonotope
13:   $\text{map} = \text{find}(lb < 0 \ \& \ ub > 0)$   $\triangleright$  construct computation map
14:   $m = \text{length}(\text{map})$   $\triangleright$  minimized number of step operations
15:  for  $i = 1 : m$  do
16:    if  $\text{method} = \text{exact}$  then
17:       $In = \text{stepReLU}(In, \text{map}(i))$   $\triangleright$  stepReLU operation
18:    else if  $\text{method} = \text{approx}$  then
19:       $In = \text{approxStepReLU}(In, \text{map}(i))$   $\triangleright$  approxStepReLU operation
20:     $R_1 = In$ 
21:  procedure  $\tilde{R} = \text{STEPReLU}(\tilde{I}, i)$ 
22:     $\tilde{R} = \emptyset$ ,  $\tilde{I} = [\tilde{\Theta}_1 \dots \tilde{\Theta}_k]$   $\triangleright$  intermediate star input and output sets
23:    for  $j = 1 : k$  do
24:       $[lb_i, ub_i] = \tilde{\Theta}_j.\text{getRange}(i)$   $\triangleright$  get exact range of the  $j^{\text{th}}$  input
25:       $R_1 = \emptyset$ ,  $M = [e_1 \ e_2 \ \dots \ e_{i-1} \ 0 \ e_{i+1} \ \dots \ e_n]$ 
26:      if  $lb_i \geq 0$  then  $R_1 = \tilde{\Theta}_j = \langle \tilde{c}_j, \tilde{V}_j, \tilde{P}_j, \tilde{Z}_j \rangle$ 
27:      if  $ub_i \leq 0$  then  $R_1 = M * \tilde{\Theta}_j = \langle M\tilde{c}_j, M\tilde{V}_j, \tilde{P}_j, M\tilde{Z}_j \rangle$ 
28:      if  $lb_i < 0 \ \& \ ub_i > 0$  then
29:         $\tilde{\Theta}'_j = \tilde{\Theta}_j \wedge x[i] \geq 0 = \langle \tilde{c}_j, \tilde{V}_j, \tilde{P}'_j, \tilde{Z}_j \rangle$ ,
30:         $\tilde{\Theta}''_j = \tilde{\Theta}_j \wedge x[i] < 0 = \langle \tilde{c}_j, \tilde{V}_j, \tilde{P}''_j, \tilde{Z}_j \rangle$ 
31:         $R_1 = \tilde{\Theta}'_j \cup M * \tilde{\Theta}''_j$ 
32:       $\tilde{R} = \tilde{R} \cup R_1$ 
33:  procedure  $\tilde{R} = \text{APPROXSTEPReLU}(\tilde{I}, i)$ 
34:     $\tilde{I} = \tilde{\Theta} = \langle \tilde{c}, \tilde{V}, \tilde{P}, \tilde{Z} \rangle$ 
35:     $[l, u] = \tilde{\Theta}.\text{getRange}(i)$   $\triangleright$  get actual range of the  $i^{\text{th}}$  input
36:     $M = [e_1 \ e_2 \ \dots \ e_{i-1} \ 0 \ e_{i+1} \ \dots \ e_n]$ 
37:    if  $l \geq 0$  then  $\tilde{R} = \tilde{\Theta} = \langle \tilde{c}, \tilde{V}, \tilde{P}, \tilde{Z} \rangle$ 
38:    if  $u \leq 0$  then  $\tilde{R} = M * \tilde{\Theta} = \langle M\tilde{c}, M\tilde{V}, \tilde{P}, M\tilde{Z} \rangle$ 
39:    if  $l < 0 \ \& \ u > 0$  then
40:       $\tilde{P}(\alpha) \triangleq \tilde{C}\alpha \leq \tilde{d}$ ,  $\alpha = [\alpha_1, \alpha_2, \dots, \alpha_m]^T$   $\triangleright$  input set's predicate
41:       $\alpha' = [\alpha_1, \dots, \alpha_m, \alpha_{m+1}]^T$   $\triangleright$  new variable  $\alpha_{m+1}$ 
42:       $C_1 = [0 \ 0 \ \dots \ 0 \ -1]$ ,  $d_1 = 0$   $\triangleright \alpha_{m+1} \geq 0 \Leftrightarrow C_1\alpha' \leq d_1$ 
43:       $C_2 = [V(\tilde{i}, :) - 1]$ ,  $d_2 = -\tilde{c}[\tilde{i}]$   $\triangleright \alpha_{m+1} \geq x[\tilde{i}] \Leftrightarrow C_2\alpha' \leq d_2$ 
44:       $C_3 = [\frac{-u}{u-l} \times V(\tilde{i}, :) \ 1]$ ,  $d_3 = \frac{ul}{u-l} \times (1 - \tilde{c}[\tilde{i}])$   $\triangleright \alpha_{m+1} \leq \frac{u(x[\tilde{i}]-l)}{u-l} \Leftrightarrow C_3\alpha' \leq d_3$ 
45:       $C_0 = [\tilde{C} \ 0_{m \times 1}]$ ,  $d_0 = \tilde{d}$ 
46:       $C' = [C_0; C_1; C_2; C_3]$ ,  $d' = [d_0; d_1; d_2; d_3]$ 
47:       $P'(\alpha') \triangleq C'\alpha' \leq d'$   $\triangleright$  output set's predicate
48:       $c' = M\tilde{c}$ ,  $V' = M\tilde{V}$ ,  $V' = [V' \ e_i]$   $\triangleright y[\tilde{i}] = \text{ReLU}(x[\tilde{i}]) = \alpha_{m+1}$ 
49:       $\tilde{R} = \langle c', V', P', \tilde{Z} \rangle$ 

```

---

**Algorithm 4.2** Improved star-based reachability for a satlin layer.**Input:**  $I = [\Theta_1 \dots \Theta_N]$ ,  $W$ ,  $b$   $\triangleright$  star input sets, weight matrix, bias vector**Output:**  $R$   $\triangleright$  exact reachable set

```

1: procedure  $R = \text{LAYERREACH}(I, W, b, \text{method})$ 
2:    $R = \emptyset$ 
3:   parfor  $i = 1 : N$  do  $\triangleright$  parallel for loop
4:      $I_1 = W * \Theta_i + b = \langle Wc_i + b, WV_i, P_i, WZ_i \rangle$ 
5:      $R_1 = \text{reachSatlin}(I_1, \text{method})$ ,  $R = R \cup R_1$ 
6:   end parfor
7:   procedure  $R_1 = \text{REACHSATLIN}(I_1, \text{method})$ 
8:      $In = I_1$ 
9:      $[lb, ub] = In.Z.\text{getRanges}$   $\triangleright$  estimate ranges of all input variables
10:     $\text{map} = \text{find}(ub \leq 0)$   $\triangleright$  list of neglected neurons
11:     $In.c(\text{map}, 1) = 0$ ,  $In.V(\text{map}, :) = 0$   $\triangleright$  update the input star set
12:     $In.Z.l(\text{map}, 1) = 0$ ,  $In.Z.G(\text{map}, :) = 0$   $\triangleright$  update the outer-zonotope
13:     $\text{map} = \text{find}(lb \geq 1)$   $\triangleright$  list of neglected neurons
14:     $In.c(\text{map}, 1) = 1$ ,  $In.V(\text{map}, :) = 0$   $\triangleright$  update the input star set
15:     $In.Z.l(\text{map}, 1) = 1$ ,  $In.Z.G(\text{map}, :) = 0$   $\triangleright$  update the outer-zonotope
16:     $\text{map} = \text{find}(lb < 1 \ || \ ub > 0)$   $\triangleright$  construct computation map
17:     $m = \text{length}(\text{map})$   $\triangleright$  minimized number of step operations
18:    for  $i = 1 : m$  do
19:      if  $\text{method} = \text{exact}$  then
20:         $In = \text{stepSatlin}(In, \text{map}(i))$   $\triangleright$  stepSatlin operation
21:      else if  $\text{method} = \text{approx}$  then
22:         $In = \text{approxStepSatlin}(In, \text{map}(i))$   $\triangleright$  approxStepSatlin operation
23:       $R_1 = In$ 
24:    procedure  $\tilde{R} = \text{STEPSATLIN}(\tilde{I}, i)$ 
25:     $\tilde{R} = \emptyset$ ,  $\tilde{I} = [\tilde{\Theta}_1 \dots \tilde{\Theta}_k]$   $\triangleright$  intermediate star input and output sets
26:    for  $j = 1 : k$  do
27:       $[lb_i, ub_i] = \tilde{\Theta}_j.\text{getRange}(i)$   $\triangleright$  get exact range of the  $j^{\text{th}}$  input
28:       $R_1 = \emptyset$ ,  $M = [e_1 \ e_2 \ \dots \ e_{i-1} \ 0 \ e_{i+1} \ \dots \ e_n]$ 
29:      if  $lb_i \geq 0$  &  $ub_i \leq 1$  then  $R_1 = \tilde{\Theta}_j = \langle \tilde{c}_j, \tilde{V}_j, \tilde{P}_j, \tilde{Z}_j \rangle$ 
30:      if  $lb_i \geq 1$  then
31:         $\tilde{c}'_j = \tilde{c}_j$ ,  $\tilde{c}'_j(i) = 1$ ,  $\tilde{V}'_j = \tilde{V}_j$ ,  $\tilde{V}'_j(i, :) = 0$ 
32:         $\tilde{Z}'_j = \tilde{Z}_j$ ,  $\tilde{Z}'_j.l(i) = 1$ ,  $\tilde{Z}'_j.G(i, :) = 0$ 
33:         $R_1 = \langle \tilde{c}'_j, \tilde{V}'_j, \tilde{P}'_j, \tilde{Z}'_j \rangle$ 
34:      if  $ub_i \leq 0$  then  $R_1 = M * \tilde{\Theta}_j = \langle M\tilde{c}_j, M\tilde{V}_j, \tilde{P}_j, M\tilde{Z}_j \rangle$ 
35:      if  $lb_i < 0$  &  $ub_i \leq 1$  then
36:         $\tilde{\Theta}'_j = \tilde{\Theta}_j \wedge x[i] \geq 0 = \langle \tilde{c}_j, \tilde{V}_j, \tilde{P}'_j, \tilde{Z}_j \rangle$ ,
37:         $\tilde{\Theta}''_j = \tilde{\Theta}_j \wedge x[i] < 0 = \langle \tilde{c}_j, \tilde{V}_j, \tilde{P}''_j, \tilde{Z}_j \rangle$ 
38:         $R_1 = \tilde{\Theta}'_j \cup M * \tilde{\Theta}''_j$ 
39:      if  $0 \leq lb_i \leq 1$  &  $ub_i \geq 1$  then
40:         $\tilde{\Theta}'_j = \tilde{\Theta}_j \wedge x[i] \leq 1 = \langle \tilde{c}_j, \tilde{V}_j, \tilde{P}'_j, \tilde{Z}_j \rangle$ ,
41:         $\tilde{\Theta}''_j = \tilde{\Theta}_j \wedge x[i] \geq 1 = \langle \tilde{c}'_j, \tilde{V}''_j, \tilde{P}''_j, \tilde{Z}''_j \rangle$ 
42:         $\tilde{c}''_j = \tilde{c}_j$ ,  $\tilde{c}''_j(i) = 1$ ,  $\tilde{V}''_j = \tilde{V}_j$ ,  $\tilde{V}''_j(i, :) = 0$ ,
43:         $\tilde{Z}''_j = \tilde{Z}_j$ ,  $\tilde{Z}''_j.l(i) = 1$ ,  $\tilde{Z}''_j.G(i, :) = 0$ 
44:         $R_1 = \tilde{\Theta}'_j \cup \tilde{\Theta}''_j$ 
45:      if  $lb_i \leq 0$  &  $ub_i \geq 1$  then
46:         $\tilde{\Theta}'_j = \tilde{\Theta}_j \wedge 0 \leq x[i] \leq 1 = \langle \tilde{c}_j, \tilde{V}_j, \tilde{P}'_j, \tilde{Z}_j \rangle$ 
47:         $\tilde{\Theta}''_j = \tilde{\Theta}_j \wedge x[i] < 0 = \langle \tilde{c}_j, \tilde{V}_j, \tilde{P}''_j, \tilde{Z}_j \rangle$ 
48:         $\tilde{\Theta}'''_j = \tilde{\Theta}_j \wedge x[i] \geq 1 = \langle \tilde{c}'''_j, \tilde{V}'''_j, \tilde{P}'''_j, \tilde{Z}'''_j \rangle$ 
49:         $\tilde{c}'''_j = \tilde{c}_j$ ,  $\tilde{c}'''_j(i) = 1$ ,  $\tilde{V}'''_j = \tilde{V}_j$ ,  $\tilde{V}'''_j(i, :) = 0$ ,
50:         $\tilde{Z}'''_j = \tilde{Z}_j$ ,  $\tilde{Z}'''_j.l(i) = 1$ ,  $\tilde{Z}'''_j.G(i, :) = 0$ 
51:         $R_1 = \tilde{\Theta}'_j \cup M * \tilde{\Theta}''_j \cup \tilde{\Theta}'''_j$ 
52:     $\tilde{R} = \tilde{R} \cup R_1$ 

```

---

```

53: procedure  $\tilde{R} = \text{APPROXSTEP SATLIN}(\tilde{I}, i)$ 
54:    $\tilde{I} = \tilde{\Theta} = \langle \tilde{c}, \tilde{V}, \tilde{P}, \tilde{Z} \rangle$ 
55:    $[l, u] = \tilde{\Theta}.getRange(i) \triangleright$  get actual range of the  $i^{\text{th}}$  input
56:    $M = [e_1 \ e_2 \ \dots \ e_{i-1} \ 0 \ e_{i+1} \ \dots \ e_n]$ 
57:   if  $l \geq 0 \ \& \ u \leq 1$  then  $\tilde{R} = \tilde{\Theta} = \langle \tilde{c}, \tilde{V}, \tilde{P}, \tilde{Z} \rangle$ 
58:   if  $l \geq 1$  then
59:      $\tilde{c}' = \tilde{c}, \tilde{c}(i) = 1, \tilde{V}' = \tilde{V}, \tilde{V}(i, :) = 0, \tilde{Z}' = \tilde{Z}, \tilde{Z}.l(i) = 1, \tilde{Z}.G(i, :) = 0$ 
60:      $\tilde{R} = \tilde{\Theta}' = \langle \tilde{c}', \tilde{V}', \tilde{P}, \tilde{Z}' \rangle$ 
61:   if  $u \leq 0$  then  $\tilde{R} = M * \tilde{\Theta} = \langle M\tilde{c}, M\tilde{V}, \tilde{P}, M\tilde{Z} \rangle$ 
62:   if  $l < 0 \ \& \ 0 < u \leq 1$  then
63:      $\tilde{P}(\alpha) \triangleq \tilde{C}\alpha \leq \tilde{d}, \alpha = [\alpha_1, \alpha_2, \dots, \alpha_m]^T \triangleright$  input set's predicate
64:      $\alpha' = [\alpha_1, \dots, \alpha_m, \alpha_{m+1}]^T \triangleright$  new variable  $\alpha_{m+1}$ 
65:      $C_1 = [0 \ 0 \ \dots \ 0 \ -1], d_1 = 0 \triangleright \alpha_{m+1} \geq 0 \Leftrightarrow C_1\alpha' \leq d_1$ 
66:      $C_2 = [V(\tilde{i}, :) - 1], d_2 = -\tilde{c}[i] \triangleright \alpha_{m+1} \geq x[i] \Leftrightarrow C_2\alpha' \leq d_2$ 
67:      $C_3 = [\frac{-u}{u-l} \times V(\tilde{i}, :) \ 1], d_3 = \frac{ul}{u-l} \times (1 - \tilde{c}[i]) \triangleright \alpha_{m+1} \leq \frac{u(x[i]-l)}{u-l} \Leftrightarrow C_3\alpha' \leq d_3$ 
68:      $C_0 = [\tilde{C} \ 0_{m \times 1}], d_0 = \tilde{d}$ 
69:      $C' = [C_0; C_1; C_2; C_3], d' = [d_0; d_1; d_2; d_3]$ 
70:      $P'(\alpha') \triangleq C'\alpha' \leq d' \triangleright$  output set's predicate
71:      $c' = M\tilde{c}, V' = M\tilde{V}, V' = [V' \ e_i] \triangleright y[i] = satlin(x[i]) = \alpha_{m+1}$ 
72:      $\tilde{R} = \langle c', V', P', \tilde{Z} \rangle$ 
73:   if  $0 \leq l < 1 \ \& \ u > 1$  then
74:      $\tilde{P}(\alpha) \triangleq \tilde{C}\alpha \leq \tilde{d}, \alpha = [\alpha_1, \alpha_2, \dots, \alpha_m]^T \triangleright$  input set's predicate
75:      $\alpha' = [\alpha_1, \dots, \alpha_m, \alpha_{m+1}]^T \triangleright$  new variable  $\alpha_{m+1}$ 
76:      $C_1 = [0 \ 0 \ \dots \ 0 \ 1], d_1 = 1 \triangleright \alpha_{m+1} \leq 1 \Leftrightarrow C_1\alpha' \leq d_1$ 
77:      $C_2 = [-V(\tilde{i}, :) \ 1], d_2 = \tilde{c}[i] \triangleright \alpha_{m+1} \leq x[i] \Leftrightarrow C_2\alpha' \leq d_2$ 
78:      $C_3 = [\frac{1-l}{u-l} \times V(\tilde{i}, :) \ -1], d_3 = \frac{l(1-l)}{u-l} \tilde{c}(i) - l \triangleright \alpha_{m+1} \geq \frac{(1-l)(x[i]-l)}{u-l} + l \Leftrightarrow C_3\alpha' \leq d_3$ 
79:      $C_0 = [\tilde{C} \ 0_{m \times 1}], d_0 = \tilde{d}$ 
80:      $C' = [C_0; C_1; C_2; C_3], d' = [d_0; d_1; d_2; d_3]$ 
81:      $P'(\alpha') \triangleq C'\alpha' \leq d' \triangleright$  output set's predicate
82:      $c' = M\tilde{c}, V' = M\tilde{V}, V' = [V' \ e_i] \triangleright y[i] = satlin(x[i]) = \alpha_{m+1}$ 
83:      $\tilde{R} = \langle c', V', P', \tilde{Z} \rangle$ 
84:   if  $l < 0 \ \& \ u > 1$  then
85:      $\tilde{P}(\alpha) \triangleq \tilde{C}\alpha \leq \tilde{d}, \alpha = [\alpha_1, \alpha_2, \dots, \alpha_m]^T \triangleright$  input set's predicate
86:      $\alpha' = [\alpha_1, \dots, \alpha_m, \alpha_{m+1}]^T \triangleright$  new variable  $\alpha_{m+1}$ 
87:      $C_1 = [0 \ 0 \ \dots \ 0 \ -1], d_1 = 0 \triangleright \alpha_{m+1} \geq 0 \Leftrightarrow C_1\alpha' \leq d_1$ 
88:      $C_2 = [0 \ 0 \ \dots \ 0 \ 1], d_2 = 1 \triangleright \alpha_{m+1} \leq 1 \Leftrightarrow C_2\alpha' \leq d_2$ 
89:      $C_3 = [-V(\tilde{i}, :) \ 1], d_3 = \frac{\tilde{c}[i]}{1-l} - \frac{l}{1-l} \triangleright \alpha_{m+1} \leq \frac{x[i]}{1-l} - \frac{l}{1-l} \Leftrightarrow C_3\alpha' \leq d_3$ 
90:      $C_4 = [\frac{1}{u} \times V(\tilde{i}, :) \ -1], d_4 = -\frac{1}{u} \tilde{c}(i) \triangleright \alpha_{m+1} \geq \frac{x[i]}{u} \Leftrightarrow C_4\alpha' \leq d_4$ 
91:      $C_0 = [\tilde{C} \ 0_{m \times 1}], d_0 = \tilde{d}$ 
92:      $C' = [C_0; C_1; C_2; C_3; C_4], d' = [d_0; d_1; d_2; d_3; d_4]$ 
93:      $P'(\alpha') \triangleq C'\alpha' \leq d' \triangleright$  output set's predicate
94:      $c' = M\tilde{c}, V' = M\tilde{V}, V' = [V' \ e_i] \triangleright y[i] = satlin(x[i]) = \alpha_{m+1}$ 
95:      $\tilde{R} = \langle c', V', P', \tilde{Z} \rangle$ 

```

---

**Algorithm 4.3** Improved star-based reachability for a satlins layer.**Input:**  $I = [\Theta_1 \dots \Theta_N]$ ,  $W$ ,  $b$   $\triangleright$  star input sets, weight matrix, bias vector**Output:**  $R$   $\triangleright$  exact reachable set

```

1: procedure  $R = \text{LAYERREACH}(I, W, b, \text{method})$ 
2:    $R = \emptyset$ 
3:   parfor  $i = 1 : N$  do  $\triangleright$  parallel for loop
4:      $I_1 = W * \Theta_i + b = \langle Wc_i + b, WV_i, P_i, WZ_i \rangle$ 
5:      $R_1 = \text{reachSatlin}(I_1, \text{method})$ ,  $R = R \cup R_1$ 
6:   end parfor
7:   procedure  $R_1 = \text{REACHSATLINS}(I_1, \text{method})$ 
8:      $In = I_1$ 
9:      $[lb, ub] = In.Z.\text{getRanges}$   $\triangleright$  estimate ranges of all input variables
10:     $map = \text{find}(ub \leq -1)$   $\triangleright$  list of neglected neurons
11:     $In.c(map, 1) = -1$ ,  $In.V(map, :) = 0$   $\triangleright$  update the input star set
12:     $In.Z.l(map, 1) = -1$ ,  $In.Z.G(map, :) = 0$   $\triangleright$  update the outer-zonotope
13:     $map = \text{find}(lb \geq 1)$   $\triangleright$  list of neglected neurons
14:     $In.c(map, 1) = 1$ ,  $In.V(map, :) = 0$   $\triangleright$  update the input star set
15:     $In.Z.l(map, 1) = 1$ ,  $In.Z.G(map, :) = 0$   $\triangleright$  update the outer-zonotope
16:     $map = \text{find}(lb < 1 \ || \ ub > -1)$   $\triangleright$  construct computation map
17:     $m = \text{length}(map)$   $\triangleright$  minimized number of step operations
18:    for  $i = 1 : m$  do
19:      if  $\text{method} = \text{exact}$  then
20:         $In = \text{stepSatlins}(In, \text{map}(i))$   $\triangleright$  stepSatlins operation
21:      else if  $\text{method} = \text{approx}$  then
22:         $In = \text{approxStepSatlins}(In, \text{map}(i))$   $\triangleright$  approxStepSatlins operation
23:       $R_1 = In$ 
24:    procedure  $\tilde{R} = \text{STEPSATLINS}(\tilde{I}, i)$ 
25:     $\tilde{R} = \emptyset$ ,  $\tilde{I} = [\tilde{\Theta}_1 \dots \tilde{\Theta}_k]$   $\triangleright$  intermediate star input and output sets
26:    for  $j = 1 : k$  do
27:       $[lb_i, ub_i] = \tilde{\Theta}_j.\text{getRange}(i)$   $\triangleright$  get exact range of the  $j^{\text{th}}$  input
28:       $R_1 = \emptyset$ ,  $M = [e_1 \ e_2 \ \dots \ e_{i-1} \ 0 \ e_{i+1} \ \dots \ e_n]$ 
29:      if  $lb_i \geq -1$  &  $ub_i \leq 1$  then  $R_1 = \tilde{\Theta}_j = \langle \tilde{c}_j, \tilde{V}_j, \tilde{P}_j, \tilde{Z}_j \rangle$ 
30:      if  $lb_i \geq 1$  then
31:         $\tilde{c}'_j = \tilde{c}_j$ ,  $\tilde{c}'_j(i) = 1$ ,  $\tilde{V}'_j = \tilde{V}_j$ ,  $\tilde{V}'_j(i, :) = 0$ 
32:         $\tilde{Z}'_j = \tilde{Z}_j$ ,  $\tilde{Z}'_j.l(i) = 1$ ,  $\tilde{Z}'_j.G(i, :) = 0$ ,  $R_1 = \tilde{\Theta}'_j = \langle \tilde{c}'_j, \tilde{V}'_j, \tilde{P}_j, \tilde{Z}'_j \rangle$ 
33:      if  $ub_i \leq -1$  then
34:         $\tilde{c}'_j = \tilde{c}_j$ ,  $\tilde{c}'_j(i) = -1$ ,  $\tilde{V}'_j = \tilde{V}_j$ ,  $\tilde{V}'_j(i, :) = 0$ 
35:         $\tilde{Z}'_j = \tilde{Z}_j$ ,  $\tilde{Z}'_j.l(i) = -1$ ,  $\tilde{Z}'_j.G(i, :) = 0$ ,  $R_1 = \tilde{\Theta}'_j = \langle \tilde{c}'_j, \tilde{V}'_j, \tilde{P}_j, \tilde{Z}'_j \rangle$ 
36:      if  $lb_i < -1$  &  $ub_i \leq 1$  then
37:         $\tilde{\Theta}'_j = \tilde{\Theta}_j \wedge x[i] \geq -1 = \langle \tilde{c}_j, \tilde{V}_j, \tilde{P}'_j, \tilde{Z}_j \rangle$ ,
38:         $\tilde{\Theta}''_j = \tilde{\Theta}_j \wedge x[i] < -1 = \langle \tilde{c}''_j, \tilde{V}''_j, \tilde{P}''_j, \tilde{Z}''_j \rangle$ 
39:         $\tilde{c}''_j = \tilde{c}_j$ ,  $\tilde{c}''_j(i) = -1$ ,  $\tilde{V}''_j = \tilde{V}_j$ ,  $\tilde{V}''_j(i, :) = 0$ 
40:         $\tilde{Z}''_j = \tilde{Z}_j$ ,  $\tilde{Z}''_j.l(i) = -1$ ,  $\tilde{Z}''_j.G(i, :) = 0$ ,  $R_1 = \tilde{\Theta}'_j \cup \tilde{\Theta}''_j$ 
41:      if  $-1 \leq lb_i \leq 1$  &  $ub_i \geq 1$  then
42:         $\tilde{\Theta}'_j = \tilde{\Theta}_j \wedge x[i] \leq 1 = \langle \tilde{c}_j, \tilde{V}_j, \tilde{P}'_j, \tilde{Z}_j \rangle$ ,
43:         $\tilde{\Theta}''_j = \tilde{\Theta}_j \wedge x[i] \geq 1 = \langle \tilde{c}''_j, \tilde{V}''_j, \tilde{P}''_j, \tilde{Z}''_j \rangle$ 
44:         $\tilde{c}''_j = \tilde{c}_j$ ,  $\tilde{c}''_j(i) = 1$ ,  $\tilde{V}''_j = \tilde{V}_j$ ,  $\tilde{V}''_j(i, :) = 0$ ,
45:         $\tilde{Z}''_j = \tilde{Z}_j$ ,  $\tilde{Z}''_j.l(i) = 1$ ,  $\tilde{Z}''_j.G(i, :) = 0$ ,  $R_1 = \tilde{\Theta}'_j \cup \tilde{\Theta}''_j$ 
46:      if  $lb_i < -1$  &  $ub_i > 1$  then
47:         $\tilde{\Theta}'_j = \tilde{\Theta}_j \wedge -1 \leq x[i] \leq 1 = \langle \tilde{c}_j, \tilde{V}_j, \tilde{P}'_j, \tilde{Z}_j \rangle$ 
48:         $\tilde{\Theta}''_j = \tilde{\Theta}_j \wedge x[i] < -1 = \langle \tilde{c}''_j, \tilde{V}''_j, \tilde{P}''_j, \tilde{Z}''_j \rangle$ 
49:         $\tilde{c}''_j = \tilde{c}_j$ ,  $\tilde{c}''_j(i) = -1$ ,  $\tilde{V}''_j = \tilde{V}_j$ ,  $\tilde{V}''_j(i, :) = 0$ ,
50:         $\tilde{Z}''_j = \tilde{Z}_j$ ,  $\tilde{Z}''_j.l(i) = -1$ ,  $\tilde{Z}''_j.G(i, :) = 0$ 
51:         $\tilde{\Theta}'''_j = \tilde{\Theta}_j \wedge x[i] \geq 1 = \langle \tilde{c}'''_j, \tilde{V}'''_j, \tilde{P}'''_j, \tilde{Z}'''_j \rangle$ 
52:         $\tilde{c}'''_j = \tilde{c}_j$ ,  $\tilde{c}'''_j(i) = 1$ ,  $\tilde{V}'''_j = \tilde{V}_j$ ,  $\tilde{V}'''_j(i, :) = 0$ ,
53:         $\tilde{Z}'''_j = \tilde{Z}_j$ ,  $\tilde{Z}'''_j.l(i) = 1$ ,  $\tilde{Z}'''_j.G(i, :) = 0$ ,  $R_1 = \tilde{\Theta}'_j \cup \tilde{\Theta}''_j \cup \tilde{\Theta}'''_j$ 
54:     $\tilde{R} = \tilde{R} \cup R_1$ 

```



---

```

55: procedure  $\tilde{R} = \text{APPROXSTEP SATLINS}(\tilde{I}, i)$ 
56:    $\tilde{I} = \tilde{\Theta} = \langle \tilde{c}, \tilde{V}, \tilde{P}, \tilde{Z} \rangle$ 
57:    $[l, u] = \tilde{\Theta}.\text{getRange}(i)$   $\triangleright$  get actual range of the  $i^{\text{th}}$  input
58:    $M = [e_1 \ e_2 \ \dots \ e_{i-1} \ 0 \ e_{i+1} \ \dots \ e_n]$ 
59:   if  $l \geq -1$  &  $u \leq 1$  then  $\tilde{R} = \tilde{\Theta} = \langle \tilde{c}, \tilde{V}, \tilde{P}, \tilde{Z} \rangle$ 
60:   if  $l \geq 1$  then
61:      $\tilde{c}' = \tilde{c}, \tilde{c}(i) = 1, \tilde{V}' = \tilde{V}, \tilde{V}(i, :) = 0, \tilde{Z}' = \tilde{Z}, \tilde{Z}.l(i) = 1, \tilde{Z}.G(i, :) = 0$ 
62:      $\tilde{R} = \tilde{\Theta}' = \langle \tilde{c}', \tilde{V}', \tilde{P}, \tilde{Z}' \rangle$ 
63:   if  $u \leq -1$  then
64:      $\tilde{c}' = \tilde{c}, \tilde{c}(i) = -1, \tilde{V}' = \tilde{V}, \tilde{V}(i, :) = 0, \tilde{Z}' = \tilde{Z}, \tilde{Z}.l(i) = -1, \tilde{Z}.G(i, :) = 0$ 
65:      $\tilde{R} = \tilde{\Theta}' = \langle \tilde{c}', \tilde{V}', \tilde{P}, \tilde{Z}' \rangle$ 
66:   if  $l < -1$  &  $0 < u \leq 1$  then
67:      $\tilde{P}(\alpha) \triangleq \tilde{C}\alpha \leq \tilde{d}, \alpha = [\alpha_1, \alpha_2, \dots, \alpha_m]^T$   $\triangleright$  input set's predicate
68:      $\alpha' = [\alpha_1, \dots, \alpha_m, \alpha_{m+1}]^T$   $\triangleright$  new variable  $\alpha_{m+1}$ 
69:      $C_1 = [0 \ 0 \ \dots \ 0 \ -1], d_1 = 1$   $\triangleright \alpha_{m+1} \geq -1 \Leftrightarrow C_1\alpha' \leq d_1$ 
70:      $C_2 = [\tilde{V}(i, :) \ -1], d_2 = -\tilde{c}(i)$   $\triangleright \alpha_{m+1} \geq x[i] \Leftrightarrow C_2\alpha' \leq d_2$ 
71:      $\triangleright \alpha_{m+1} \leq \frac{(u+1)(x[i]-u)}{u-l} + u \Leftrightarrow C_3\alpha' \leq d_3$ 
72:      $C_3 = [\frac{-u-1}{u-l} \times \tilde{V}(i, :) \ 1], d_3 = \frac{-u(u+1)}{u-l} \times \tilde{c}[i] + u$ 
73:      $C_0 = [\tilde{C} \ 0_{m \times 1}], d_0 = \tilde{d}$ 
74:      $C' = [C_0; C_1; C_2; C_3], d' = [d_0; d_1; d_2; d_3]$ 
75:      $P'(\alpha') \triangleq C'\alpha' \leq d'$   $\triangleright$  output set's predicate
76:      $c' = M\tilde{c}, V' = M\tilde{V}, V' = [V' \ e_i]$   $\triangleright y[i] = \text{satlins}(x[i]) = \alpha_{m+1}$ 
77:      $\tilde{R} = \langle c', V', P', \tilde{Z} \rangle$ 
78:   if  $-1 \leq l < 1$  &  $u > 1$  then
79:      $\tilde{P}(\alpha) \triangleq \tilde{C}\alpha \leq \tilde{d}, \alpha = [\alpha_1, \alpha_2, \dots, \alpha_m]^T$   $\triangleright$  input set's predicate
80:      $\alpha' = [\alpha_1, \dots, \alpha_m, \alpha_{m+1}]^T$   $\triangleright$  new variable  $\alpha_{m+1}$ 
81:      $C_1 = [0 \ 0 \ \dots \ 0 \ 1], d_1 = 1$   $\triangleright \alpha_{m+1} \leq 1 \Leftrightarrow C_1\alpha' \leq d_1$ 
82:      $C_2 = [-\tilde{V}(i, :) \ 1], d_2 = \tilde{c}(i)$   $\triangleright \alpha_{m+1} \leq x[i] \Leftrightarrow C_2\alpha' \leq d_2$ 
83:      $C_3 = [\frac{1-l}{u-l} \times \tilde{V}(i, :) \ -1], d_3 = \frac{l(1-l)}{u-l} \tilde{c}(i) - l$   $\triangleright \alpha_{m+1} \geq \frac{(1-l)(x[i]-l)}{u-l} + l \Leftrightarrow C_3\alpha' \leq d_3$ 
84:      $C_0 = [\tilde{C} \ 0_{m \times 1}], d_0 = \tilde{d}$ 
85:      $C' = [C_0; C_1; C_2; C_3], d' = [d_0; d_1; d_2; d_3]$ 
86:      $P'(\alpha') \triangleq C'\alpha' \leq d'$   $\triangleright$  output set's predicate
87:      $c' = M\tilde{c}, V' = M\tilde{V}, V' = [V' \ e_i]$   $\triangleright y[i] = \text{satlins}(x[i]) = \alpha_{m+1}$ 
88:      $\tilde{R} = \langle c', V', P', \tilde{Z} \rangle$ 
89:   if  $l < -1$  &  $u > 1$  then
90:      $\tilde{P}(\alpha) \triangleq \tilde{C}\alpha \leq \tilde{d}, \alpha = [\alpha_1, \alpha_2, \dots, \alpha_m]^T$   $\triangleright$  input set's predicate
91:      $\alpha' = [\alpha_1, \dots, \alpha_m, \alpha_{m+1}]^T$   $\triangleright$  new variable  $\alpha_{m+1}$ 
92:      $C_1 = [0 \ 0 \ \dots \ 0 \ -1], d_1 = 1$   $\triangleright \alpha_{m+1} \geq -1 \Leftrightarrow C_1\alpha' \leq d_1$ 
93:      $C_2 = [0 \ 0 \ \dots \ 0 \ 1], d_2 = 1$   $\triangleright \alpha_{m+1} \leq 1 \Leftrightarrow C_2\alpha' \leq d_2$ 
94:      $C_3 = [-\frac{2}{1-l} \tilde{V}(i, :) \ 1], d_3 = \frac{2\tilde{c}[i]}{1-l} - 1$   $\triangleright \alpha_{m+1} \leq \frac{2x[i]}{1-l} - 1 \Leftrightarrow C_3\alpha' \leq d_3$ 
95:      $C_4 = [\frac{2}{u+1} \times \tilde{V}(i, :) \ -1], d_4 = -\frac{2}{u+1} \tilde{c}(i) + 1$   $\triangleright \alpha_{m+1} \geq \frac{2(x[i]+1)}{u+1} - 1 \Leftrightarrow C_4\alpha' \leq d_4$ 
96:      $C_0 = [\tilde{C} \ 0_{m \times 1}], d_0 = \tilde{d}$ 
97:      $C' = [C_0; C_1; C_2; C_3; C_4], d' = [d_0; d_1; d_2; d_3; d_4]$ 
98:      $P'(\alpha') \triangleq C'\alpha' \leq d'$   $\triangleright$  output set's predicate
99:      $c' = M\tilde{c}, V' = M\tilde{V}, V' = [V' \ e_i]$   $\triangleright y[i] = \text{satlins}(x[i]) = \alpha_{m+1}$ 
100:     $\tilde{R} = \langle c', V', P', \tilde{Z} \rangle$ 

```

---

### 4.3. Reachability of a leaky ReLU layer

The leaky ReLU activation function is defined by:

$$f(\gamma, x) = \begin{cases} \gamma x & \text{if } x \leq 0 \\ x & \text{if } -1 \leq x \leq 0 \end{cases}$$

The reachability algorithms for a leaky ReLU layer is given in Algorithm 4.4. These algorithms are similar to the ones for a ReLU layer except for the case that when the input to a specific neuron is smaller than zero, the output is proportional to the input with a coefficient  $\gamma$  instead of being set by zero.

## 5. Evaluation

In this section, we evaluate the improved star set reachability algorithms compared to Marabou [KHI<sup>+</sup>19], which is an improvement of Reluplex [KBD<sup>+</sup>17], the zonotope method [SGM<sup>+</sup>18], and the abstract domain method [SGPV19] implemented in our NNV tool [TYL<sup>+</sup>20]. The implementation of the zonotope and new abstract domain methods allows us to visualize the over-approximate reachable set of these approaches to intuitively evaluate their conservativeness. All results presented in this section and their corresponding scripts are available online.<sup>1</sup>

### 5.1. Safety verification for ACAS Xu DNNs

The ACAS Xu networks are DNN-based advisory controllers that map the sensor measurements to advisories in the Airborne Collision Avoidance System X [JKO18]. It is a set of 45 feedforward neural networks that map input variables to actions for horizontal maneuvers. The output means the command to the UAV to follow, and this can be: clear of conflict (COC), weak left, weak right, strong left, or strong right. All the networks have 6 fully connected layers with a total of 300 neurons, 5 inputs and 5 outputs, with all ReLU activation functions. The inputs are:

- $\rho$ : distance from ownship to intruder (feet)
- $\theta$ : angle to intruder relative to ownship heading direction (radians)
- $\psi$ : heading angle of intruder relative to ownship heading direction (radians)
- $v_{own}$ : speed of ownship (feet per second)
- $v_{int}$ : speed of intruder (feet per second)

Two other variables,  $\tau$ , time until loss of vertical separation (seconds), and  $a_{prev}$ , previous advisory, are discretized and used to generate the 45 neural networks mentioned.

The following safety properties are used to re-evaluate the performance of different methods.

- **Property  $\phi_3$ .**
  - If the intruder is directly ahead and is moving towards the ownship, the score for COC will not be minimal.
  - The desired output property is that the score for COC is not the minimal score.
  - It has 5 input constraints:  $1500 \leq \rho \leq 1800$ ,  $\theta \leq |0.06|$ ,  $\psi \geq 3.10$ ,  $v_{own} \geq 980$ ,  $v_{int} \geq 960$ .
- **Property  $\phi_4$ .**
  - If the intruder is directly ahead and is moving away from the ownship but at a lower speed than that of the ownship, the score for COC will not be minimal.
  - The desired output property is that the score for COC is not the minimal score.
  - It has 5 input constraints:  $1500 \leq \rho \leq 1800$ ,  $\theta \leq |0.06|$ ,  $\psi = 0$ ,  $v_{own} \geq 1000$ ,  $700 \leq v_{int} \leq 800$ .

<sup>1</sup>[https://github.com/verivital/nnv/tree/faoc2021/code/nnv/examples/Submission/FM2019\\_Journal](https://github.com/verivital/nnv/tree/faoc2021/code/nnv/examples/Submission/FM2019_Journal) which is also archived via Zenodo at <https://doi.org/10.5281/zenodo.4765886>.

**Algorithm 4.4** Improved star-based reachability for a leaky ReLU layer.

---

**Input:**  $I = [\Theta_1 \cdots \Theta_N]$ ,  $W$ ,  $b$   $\triangleright$  star input sets, weight matrix, bias vector  
**Output:**  $R$   $\triangleright$  exact reachable set

```

1: procedure  $R = \text{LAYERREACH}(I, W, b, \gamma, \text{method})$ 
2:    $R = \emptyset$ 
3:   parfor  $i = 1 : N$  do  $\triangleright$  parallel for loop
4:      $I_1 = W * \Theta_i + b = \langle Wc_i + b, WV_i, P_i, WZ_i \rangle$ 
5:      $R_1 = \text{reachLeakyReLU}(I_1, \gamma, \text{method})$ ,  $R = R \cup R_1$ 
6:   end parfor
7:   procedure  $R_1 = \text{REACHLEAKYReLU}(I_1, \gamma, \text{method})$ 
8:      $In = I_1$ ,  $M_i = [e_1 \ e_2 \ \cdots \ e_{i-1} \ \gamma \times e_i \ e_{i+1} \ \cdots \ e_n]$ 
9:      $[lb, ub] = \text{In.Z.getRanges}$   $\triangleright$  estimate ranges of all input variables
10:     $\text{map} = \text{find}(ub < 0)$   $\triangleright$  list of neglected neurons
11:     $\text{In.c} = M_{\text{map}} \text{In.c}$ ,  $\text{In.V} = M_{\text{map}} \text{In.V}$   $\triangleright$  update the input star set
12:     $\text{In.Z.l} = M_{\text{map}} \text{In.Z.l}$ ,  $\text{In.Z.G} = M_{\text{map}} \text{In.Z.G}$   $\triangleright$  update the outer-zonotope
13:     $\text{map} = \text{find}(lb < 0 \ \& \ ub > 0)$   $\triangleright$  construct computation map
14:     $m = \text{length}(\text{map})$   $\triangleright$  minimized number of step operations
15:    for  $i = 1 : m$  do
16:      if  $\text{method} = \text{exact}$  then
17:         $\text{In} = \text{stepLeakyReLU}(\text{In}, \text{map}(i))$   $\triangleright$  stepLeakyReLU operation
18:      else if  $\text{method} = \text{approx}$  then
19:         $\text{In} = \text{approxStepLeakyReLU}(\text{In}, \text{map}(i))$   $\triangleright$  approxStepLeakyReLU operation
20:     $R_1 = \text{In}$ 
21:  procedure  $\tilde{R} = \text{STEPLEAKYReLU}(\tilde{I}, i, \gamma)$ 
22:     $\tilde{R} = \emptyset$ ,  $\tilde{I} = [\tilde{\Theta}_1 \cdots \tilde{\Theta}_k]$   $\triangleright$  intermediate star input and output sets
23:    for  $j = 1 : k$  do
24:       $[lb_i, ub_i] = \tilde{\Theta}_j.\text{getRange}(i)$   $\triangleright$  get exact range of the  $j^{\text{th}}$  input
25:       $R_1 = \emptyset$ ,  $M = [e_1 \ e_2 \ \cdots \ e_{i-1} \ \gamma \times e_i \ e_{i+1} \ \cdots \ e_n]$ 
26:      if  $lb_i \geq 0$  then  $R_1 = \tilde{\Theta}_j = \langle \tilde{c}_j, \tilde{V}_j, \tilde{P}_j, \tilde{Z}_j \rangle$ 
27:      if  $ub_i \leq 0$  then  $R_1 = M * \tilde{\Theta}_j = \langle M\tilde{c}_j, M\tilde{V}_j, \tilde{P}_j, M\tilde{Z}_j \rangle$ 
28:      if  $lb_i < 0 \ \& \ ub_i > 0$  then
29:         $\tilde{\Theta}'_j = \tilde{\Theta}_j \wedge x[i] \geq 0 = \langle \tilde{c}_j, \tilde{V}_j, \tilde{P}'_j, \tilde{Z}_j \rangle$ ,
30:         $\tilde{\Theta}''_j = \tilde{\Theta}_j \wedge x[i] < 0 = \langle \tilde{c}_j, \tilde{V}_j, \tilde{P}''_j, \tilde{Z}_j \rangle$ 
31:         $R_1 = \tilde{\Theta}'_j \cup M * \tilde{\Theta}''_j$ 
32:     $\tilde{R} = \tilde{R} \cup R_1$ 
33:  procedure  $\tilde{R} = \text{APPROXSTEPLEAKYReLU}(\tilde{I}, i, \gamma)$ 
34:     $\tilde{I} = \tilde{\Theta} = \langle \tilde{c}, \tilde{V}, \tilde{P}, \tilde{Z} \rangle$ 
35:     $[l, u] = \tilde{\Theta}.\text{getRange}(i)$   $\triangleright$  get actual range of the  $i^{\text{th}}$  input
36:     $M = [e_1 \ e_2 \ \cdots \ e_{i-1} \ \gamma \times e_i \ e_{i+1} \ \cdots \ e_n]$ 
37:    if  $l \geq 0$  then  $\tilde{R} = \tilde{\Theta} = \langle \tilde{c}, \tilde{V}, \tilde{P}, \tilde{Z} \rangle$ 
38:    if  $u \leq 0$  then  $\tilde{R} = M * \tilde{\Theta} = \langle M\tilde{c}, M\tilde{V}, \tilde{P}, M\tilde{Z} \rangle$ 
39:    if  $l < 0 \ \& \ u > 0$  then
40:       $\tilde{P}(\alpha) \triangleq \tilde{C}\alpha \leq \tilde{d}$ ,  $\alpha = [\alpha_1, \alpha_2, \dots, \alpha_m]^T$   $\triangleright$  input set's predicate
41:       $\alpha' = [\alpha_1, \dots, \alpha_m, \alpha_{m+1}]^T$   $\triangleright$  new variable  $\alpha_{m+1}$ 
42:       $C_1 = [\tilde{V}(i, :) - 1]$ ,  $d_1 = -\gamma\tilde{c}(i)$   $\triangleright \alpha_{m+1} \geq \gamma x[i] \Leftrightarrow C_1\alpha' \leq d_1$ 
43:       $C_2 = [\tilde{V}(i, :) - 1]$ ,  $d_2 = -\tilde{c}[i]$   $\triangleright \alpha_{m+1} \geq x[i] \Leftrightarrow C_2\alpha' \leq d_2$ 
44:       $C_3 = [\frac{-u}{u-l} \times \tilde{V}(i, :) \ 1]$ ,  $d_3 = \frac{ul}{u-l} \times (1 - \tilde{c}[i])$   $\triangleright \alpha_{m+1} \leq \frac{u(x[i]-l)}{u-l} \Leftrightarrow C_3\alpha' \leq d_3$ 
45:       $C_0 = [\tilde{C} \ 0_{m \times 1}]$ ,  $d_0 = \tilde{d}$ 
46:       $C' = [C_0; C_1; C_2; C_3]$ ,  $d' = [d_0; d_1; d_2; d_3]$ 
47:       $P'(\alpha') \triangleq C'\alpha' \leq d'$   $\triangleright$  output set's predicate
48:       $M = [e_1 \ e_2 \ \cdots \ e_{i-1} \ 0 \ e_{i+1} \ \cdots \ e_n]$ 
49:       $c' = M\tilde{c}$ ,  $V' = M\tilde{V}$ ,  $V' = [V' \ e_i]$   $\triangleright y[i] = \text{leakyReLU}(x[i]) = \alpha_{m+1}$ 
50:       $\tilde{R} = \langle c', V', P', \tilde{Z} \rangle$ 

```

---

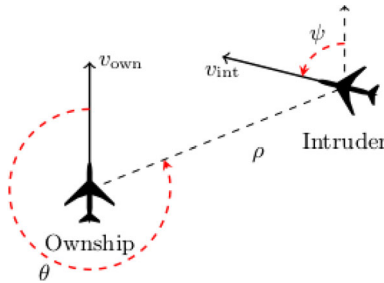


Fig. 3. Vertical view of a generic example of the ACAS Xu benchmark set.[KBD<sup>+</sup>17]

Properties  $\phi_3$  and  $\phi_4$  require that the *COC* is not the minimal score. Therefore, if *COC* is the minimal score, the networks violate their safety specification. Mathematically, the unsafe specification can be written as  $COC \leq WeakLeft \wedge COC \leq WeakRight \wedge COC \leq StrongLeft \wedge COC \leq StrongRight$ .

Our experiments are done on a laptop with the following configuration: Intel Core i7-8850H CPU @ 2.6GHz 8 core Processor, 32 GB Memory, and 64-bit Windows 10 OS.<sup>2</sup> The verification results are presented in Tables 1 and 2. We used 6 cores for the exact reachability analysis of the ACAS Xu networks using the polyhedron- and star-based approaches, and only 1 core for the over-approximate reachability analysis approaches.

**Verification results and timing performance** Safety verification using star-based reachability algorithms consists of two major steps. The first step constructs the whole reachable set of the networks. The second step checks the intersection of the constructed reachable set with the unsafe region. The verification time (VT) in our approach is the sum of the reachable set computation time (RT) and the safety checking time (ST). The reachable set computation time dominates (averagely 95% of) the verification time in all cases and the verification time varies for different properties. In the following, we use ‘ $\times$ ’ to state the timing improvement of our approach compared with Marabou or Reluplex.

**Star set approach** The experimental results show that the improved exact star method is on average  $12\times$  times faster than Marabou. It is also  $27\times$  and  $29.8\times$  times faster than Reluplex on property  $\phi_3$  and  $\phi_4$  respectively when we use parallel computing. Impressively, the approximate star method can achieve on average  $769\times$  (or  $1408\times$ ) (for property  $\phi_3$ ) and  $544\times$  (or  $961\times$ ) (for property  $\phi_4$ ) faster than Marabou (or Reluplex). This notable improvement illustrates the efficiency of star set in the reachability analysis and verification of piecewise linear DNNs where the affine mapping and half-space intersection operations can be done quickly. The improvement also come from the utilization of the zonotope pre-filtering step. In comparison with the original star set method, the improved exact method reduces the total verification time by  $5\times$  for property  $\phi_3$  and by  $2.6\times$  for property  $\phi_4$  (the original exact method verifies 45 networks with 7457 seconds for  $\phi_3$  and 1157 seconds for  $\phi_4$ ). Similarly, the improved approximate method reduces the total verification time by  $3.7\times$  for property  $\phi_3$  and by  $2.5\times$  for property  $\phi_4$  (the original approximate method verifies 45 networks with 52 seconds for  $\phi_3$  and 30 seconds for  $\phi_4$ ). We note that due to these essential characteristics of a star set, the exact star-based method is also much more efficient and scalable than the polyhedron-based approach [TML<sup>+</sup>19a] whose the verification results are not presented in this paper. For example, when verifying property  $\phi_3$ , the polyhedron method reaches time-out set as 30 minutes on more than 20 networks.

Figure 4 describes the benefits of parallel computing which can be exploited naturally with the star set approach. The figure shows that when a single core is used for verifying property  $\phi_4$  on  $N_{13}$ , our approach takes 525.8 seconds which is  $2.2\times$  faster than Reluplex (with 1150 seconds). With only 2 cores, our verification time drops quickly to 86.7 seconds which is  $13.3\times$  faster than Reluplex. When 4 cores are used, the verification time decreases to 42 seconds which is  $27.4\times$  faster than Reluplex. The figure shows that fact that when we use more cores for verification, the verification time may not be improved much. This is because the communication overhead between different cores becomes larger which affects directly to the verification time.

<sup>2</sup>In [TML<sup>+</sup>19b], this experiment was done using Amazon Web Services Elastic Computing Cloud (EC2), on a powerful m5a.24xlarge instance with 96 cores and 384 GB of memory.

**Table 1.** Verification results for property  $\phi_3$  on 45 ACAS Xu networks in which  $VT$  is the verification time in seconds. The exact-star method is on average  $12\times$  faster than Marabou ( $27\times$  faster than Reluplex which is not presented here). The approx-star is on average  $769\times$  faster than Marabou ( $1408\times$  faster than Reluplex). It can verify 29/45 networks while the zonotope can verify only 2/45 networks, and the new abstract domain approaches cannot verify any networks.

ID	Marabou		Exact-Star		Approx-Star			Zonotope			Abstract Domain			
	Res.	VT	Res.	VT	Imp.	Res.	VT	Imp.	Res.	VT	Imp.	Res.	VT	Imp.
$N_{11}$	SAFE	3636	SAFE	383.68	<b>9.5</b> $\times$	UNK	1.16	<b>3135</b> $\times$	UNK	0.04	<b>90925</b> $\times$	UNK	0.10	<b>36370</b> $\times$
$N_{12}$	SAFE	3418	SAFE	244.63	<b>14</b> $\times$	UNK	0.86	<b>3974</b> $\times$	UNK	0.02	<b>170900</b> $\times$	UNK	0.06	<b>56957</b> $\times$
$N_{13}$	SAFE	1795	SAFE	55.96	<b>32.1</b> $\times$	UNK	0.95	<b>1889</b> $\times$	UNK	0.02	<b>89750</b> $\times$	UNK	0.07	<b>25643</b> $\times$
$N_{14}$	SAFE	187	SAFE	13.68	<b>13.7</b> $\times$	SAFE	0.17	<b>1100</b> $\times$	UNK	0.02	<b>9350</b> $\times$	UNK	0.08	<b>2338</b> $\times$
$N_{15}$	SAFE	124	SAFE	18.25	<b>6.8</b> $\times$	SAFE	0.24	<b>517</b> $\times$	UNK	0.02	<b>6200</b> $\times$	UNK	0.06	<b>2067</b> $\times$
$N_{16}$	SAFE	37	SAFE	4.60	<b>8</b> $\times$	SAFE	0.09	<b>411</b> $\times$	UNK	0.02	<b>1850</b> $\times$	UNK	0.07	<b>529</b> $\times$
$N_{17}$	UNSAFE	5	UNSAFE	2.17	<b>2.3</b> $\times$	UNK	0.07	<b>71</b> $\times$	UNK	0.02	<b>250</b> $\times$	UNK	0.05	<b>100</b> $\times$
$N_{18}$	UNSAFE	2	UNSAFE	1.73	<b>1.2</b> $\times$	UNK	0.05	<b>40</b> $\times$	UNK	0.02	<b>100</b> $\times$	UNK	0.04	<b>50</b> $\times$
$N_{19}$	UNSAFE	2	UNSAFE	1.55	<b>1.3</b> $\times$	UNK	0.03	<b>67</b> $\times$	UNK	0.01	<b>200</b> $\times$	UNK	0.05	<b>40</b> $\times$
$N_{21}$	SAFE	1001	SAFE	67.63	<b>14.8</b> $\times$	UNK	0.55	<b>1820</b> $\times$	UNK	0.02	<b>50050</b> $\times$	UNK	0.06	<b>16683</b> $\times$
$N_{22}$	SAFE	480	SAFE	22.24	<b>21.6</b> $\times$	UNK	0.36	<b>1333</b> $\times$	UNK	0.02	<b>24000</b> $\times$	UNK	0.05	<b>9600</b> $\times$
$N_{23}$	SAFE	1103	SAFE	38.68	<b>28.5</b> $\times$	UNK	0.66	<b>1671</b> $\times$	UNK	0.02	<b>55150</b> $\times$	UNK	0.05	<b>22060</b> $\times$
$N_{24}$	SAFE	22	SAFE	1.71	<b>12.9</b> $\times$	SAFE	0.05	<b>440</b> $\times$	UNK	0.02	<b>1100</b> $\times$	UNK	0.05	<b>440</b> $\times$
$N_{25}$	SAFE	58	SAFE	8.45	<b>6.9</b> $\times$	SAFE	0.26	<b>223</b> $\times$	UNK	0.02	<b>2900</b> $\times$	UNK	0.06	<b>967</b> $\times$
$N_{26}$	SAFE	23	SAFE	1.64	<b>14</b> $\times$	SAFE	0.10	<b>230</b> $\times$	UNK	0.02	<b>1150</b> $\times$	UNK	0.06	<b>383</b> $\times$
$N_{27}$	SAFE	57	SAFE	4.16	<b>13.7</b> $\times$	SAFE	0.12	<b>475</b> $\times$	UNK	0.02	<b>2850</b> $\times$	UNK	0.07	<b>814</b> $\times$
$N_{28}$	SAFE	13	SAFE	1.74	<b>7.5</b> $\times$	SAFE	0.06	<b>217</b> $\times$	UNK	0.02	<b>650</b> $\times$	UNK	0.05	<b>260</b> $\times$
$N_{29}$	SAFE	2	SAFE	1.21	<b>1.7</b> $\times$	SAFE	0.02	<b>100</b> $\times$	SAFE	0.01	<b>200</b> $\times$	UNK	0.03	<b>67</b> $\times$
$N_{31}$	SAFE	97	SAFE	19.99	<b>4.9</b> $\times$	SAFE	0.22	<b>441</b> $\times$	UNK	0.02	<b>4850</b> $\times$	UNK	0.05	<b>1940</b> $\times$
$N_{32}$	SAFE	1973	SAFE	210.03	<b>9.4</b> $\times$	UNK	0.86	<b>2294</b> $\times$	UNK	0.02	<b>98650</b> $\times$	UNK	0.07	<b>28186</b> $\times$
$N_{33}$	SAFE	344	SAFE	35.06	<b>9.8</b> $\times$	SAFE	0.62	<b>555</b> $\times$	UNK	0.02	<b>17200</b> $\times$	UNK	0.05	<b>6880</b> $\times$
$N_{34}$	SAFE	106	SAFE	8.64	<b>12.3</b> $\times$	UNK	0.78	<b>136</b> $\times$	UNK	0.02	<b>5300</b> $\times$	UNK	0.07	<b>1514</b> $\times$
$N_{35}$	SAFE	41	SAFE	3.96	<b>10.4</b> $\times$	SAFE	0.13	<b>315</b> $\times$	UNK	0.02	<b>2050</b> $\times$	UNK	0.05	<b>820</b> $\times$
$N_{36}$	SAFE	120	SAFE	7.90	<b>15.2</b> $\times$	UNK	0.34	<b>353</b> $\times$	UNK	0.02	<b>6000</b> $\times$	UNK	0.07	<b>1714</b> $\times$
$N_{37}$	SAFE	8	SAFE	1.11	<b>7.2</b> $\times$	SAFE	0.04	<b>200</b> $\times$	UNK	0.02	<b>400</b> $\times$	UNK	0.06	<b>133</b> $\times$
$N_{38}$	SAFE	57	SAFE	2.87	<b>20</b> $\times$	SAFE	0.17	<b>335</b> $\times$	UNK	0.02	<b>2850</b> $\times$	UNK	0.06	<b>950</b> $\times$
$N_{39}$	SAFE	55	SAFE	4.83	<b>11.4</b> $\times$	SAFE	0.11	<b>500</b> $\times$	UNK	0.02	<b>2750</b> $\times$	UNK	0.05	<b>1100</b> $\times$
$N_{41}$	SAFE	94	SAFE	8.77	<b>10.7</b> $\times$	UNK	0.22	<b>427</b> $\times$	UNK	0.02	<b>4700</b> $\times$	UNK	0.04	<b>2350</b> $\times$
$N_{42}$	SAFE	1724	SAFE	83.21	<b>20.7</b> $\times$	UNK	0.52	<b>3315</b> $\times$	UNK	0.02	<b>86200</b> $\times$	UNK	0.06	<b>28733</b> $\times$
$N_{43}$	SAFE	1257	SAFE	121.59	<b>10.3</b> $\times$	UNK	0.82	<b>1533</b> $\times$	UNK	0.02	<b>62850</b> $\times$	UNK	0.06	<b>20950</b> $\times$
$N_{44}$	SAFE	19	SAFE	2.19	<b>8.7</b> $\times$	SAFE	0.07	<b>271</b> $\times$	UNK	0.02	<b>950</b> $\times$	UNK	0.06	<b>317</b> $\times$
$N_{45}$	SAFE	14	SAFE	1.40	<b>10</b> $\times$	SAFE	0.06	<b>233</b> $\times$	UNK	0.02	<b>700</b> $\times$	UNK	0.04	<b>350</b> $\times$
$N_{46}$	SAFE	118	SAFE	11.27	<b>10.5</b> $\times$	SAFE	1.01	<b>1168</b> $\times$	UNK	0.02	<b>5900</b> $\times$	UNK	0.07	<b>1686</b> $\times$
$N_{47}$	SAFE	54	SAFE	3.50	<b>15.4</b> $\times$	SAFE	0.09	<b>600</b> $\times$	UNK	0.02	<b>2700</b> $\times$	UNK	0.06	<b>900</b> $\times$
$N_{48}$	SAFE	36	SAFE	2.66	<b>13.5</b> $\times$	SAFE	0.07	<b>514</b> $\times$	UNK	0.02	<b>1800</b> $\times$	UNK	0.06	<b>600</b> $\times$
$N_{49}$	SAFE	33	SAFE	4.03	<b>8.2</b> $\times$	SAFE	0.19	<b>174</b> $\times$	UNK	0.02	<b>1650</b> $\times$	UNK	0.06	<b>550</b> $\times$
$N_{51}$	SAFE	764	SAFE	34.87	<b>21.9</b> $\times$	UNK	0.47	<b>1626</b> $\times$	UNK	0.02	<b>38200</b> $\times$	UNK	0.05	<b>15280</b> $\times$
$N_{52}$	SAFE	120	SAFE	7.28	<b>16.5</b> $\times$	SAFE	0.22	<b>546</b> $\times$	UNK	0.02	<b>6000</b> $\times$	UNK	0.03	<b>4000</b> $\times$
$N_{53}$	SAFE	146	SAFE	8.53	<b>17.1</b> $\times$	SAFE	0.39	<b>374</b> $\times$	UNK	0.02	<b>7300</b> $\times$	UNK	0.05	<b>2920</b> $\times$
$N_{54}$	SAFE	42	SAFE	3.02	<b>13.9</b> $\times$	SAFE	0.14	<b>300</b> $\times$	UNK	0.02	<b>2100</b> $\times$	UNK	0.06	<b>700</b> $\times$
$N_{55}$	SAFE	27	SAFE	4.36	<b>6.2</b> $\times$	SAFE	0.24	<b>113</b> $\times$	UNK	0.02	<b>1350</b> $\times$	UNK	0.06	<b>450</b> $\times$
$N_{56}$	SAFE	81	SAFE	4.75	<b>17.1</b> $\times$	SAFE	0.24	<b>338</b> $\times$	UNK	0.02	<b>4050</b> $\times$	UNK	0.06	<b>1350</b> $\times$
$N_{57}$	SAFE	5	SAFE	0.89	<b>5.6</b> $\times$	SAFE	0.02	<b>250</b> $\times$	SAFE	0.02	<b>250</b> $\times$	UNK	0.04	<b>125</b> $\times$
$N_{58}$	SAFE	157	SAFE	9.14	<b>17.2</b> $\times$	SAFE	0.18	<b>872</b> $\times$	UNK	0.02	<b>7850</b> $\times$	UNK	0.06	<b>2617</b> $\times$
$N_{59}$	SAFE	8	SAFE	1.05	<b>7.6</b> $\times$	SAFE	0.05	<b>160</b> $\times$	UNK	0.01	<b>800</b> $\times$	UNK	0.02	<b>400</b> $\times$
<b>VT</b>		<b>19466</b>		<b>1480.60</b>			<b>14.01</b>			<b>0.84</b>			<b>2.56</b>	
<b>SAFE</b>	<b>42/45</b>		<b>42/45</b>			<b>29/45</b>			<b>2/45</b>			<b>0/45</b>		
<b>Imp.</b>					<b>12</b> $\times$			<b>769</b> $\times$			<b>19622</b> $\times$			<b>6731</b> $\times$

Bold values indicate runtime improvement and summary statistics of the methods

**Table 2.** Verification results for property  $\phi_4$  on 45 ACAS Xu networks in which  $VT$  is the verification time in seconds. The exact-star method is on average  $12\times$  faster than Marbou (and  $29.8\times$  faster than Reluplex). The approx-star is on average  $544\times$  faster than Marabou (and  $961\times$  faster than Reluplex). It can verify 32/45 networks while the zonotope can verify only 1/45 networks, and the new abstract domain approaches cannot verify any networks.

ID	Marabou		Exact-Star			Approx-Star			Zonotope			Abstract Domain		
	Res.	VT	Res.	VT	Imp.	Res.	VT	Imp.	Res.	VT	Imp.	Res.	VT	Imp.
$N_{11}$	SAFE	1929.00	SAFE	76.95	<b>25</b> $\times$	UNK	0.42	<b>4643</b> $\times$	UNK	0.02	<b>111334</b> $\times$	UNK	0.04	<b>47639</b> $\times$
$N_{12}$	SAFE	2012.00	SAFE	47.54	<b>42</b> $\times$	UNK	0.54	<b>3709</b> $\times$	UNK	0.02	<b>120003</b> $\times$	UNK	0.05	<b>37812</b> $\times$
$N_{13}$	SAFE	956.00	SAFE	36.29	<b>26</b> $\times$	UNK	0.61	<b>1571</b> $\times$	UNK	0.02	<b>47960</b> $\times$	UNK	0.06	<b>16898</b> $\times$
$N_{14}$	SAFE	79.00	SAFE	3.51	<b>23</b> $\times$	UNK	0.14	<b>571</b> $\times$	UNK	0.02	<b>4455</b> $\times$	UNK	0.06	<b>1360</b> $\times$
$N_{15}$	SAFE	299.00	SAFE	26.58	<b>11</b> $\times$	UNK	0.21	<b>1434</b> $\times$	UNK	0.02	<b>18200</b> $\times$	UNK	0.03	<b>8772</b> $\times$
$N_{16}$	SAFE	201.00	SAFE	12.80	<b>16</b> $\times$	UNK	0.24	<b>850</b> $\times$	UNK	0.02	<b>12201</b> $\times$	UNK	0.05	<b>3965</b> $\times$
$N_{17}$	UNSAFE	2.00	UNSAFE	2.03	<b>1</b> $\times$	UNK	0.05	<b>37</b> $\times$	UNK	0.01	<b>141</b> $\times$	UNK	0.04	<b>45</b> $\times$
$N_{18}$	UNSAFE	2.00	UNSAFE	2.10	<b>1</b> $\times$	UNK	0.06	<b>32</b> $\times$	UNK	0.01	<b>133</b> $\times$	UNK	0.04	<b>50</b> $\times$
$N_{19}$	UNSAFE	1.00	UNSAFE	1.94	<b>1</b> $\times$	UNK	0.04	<b>23</b> $\times$	UNK	0.02	<b>65</b> $\times$	UNK	0.03	<b>39</b> $\times$
$N_{21}$	SAFE	239.00	SAFE	15.69	<b>15</b> $\times$	UNK	0.42	<b>566</b> $\times$	UNK	0.02	<b>14571</b> $\times$	UNK	0.05	<b>4535</b> $\times$
$N_{22}$	SAFE	273.00	SAFE	14.79	<b>18</b> $\times$	UNK	0.40	<b>687</b> $\times$	UNK	0.02	<b>16097</b> $\times$	UNK	0.06	<b>4615</b> $\times$
$N_{23}$	SAFE	47.00	SAFE	3.23	<b>15</b> $\times$	SAFE	0.12	<b>406</b> $\times$	UNK	0.02	<b>2871</b> $\times$	UNK	0.03	<b>1447</b> $\times$
$N_{24}$	SAFE	23.00	SAFE	3.55	<b>6</b> $\times$	SAFE	0.22	<b>106</b> $\times$	UNK	0.02	<b>1312</b> $\times$	UNK	0.03	<b>702</b> $\times$
$N_{25}$	SAFE	95.00	SAFE	11.95	<b>8</b> $\times$	SAFE	0.42	<b>226</b> $\times$	UNK	0.02	<b>5401</b> $\times$	UNK	0.05	<b>1840</b> $\times$
$N_{26}$	SAFE	71.00	SAFE	5.65	<b>13</b> $\times$	SAFE	0.37	<b>192</b> $\times$	UNK	0.02	<b>3749</b> $\times$	UNK	0.07	<b>1090</b> $\times$
$N_{27}$	SAFE	25.00	SAFE	2.47	<b>10</b> $\times$	SAFE	0.12	<b>208</b> $\times$	UNK	0.02	<b>1413</b> $\times$	UNK	0.07	<b>376</b> $\times$
$N_{28}$	SAFE	121.00	SAFE	8.41	<b>14</b> $\times$	UNK	1.24	<b>97</b> $\times$	UNK	0.02	<b>6042</b> $\times$	UNK	0.07	<b>1764</b> $\times$
$N_{29}$	SAFE	6.00	SAFE	1.29	<b>5</b> $\times$	SAFE	0.04	<b>160</b> $\times$	UNK	0.01	<b>429</b> $\times$	UNK	0.03	<b>215</b> $\times$
$N_{31}$	SAFE	220.00	SAFE	14.23	<b>15</b> $\times$	SAFE	0.32	<b>686</b> $\times$	UNK	0.02	<b>12501</b> $\times$	UNK	0.05	<b>4407</b> $\times$
$N_{32}$	SAFE	237.00	SAFE	27.10	<b>9</b> $\times$	SAFE	0.21	<b>1117</b> $\times$	UNK	0.02	<b>14725</b> $\times$	UNK	0.03	<b>7339</b> $\times$
$N_{33}$	SAFE	36.00	SAFE	3.70	<b>10</b> $\times$	SAFE	0.12	<b>308</b> $\times$	UNK	0.02	<b>2271</b> $\times$	UNK	0.03	<b>1228</b> $\times$
$N_{34}$	SAFE	32.00	SAFE	4.18	<b>8</b> $\times$	SAFE	0.13	<b>248</b> $\times$	UNK	0.02	<b>2020</b> $\times$	UNK	0.03	<b>958</b> $\times$
$N_{35}$	SAFE	212.00	SAFE	16.68	<b>13</b> $\times$	SAFE	0.71	<b>300</b> $\times$	UNK	0.02	<b>11232</b> $\times$	UNK	0.05	<b>4514</b> $\times$
$N_{36}$	SAFE	50.00	SAFE	5.57	<b>9</b> $\times$	SAFE	0.31	<b>163</b> $\times$	UNK	0.02	<b>2713</b> $\times$	UNK	0.07	<b>702</b> $\times$
$N_{37}$	SAFE	49.00	SAFE	4.07	<b>12</b> $\times$	SAFE	0.11	<b>436</b> $\times$	UNK	0.02	<b>2314</b> $\times$	UNK	0.07	<b>753</b> $\times$
$N_{38}$	SAFE	48.00	SAFE	2.98	<b>16</b> $\times$	UNK	0.28	<b>174</b> $\times$	UNK	0.02	<b>2124</b> $\times$	UNK	0.05	<b>1062</b> $\times$
$N_{39}$	SAFE	144.00	SAFE	13.04	<b>11</b> $\times$	SAFE	0.44	<b>330</b> $\times$	UNK	0.02	<b>6914</b> $\times$	UNK	0.06	<b>2307</b> $\times$
$N_{41}$	SAFE	30.00	SAFE	2.73	<b>11</b> $\times$	SAFE	0.07	<b>405</b> $\times$	SAFE	0.02	<b>1933</b> $\times$	UNK	0.03	<b>1182</b> $\times$
$N_{42}$	SAFE	77.00	SAFE	4.96	<b>16</b> $\times$	SAFE	0.32	<b>240</b> $\times$	UNK	0.02	<b>4698</b> $\times$	UNK	0.05	<b>1632</b> $\times$
$N_{43}$	SAFE	163.00	SAFE	9.79	<b>17</b> $\times$	SAFE	0.32	<b>503</b> $\times$	UNK	0.02	<b>9874</b> $\times$	UNK	0.05	<b>3491</b> $\times$
$N_{44}$	SAFE	52.00	SAFE	4.94	<b>11</b> $\times$	UNK	0.57	<b>91</b> $\times$	UNK	0.02	<b>2786</b> $\times$	UNK	0.06	<b>920</b> $\times$
$N_{45}$	SAFE	53.00	SAFE	3.91	<b>14</b> $\times$	SAFE	0.22	<b>242</b> $\times$	UNK	0.02	<b>3086</b> $\times$	UNK	0.07	<b>788</b> $\times$
$N_{46}$	SAFE	18.00	SAFE	7.96	<b>2</b> $\times$	SAFE	0.30	<b>59</b> $\times$	UNK	0.02	<b>983</b> $\times$	UNK	0.06	<b>327</b> $\times$
$N_{47}$	SAFE	16.00	SAFE	1.26	<b>13</b> $\times$	SAFE	0.09	<b>169</b> $\times$	UNK	0.02	<b>968</b> $\times$	UNK	0.06	<b>288</b> $\times$
$N_{48}$	SAFE	68.00	SAFE	6.10	<b>11</b> $\times$	SAFE	0.18	<b>382</b> $\times$	UNK	0.02	<b>3787</b> $\times$	UNK	0.06	<b>1144</b> $\times$
$N_{49}$	SAFE	51.00	SAFE	8.74	<b>6</b> $\times$	SAFE	0.17	<b>295</b> $\times$	UNK	0.02	<b>3190</b> $\times$	UNK	0.05	<b>1038</b> $\times$
$N_{51}$	SAFE	174.00	SAFE	19.88	<b>9</b> $\times$	SAFE	0.22	<b>782</b> $\times$	UNK	0.02	<b>9626</b> $\times$	UNK	0.05	<b>3733</b> $\times$
$N_{52}$	SAFE	48.00	SAFE	13.37	<b>4</b> $\times$	SAFE	0.15	<b>318</b> $\times$	UNK	0.02	<b>3066</b> $\times$	UNK	0.04	<b>1306</b> $\times$
$N_{53}$	SAFE	39.00	SAFE	5.16	<b>8</b> $\times$	SAFE	0.27	<b>144</b> $\times$	UNK	0.02	<b>2301</b> $\times$	UNK	0.05	<b>823</b> $\times$
$N_{54}$	SAFE	43.00	SAFE	3.61	<b>12</b> $\times$	SAFE	0.19	<b>225</b> $\times$	UNK	0.02	<b>2497</b> $\times$	UNK	0.04	<b>1043</b> $\times$
$N_{55}$	SAFE	61.00	SAFE	5.29	<b>12</b> $\times$	SAFE	0.16	<b>371</b> $\times$	UNK	0.02	<b>3646</b> $\times$	UNK	0.06	<b>1102</b> $\times$
$N_{56}$	SAFE	55.00	SAFE	3.01	<b>18</b> $\times$	SAFE	0.19	<b>294</b> $\times$	UNK	0.02	<b>3300</b> $\times$	UNK	0.04	<b>1359</b> $\times$
$N_{57}$	SAFE	8.00	SAFE	1.13	<b>7</b> $\times$	SAFE	0.07	<b>114</b> $\times$	UNK	0.02	<b>478</b> $\times$	UNK	0.05	<b>163</b> $\times$
$N_{58}$	SAFE	48.00	SAFE	3.09	<b>16</b> $\times$	SAFE	0.28	<b>172</b> $\times$	UNK	0.02	<b>2642</b> $\times$	UNK	0.06	<b>742</b> $\times$
$N_{59}$	SAFE	57.00	SAFE	3.77	<b>15</b> $\times$	SAFE	0.14	<b>404</b> $\times$	UNK	0.02	<b>3266</b> $\times$	UNK	0.05	<b>1084</b> $\times$
<b>VT</b>		<b>8470</b>		<b>477.01</b>			<b>12.20</b>			<b>0.78</b>			<b>2.19</b>	
<b>SAFE</b>	<b>42/45</b>		<b>42/45</b>			<b>32/45</b>			<b>1/45</b>			<b>0/45</b>		
<b>Imp.</b>					<b>12</b> $\times$			<b>544</b> $\times$			<b>10785</b> $\times$			<b>3969</b> $\times$

Bold values indicate runtime improvement and summary statistics of the methods

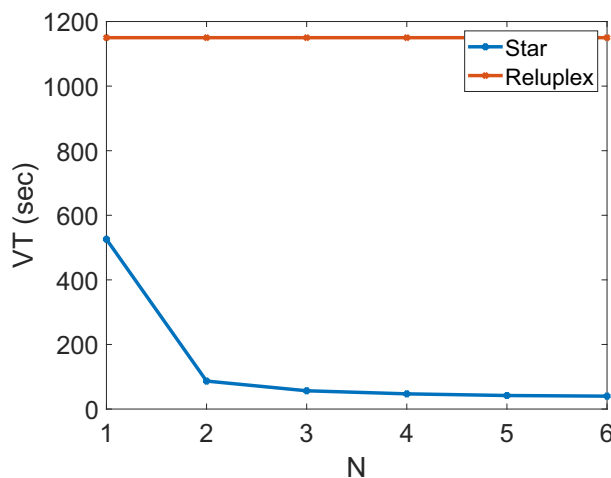


Fig. 4. Verification times for property  $\phi_4$  on  $N_{13}$  network with different number of cores

**Zonotope-based method** The experimental results show that the over-approximate, zonotope-based method [SGM<sup>+</sup>18] is faster than exact methods. In some cases, it can verify the safety of the networks with a small verification time, for example, the zonotope-based method successfully verifies property  $\phi_4$  on  $N_{41}$  network in 0.02 seconds. Although the zonotope-based method is very time-efficient, it is unable to verify the safety of most of networks due to its large over-approximation error. The zonotope approach can verify only 2/45 ( $\approx 4.44\%$ ) networks for property  $\phi_3$  and 1/45 ( $\approx 2.22\%$ ) networks for property  $\phi_4$ . In comparison with our approximate star method, we can verify 29/45 ( $\approx 64.44\%$ ) networks for property  $\phi_3$  and 32/45 ( $\approx 71.11\%$ ) networks for property  $\phi_4$ . This shows that the star set method is significantly less conservative than the zonotope approach.

**Abstract-domain based method** Similar to zonotope method, the abstract-domain is very time-efficient [SGPV19]. However, it is the most conservative approach since it cannot verify any networks for both properties. We note that in [TML<sup>+</sup>19b], we implemented an improved version of the new abstract domain method in which we still solve LP optimization problems to find the lower and upper bounds of an input to a specific neuron and use these bounds to construct the reachable set. In this paper, to have a fair comparison, we re-implement the original abstract domain method in which the lower and upper bounds are found using “back-tracking” method. We have experienced that the lower and upper bounds obtained by the new abstract-domain method may be very conservative.

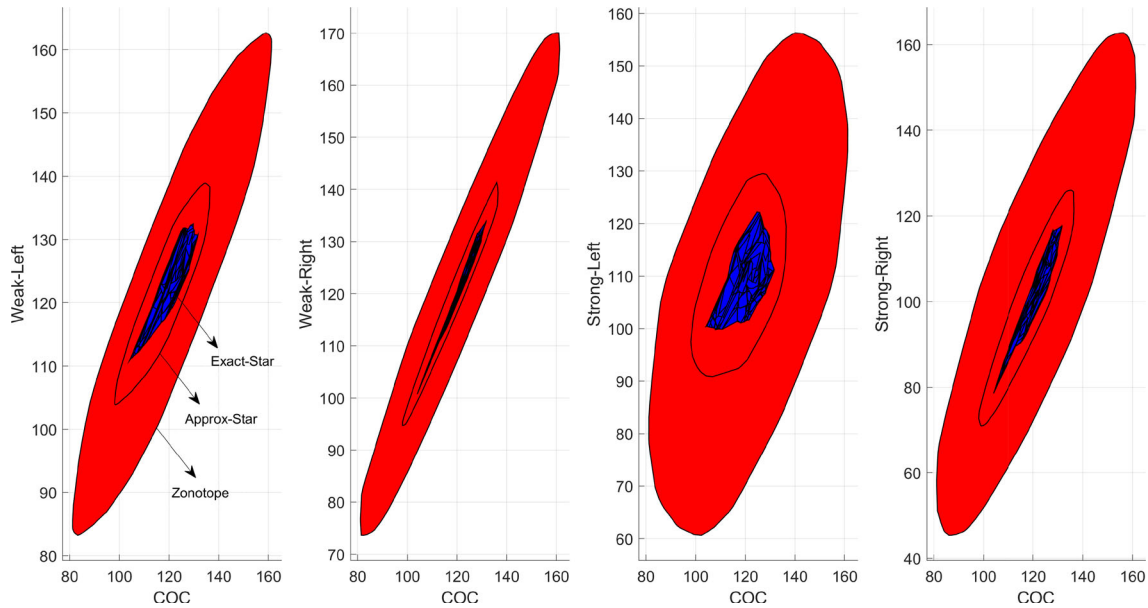


Fig. 5. Reachable sets of  $N_{41}$  network w.r.t property  $\phi_4$  with different methods

**Benefits of computing the reachable set** The computed reachable sets are useful for intuitively observing the complex behavior of the network. For example, Fig. 5 describes the behaviors of  $N_{51}$  network corresponding to property  $\phi_4$  requiring that the output  $COC$  is not the minimal score. From the figure, one can see that the  $COC > StrongRight$  and thus, property  $\phi_4$  holds on  $N_{41}$  network. Importantly, as shown in the figure, via visualization, one can intuitively observe the conservativeness of different over-approximation approaches in comparison to the exact ones which is impossible if we use *ERAN* [SGPV19], a *C-Python* implementation of the zonotope and new abstract domain methods. We note that the reachable set obtained by the new abstract domain method is neglected for visualization because it is too large. Last but not least, the reachable set is useful in the case that we need to verify a set of safety properties corresponding to the same input set. In this case, once the reachable set is obtained, it can be re-used to check different safety properties without rerunning the whole verification procedure as Reluplex does, and thus helps saving a significant amount of time.

**Complete counter example input set construction** Another strong advantage of our approach in comparison with other existing approaches is, in the case that a neural network violates its safety specification, our exact star method can construct a *complete counter input set* that leads the neural network to the unsafe region. The complete counter input set can be used as an adversarial input generator [GSS14, BIL<sup>+</sup>16] for robust training of the network. We note that finding a single counter input falsifying a safety property of a neural network can be done efficiently using only random simulations. However, constructing a complete counter input set that contains all counter inputs is very challenging because of the non-linearity of a neural network. To the best of our knowledge, our exact star-based approach is the only approach that can solve this problem. For example, assume that we want to check the following property  $\phi'_4 \triangleq \neg(COC \geq 15.8 \wedge StrongRight \leq 15.09)$  on  $N_{2.8}$  network with the same input constraints as in property  $\phi_4$ . Using the available reachable set of  $N_{2.8}$  network, we can verify that the above property  $\phi'_4$  is violated in which 60 stars in 421 stars of the reachable set reach the unsafe region. Using Theorem 3.2, we can construct a complete counter input set which is a union of 60 stars in 0.9893 seconds. This counter input set depicted in Fig. 6 is a part of the input set that contains all counter inputs that make the neural network unsafe.



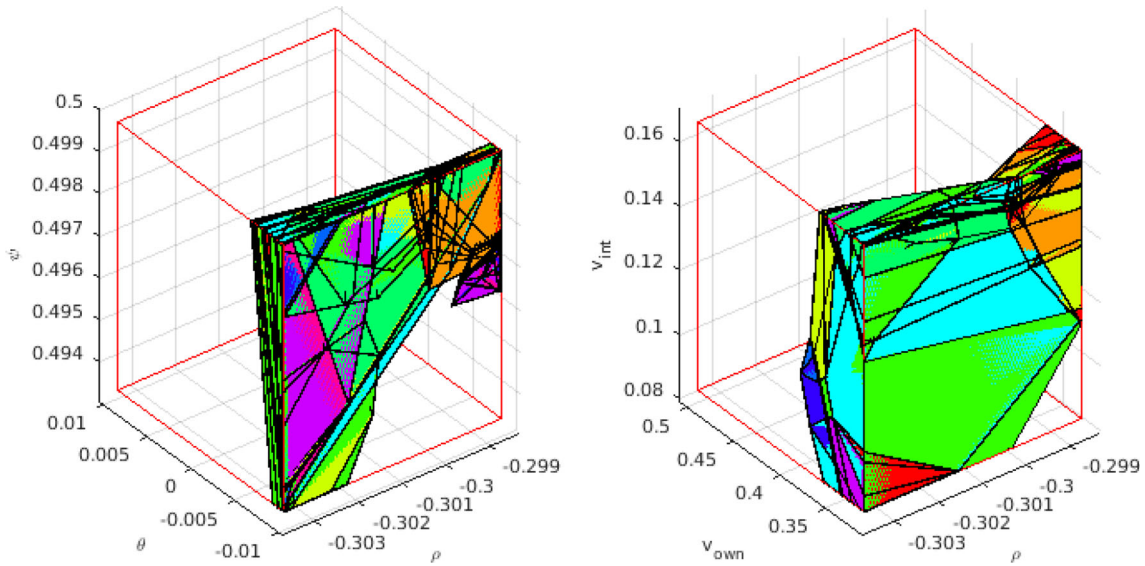


Fig. 6. The (normalized) complete counter input set for property  $\phi_4$  on  $N_{2\_8}$  network is a part of the normalized input set (red boxes)

## 5.2. Robustness certification of image classification ReLU DNNs under adversarial attacks

Robustness certification of DNNs becomes more and more important as many safety-critical applications using image classification DNNs can be fooled easily by slightly perturbing a correctly classified input. A network is said to be  $\delta$ -locally-robust at input point  $x$  if for every  $x'$  such that  $\|x - x'\|_\infty \leq \delta$ , the network assigns the same label to  $x$  and  $x'$ . In this case study, instead of proving the robustness of a network corresponding to a given robustness certification  $\delta$ , we focus on finding the maximum robustness certification value  $\delta_{max}$  that a verification method can provide a robustness guarantee for the network. We investigate this interesting problem on a set of image classification DNN with different architectures trained (with an accuracy of 98%) using the well-known MNIST data set consisting of 60,000 images of handwritten digits with a resolution of  $28 \times 28$  pixels [LeC98] (50,000 are used for training, while 10,000 are used for testing). The trained networks have 784 inputs and a single output with expected value from 0 to 9. We trained a variety of feed-forward neural networks with varying sizes. The number of neurons that we considered ranged from 140 Neurons to 4000 neurons. Throughout our experiments, we used a fixed number of neurons for each layer. We utilized mean squared error as the loss function and utilized the Adam optimization algorithm for training.<sup>3</sup> On average, it took 12 epochs for the training to converge to an accuracy greater than 95%. We find the maximum robustness verification value  $\delta_{max}$  for the networks on an image of digit one with the assumption that there is a  $\delta_{max}$ -bounded disturbance modifying the (normalized) values of the input vector  $x$  at all pixels of the image, i.e.,  $|x[i] - x'[i]| \leq \delta_{max}$ . The result are presented in Table 3. We note that the polyhedron and Reluplex approaches are not applicable for these networks because they cannot deal with a high-dimensional input space. The table shows that our approximate star approach produces larger upper bounds of the robustness values of the networks with many layers. For single layer networks, our approach gives the same results as the zonotope [SGM<sup>+</sup>18] and the abstract domain [SGPV19] methods. The exact-star method can prove that the network  $N_1$  is robust with the bounded disturbance  $\delta = 0.0058$ . When  $\delta > 0.0058$ , we ran into the “out of memory” issue in parallel computation since the number of the reachable sets becomes too large. The exact star method reaches timeout (set as 1 h) when finding the maximum robustness value for the other networks.

<sup>3</sup>The hyperparameters we utilized were  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ ,  $lr = 0.001$ .

**Table 3.** Maximum robustness values ( $\delta_{max}$ ) of image classification networks with different methods in which  $k$  is the number of hidden layers of the network,  $N$  is the total number of neurons,  $Tol$  is the tolerance error in searching.

Net	Parameters	Tol	$\delta_{max}$			
			Zonotope	Approximate-Star	Abstract-Domain	Exact-Star
$N_1$	$k = 5, N = 140$	0.0001	0.0046	0.0048	0.0025	$\geq 0.0058$
$N_2$	$k = 5, N = 250$	0.0001	0.0087	0.0101	0.0042	TimeOut
$N_3$	$k = 2, N = 1000$	0.0001	0.0072	0.0089	0.0052	TimeOut
$N_4$	$k = 1, N = 2000$	0.0001	0.0027	0.0027	0.0027	TimeOut
$N_5$	$k = 1, N = 4000$	0.0001	0.0035	0.0035	0.0035	TimeOut

### 5.3. Robustness certification of image classification DNNs with LeakyReLU, SatLin, and SatLins activation functions

We further evaluate our method on DNNs with LeakyReLU, Satlin, and Satlins activation function via verifying the robustness of image classification DNNs under brightening attack [TBXJ20]. We train three image classification networks on the MNIST data set using Matlab Deep Learning Toolbox. Each network has 784 (flattening each image with 28X28 resolution) input neurons and gives ten output neurons (for ten different classes; each neuron outputs corresponding to each of the classes; the max valued class being the prediction). The input dataset is normalized with  $mean = 33.3184$  and  $std = 1$  for networks with LReLU activation, otherwise no normalization is considered. The accuracy of each of the networks is between 96 and 98 percent.

To certify the robustness of the networks, we selected 100 images that are correctly classified by the respective networks and performed the brightening attack on these images. Then we calculated the reachable output sets for each of the perturbed images to determine if the correctly classified output always has the maximum value in the set compared to the others. The output layer activation '*softmax*' is ignored for the approximate reachable set calculations.

In a brightening attack, some pixels are changed independently in the image to make it brighter or darker to fool the network, i.e., to misclassify the image. In this case study, to implement this attack, the value of that pixel  $x_i$  is reduced to the new value  $x'_i$  if it is greater than a threshold value  $d$  such that  $0 \leq x'_i \leq \delta \cdot x_i$ . Here  $\delta$  is a very small number and indicates the size of the input set that can be created by a specific attack. When  $d$  is small, the number of pixels in the image that are attacked is large and vice versa. The verification results of the networks corresponding to different values of  $d$  and  $\delta$  are presented in Tables 4, 5 and 6.

**Table 4.** Verification results of 5-layered LeakyReLU network.

	Robustness Results (in Percent)								
	$\delta = 0.01$			$\delta = 0.05$			$\delta = 0.1$		
	Approx-Star	Abs-Dom	Approx-Zono	Approx-Star	Abs-Dom	Approx-Zono	Approx-Star	Abs-Dom	Approx-Zono
$d = 254$	91.00	91.00	91.00	90.00	88.00	87.00	86.00	82.00	82.00
$d = 250$	70.00	70.00	70.00	64.00	58.00	53.00	48.00	31.00	26.00
$d = 245$	66.00	66.00	66.00	57.00	49.00	47.00	45.00	27.00	22.00
	Verification Times (in Seconds)								
$d = 254$	1.17	49.38	1.34	3.98	55.07	1.35	18.07	57.89	0.86
$d = 250$	1.11	44.52	0.50	7.95	49.87	1.49	34.20	65.73	2.76
$d = 245$	1.13	43.12	0.65	9.25	50.15	1.80	44.23	74.60	3.17

**Table 5.** Verification results of 4-layered Satlin network.

	<b>Robustness Results (in Percent)</b>								
	$\delta = 0.01$			$\delta = 0.05$			$\delta = 0.1$		
	<i>Approx-Star</i>	<i>Abs-Dom</i>	<i>Approx-Zono</i>	<i>Approx-Star</i>	<i>Abs-Dom</i>	<i>Approx-Zono</i>	<i>Approx-Star</i>	<i>Abs-Dom</i>	<i>Approx-Zono</i>
$d = 254$	89.00	89.00	89.00	86.00	85.00	81.00	83.00	82.00	74.00
$d = 245$	60.00	60.00	59.00	48.00	47.00	26.00	37.00	33.00	7.00
$d = 234$	58.00	58.00	57.00	44.00	42.00	19.00	26.00	26.00	2.00
	<b>Verification Times (in Seconds)</b>								
$d = 254$	32.86	27.56	0.95	29.16	29.32	1.13	33.09	30.48	1.33
$d = 245$	22.05	21.67	0.83	32.16	32.17	2.56	37.68	28.54	5.50
$d = 234$	24.94	22.96	1.04	40.16	27.05	3.46	35.43	29.34	3.80

**Table 6.** Verification results of 5-layered Satlins network.

	<b>Robustness Results (in Percent)</b>								
	$\delta = 0.01$			$\delta = 0.05$			$\delta = 0.1$		
	<i>Approx-Star</i>	<i>Abs-Dom</i>	<i>Approx-Zono</i>	<i>Approx-Star</i>	<i>Abs-Dom</i>	<i>Approx-Zono</i>	<i>Approx-Star</i>	<i>Abs-Dom</i>	<i>Approx-Zono</i>
$d = 254$	86.00	86.00	82.00	83.00	83.00	66.00	80.00	80.00	51.00
$d = 245$	57.00	57.00	40.00	42.00	41.00	4.00	27.00	23.00	2.00
$d = 234$	51.00	51.00	34.00	40.00	37.00	1.00	19.00	14.00	0.00
	<b>Verification Times (in Seconds)</b>								
$d = 254$	50.14	49.36	1.55	59.87	52.35	1.69	56.43	50.81	2.26
$d = 245$	31.95	32.07	1.83	50.79	46.09	3.54	92.11	61.21	4.08
$d = 234$	31.21	30.36	1.98	56.34	49.65	3.58	80.83	65.43	3.79

From the comparative tables, we can infer that all the reachability approaches that were considered exhibit almost similar robustness performance (in terms of the correct number of class predictions) when  $\delta$  is very small, i.e., when the attack on the pixel is minimal, when the threshold value  $d$  decreases and  $\delta$  value increases the robustness performance changes. It is important to emphasize that in this case study, we solve LP to find the ranges of neurons of the networks for the new abstract domain approach (in the ACAS Xu case study, we use only the estimated range for this method). Because the ranges of neurons are optimized, the abstract domain method is more precise. However, at the same time, its verification time grows.

From the tables, we ascertain that the approximate star method is more robust compared to the abstract domain and Zonotope in terms of prediction. Interestingly, for the LeakyReLU network, the approximate star method is much faster than the abstract domain method. This impressive improvement is due to the optimization in the implementation of the approximate reachability algorithm for the LeakyReLU network in which instead of performing a sequence of approximate stepReLU operations, we perform these operations in parallel. Though the zonotope method takes less time to analyze robustness, it is very clear that it can only predict 2% of the images correctly when the brightening attack on pixels is higher (i.e., lesser  $d$  and higher  $\delta$ ). The prediction accuracy of the abstract domain approach with optimized ranges of neurons is better than the Zonotope but lesser than approx-star. Overall, the approx-star method achieves a more consistent level of robustness performance for varying disturbance levels than others. Importantly, by parallelizing the approximate stepReLU operations, we can speed up the approximate star method by several orders of magnitudes.

## 6. Conclusion and future work

We have proposed two reachability analysis algorithms for DNNs with piecewise linear activations using star sets, one that is exact (sound and complete), but has scalability limitations, and one that over-approximates (sound) with better scalability. The exact algorithm can compute and visualize the exact behaviors of DNNs. The exact method is more efficient than standard polyhedra approaches, and faster than SMT-based approaches when running on multi-core platforms. The over-approximate algorithm is much faster than the exact one, and notably, it is less conservative than recent zonotope and abstract-domain based approaches. Our algorithms are applicable for real world applications as shown in the safety verification of ACAS Xu DNNs and robustness certification of image classification DNNs. In ongoing and future work, we are extending the proposed methods for convolutional neural networks (CNNs) [TBXJ20, TPM<sup>+</sup>21] and recurrent neural networks (RNNs), improving scalability for other types of activation functions, such as tanh and sigmoid, and continuing to improve the NNV software framework [TYL<sup>+</sup>20].

## Acknowledgements

The material presented in this paper is based upon work supported by the Air Force Office of Scientific Research (AFOSR) through contract numbers FA9550-18-1-0122 and FA9550-19-1-0288, the Defense Advanced Research Projects Agency (DARPA) through contract number FA8750-18-C-0089, and the National Science Foundation (NSF) through Grant Number 1910017. The U.S. government is authorized to reproduce and distribute reprints for Governmental purposes notwithstanding any copyright notation thereon. Any opinions, findings, and conclusions or recommendations expressed in this publication are those of the authors and do not necessarily reflect the views of AFOSR, DARPA, or NSF.

**Publisher's Note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

## References

- [AKLP19] Akintunde ME, Kevorchian A, Lomuscio A, Pirovano E (2019) Verification of RNN-based neural agent-environment systems. In: Proceedings of the 33th AAAI conference on artificial intelligence (AAAI19). Honolulu, HI, USA. AAAI Press (**to appear**)
- [ALMP18] Akintunde M, Lomuscio A, Maganti L, Pirovano E (2018) Reachability analysis for neural agent-environment systems. In: Sixteenth international conference on principles of knowledge representation and reasoning
- [BD17] Bak S, Duggirala PS (2017) Simulation-equivalent reachability of large linear systems with inputs. In: International conference on computer aided verification. Springer, pp 401–420
- [BDDT<sup>+</sup>16] Bojarski M, Del Testa D, Dworakowski D, Firner B, Flepp B, Goyal P, Jackel LD, Monfort M, Muller U, Zhang J et al (2016) End to end learning for self-driving cars. arXiv preprint [arXiv:1604.07316](https://arxiv.org/abs/1604.07316)
- [BIL<sup>+</sup>16] Bastani O, Ioannou Y, Lampropoulos L, Vytiniotis D, Nori A, Criminisi A (2016) Measuring neural net robustness with constraints. In: Advances in neural information processing systems, pp 2613–2621
- [BTHJ20] Bak S, Tran H-D, Hobbs K, Johnson T (2020) Improved geometric path enumeration for verifying ReLU neural networks. In: Proceedings of the 32nd international conference on computer aided verification. Springer
- [DJST17] Dutta S, Jha S, Sanakranarayanan S, Tiwari A (2017) Output range analysis for deep neural networks. arXiv preprint [arXiv:1709.09130](https://arxiv.org/abs/1709.09130)
- [Ehl17] Ehlers R (2017) Formal verification of piece-wise linear feed-forward neural networks. In: International symposium on automated technology for verification and analysis. Springer, pp 269–286
- [GMDC<sup>+</sup>18] Gehr T, Mirman M, Drachler-Cohen D, Tsankov P, Chaudhuri S, Vechev M (2018) Ai 2: safety and robustness certification of neural networks with abstract interpretation. In: 2018 IEEE symposium on security and privacy (SP)
- [GSS14] Goodfellow IJ, Shlens J, Szegedy C (2014) Explaining and harnessing adversarial examples. arXiv preprint [arXiv:1412.6572](https://arxiv.org/abs/1412.6572)
- [HDY<sup>+</sup>12] Hinton G, Deng L, Yu D, Dahl GE, Mohamed A, Jaitly N, Senior A, Vanhoucke V, Nguyen P, Sainath TN et al (2012) Deep neural networks for acoustic modeling in speech recognition: the shared views of four research groups. *IEEE Signal Process Mag* 29(6):82–97
- [Hei20] Heilweil R (2020) *Tesla needs to fix its deadly Autopilot problem*
- [HKWW17] Huang X, Kwiatkowska M, Wang S, Wu M (2017) Safety verification of deep neural networks. In: International conference on computer aided verification. Springer, pp 3–29
- [JKO18] Julian KD, Kochenderfer MJ, Owen MP (2018) Deep neural network compression for aircraft collision avoidance systems. arXiv preprint [arXiv:1810.04240](https://arxiv.org/abs/1810.04240)
- [KBD<sup>+</sup>17] Katz G, Barrett C, Dill DL, Julian K, Kochenderfer MJ (2017) Reluplex: an efficient smt solver for verifying deep neural networks. In: International conference on computer aided verification. Springer, pp 97–117
- [KGBM04] Kvasnica M, Grieder P, Baotić M, Morari M (2004) Multi-parametric toolbox (MPT). In: International workshop on hybrid systems: computation and control. Springer, pp 448–462

- [KHI<sup>+</sup>19] Guy K, Derek AH, Duligur I, Kyle J, Christopher L, Rachel L, Parth S, Shantanu T, Haoze W, Aleksandar Z et al (2019) The marabou framework for verification and analysis of deep neural networks. In: International conference on computer aided verification. Springer, pp 443–452
- [KL18] Kouvaros P, Lomuscio A (2018) Formal verification of cnn-based perception systems. arXiv preprint [arXiv:1811.11373](https://arxiv.org/abs/1811.11373)
- [LeC98] LeCun Y (1998) The mnist database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>
- [LKB<sup>+</sup>17] Litjens G, Kooi T, Bejnordi BE, Setio AAA, Ciampi F, Ghahfoorian M, van der Laak Jeroen AWM, Van Ginneken B, Sánchez CI (2017) A survey on deep learning in medical image analysis. *Med Image Anal* 42:60–88
- [LM17] Lomuscio A, Maganti L (2017) An approach to reachability analysis for feed-forward relu neural networks. arXiv preprint [arXiv:1706.07351](https://arxiv.org/abs/1706.07351)
- [LWL<sup>+</sup>17] Liu W, Wang Z, Liu X, Zeng N, Liu Y, Alsaadi FE (2017) A survey of deep neural network architectures and their applications. *Neurocomputing* 234:11–26
- [MDF16] Moosavi-Dezfooli S-M, Fawzi A, Frossard P (2016) Deepfool: a simple and accurate method to fool deep neural networks. In: Proceedings of the IEEE conference on computer vision and pattern recognition, pp 2574–2582
- [Muo17] Muoio D (2017) The self-driving Uber in the Arizona crash was hit crossing an intersection on yellowUber crashes
- [PT10] Pulina L, Tacchella A (2010) An abstraction-refinement approach to verification of artificial neural networks. In: International conference on computer aided verification. Springer, pp 243–257
- [SGM<sup>+</sup>18] Singh G, Gehr T, Mirman M, Püschel M, Vechev M (2018) Fast and effective robustness certification. In: Advances in neural information processing systems, pp 10825–10836
- [SGPV19] Singh G, Gehr T, Püschel M, Vechev M (2019) An abstract domain for certifying neural networks. *Proc ACM Programm Lang* 3(POPL):41
- [TBXJ20] Tran H-D, Bak S, Xiang W, Johnson TT (2020) Verification of deep convolutional neural networks using imagestars. In: 32nd international conference on computer-aided verification (CAV). Springer
- [TML<sup>+</sup>19a] Tran H-D, Musau P, Lopez DM, Yang X, Nguyen LV, Xiang W, Johnson TT (2019) Parallelizable reachability analysis algorithms for feed-forward neural networks. In: 7th international conference on formal methods in software engineering (FormalISE2019), Montreal, Canada
- [TML<sup>+</sup>19b] Tran H-D, Musau P, Lopez DM, Yang X, Nguyen LV, Xiang W, Johnson TT (2019) Star-based reachability analysis for deep neural networks. In: 23rd international symposium on formal methods (FM'19). Springer
- [TPM<sup>+</sup>21] Tran H-D, Pal N, Musau P, Yang X, Hamilton NP, Lopez DM, Bak S, Johnson TT (2021) Robustness verification of semantic segmentation neural networks using relaxed reachability. In: Proceedings of the 33rd international conference on computer-aided verification. Springer
- [TYL<sup>+</sup>20] Tran H-D, Yang X, Lopez DM, Musau P, Nguyen LV, Xiang W, Bak S, Johnson TT (2020) NNV: the neural network verification tool for deep neural networks and learning-enabled cyber-physical systems. In: 32nd international conference on computer-aided verification (CAV)
- [WPW<sup>+</sup>18a] Wang S, Pei K, Whitehouse J, Yang J, Jana S (2018) Efficient formal safety analysis of neural networks. In: Advances in neural information processing systems, pp 6369–6379
- [WPW<sup>+</sup>18b] Wang S, Pei K, Whitehouse J, Yang J, Jana S (2018) Formal security analysis of neural networks using symbolic intervals. arXiv preprint [arXiv:1804.10829](https://arxiv.org/abs/1804.10829)
- [WZC<sup>+</sup>18] Weng T-W, Zhang H, Chen H, Song Z, Hsieh C-J, Boning D, Dhillon IS, Daniel L (2018) Towards fast computation of certified robustness for relu networks. arXiv preprint [arXiv:1804.09699](https://arxiv.org/abs/1804.09699)
- [XMW<sup>+</sup>18] Xiang W, Musau P, Wild AA, Lopez DM, Hamilton N, Yang X, Rosenfeld JA, Johnson TT (2018) Verification for machine learning, autonomy, and neural networks survey. *CoRR* [arXiv:1810.01989](https://arxiv.org/abs/1810.01989)
- [XTJ17] Xiang W, Tran H-D, Johnson TT (2017) Reachable set computation and safety verification for neural networks with relu activations. arXiv preprint [arXiv:1712.08163](https://arxiv.org/abs/1712.08163)
- [XTJ18] Xiang W, Tran H-D, Johnson TT (2018) Output reachable set estimation and verification for multilayer neural networks. *IEEE Trans Neural Netw Learn Syst* (99):1–7
- [XTJ19] Xiang W, Tran H-D, Johnson TT (2019) Specification-guided safety verification for feedforward neural networks. AAAI Spring symposium on verification of neural networks
- [ZWC<sup>+</sup>18] Zhang H, Weng T-W, Chen P-Y, Hsieh C-J, Daniel L (2018) Efficient neural network robustness certification with general activation functions. In: Advances in neural information processing systems, pp 4944–4953

*Received 10 April 2020*

*Accepted in revised form 29 May 2021 by Annabelle McIver, Maurice ter Beek and Cliff Jones*

*Published online 28 August 2021*