

Input-Relational Verification of Deep Neural Networks

DEBANGSHU BANERJEE, University of Illinois Urbana-Champaign, USA

CHANGMING XU, University of Illinois Urbana-Champaign, USA

GAGANDEEP SINGH, University of Illinois Urbana-Champaign and VMware Research, USA

We consider the verification of input-relational properties defined over deep neural networks (DNNs) such as robustness against universal adversarial perturbations, monotonicity, etc. Precise verification of these properties requires reasoning about multiple executions of the same DNN. We introduce a novel concept of difference tracking to compute the difference between the outputs of two executions of the same DNN at all layers. We design a new abstract domain, DiffPoly for efficient difference tracking that can scale large DNNs. DiffPoly is equipped with custom abstract transformers for common activation functions (ReLU, Tanh, Sigmoid, etc.) and affine layers and can create precise linear cross-execution constraints. We implement an input-relational verifier for DNNs called RaVeN which uses DiffPoly and linear program formulations to handle a wide range of input-relational properties. Our experimental results on challenging benchmarks show that by leveraging precise linear constraints defined over multiple executions of the DNN, RaVeN gains substantial precision over baselines on a wide range of datasets, networks, and input-relational properties.

CCS Concepts: • **Theory of computation** → **Program verification**; **Abstraction**; • **Computing methodologies** → **Neural networks**.

Additional Key Words and Phrases: Abstract Interpretation, Deep Learning, Relational Verification

1 INTRODUCTION

Deep neural networks (DNNs) have become more powerful and widespread over the past few years and have now penetrated almost all fields and application areas including safety-critical domains such as autonomous driving [10] or medical diagnosis [2], etc. Especially in these domains, the decisions generated from these DNNs are important and mistakes can have grave consequences. However, it can be hard to reason about DNNs as they are constructed in a black-box manner and have highly nonlinear behavior. As such, although the machine learning community has made great strides towards discovering and defending against DNN vulnerabilities [33, 50, 54, 60, 72, 84], these methods cannot guarantee safety. As a result, there has been a lot of work on verifying the safety properties of DNNs [3, 4, 6, 13, 14, 22, 32, 38, 39, 43, 56–58, 67, 68, 70, 75, 76, 82, 83, 86, 87, 89]. Despite this progress, existing DNN verification techniques can be imprecise for input-relational properties that arise in many practical scenarios. For example, most existing works mentioned above focus on verifying the absence of an adversarial attack (imperceptible perturbations added to an input) around a local neighborhood of test inputs. Recent work [46] has shown that attacks against individual inputs can be unrealistic as they rely on the attacker having perfect knowledge of the inputs processed by the DNN and being able to create perturbations specialized for that input. Indeed, many practical attack scenarios [46, 47, 49] involve constructing universal adversarial perturbations (UAPs) [54] that can work against a set of inputs. Other interesting input-relational properties that have become popular in recent years include monotonicity [74], and fairness [40]. Efficient verification of input-relational properties requires reasoning about the relationship between multiple executions of the same DNN. Existing verifiers lack these capabilities and as a result, are not precise. For the remainder of this paper, relational will refer to input-relational.

This Work. In this work, we propose a framework for verifying the relational properties of DNNs - RaVeN (Relational Verifier of Neural Networks). To the best of our knowledge, RaVeN is the first framework to verify a broad range of relational properties defined over multiple executions of the

Authors' addresses: Debangshu Banerjee, University of Illinois Urbana-Champaign, USA; Changming Xu, University of Illinois Urbana-Champaign, USA; Gagandeep Singh, University of Illinois Urbana-Champaign and VMware Research, USA.

same DNN. Next, we detail the key technical contributions that allow RaVeN to verify relational properties that state-of-the-art verifiers [68, 69, 88] cannot.

Main Contributions. Our main contributions are:

- A new abstract domain, DiffPoly with custom abstract transforms for affine and activation (ReLU, Sigmoid, Tanh, etc.) layers allowing us to efficiently compute precise lower and upper bounds of the difference between the outputs of a pair of DNN executions at each layer.
- A verification framework, RaVeN, which leverages the DiffPoly analysis to compute precise layerwise linear constraints over outputs from different executions of the DNN. These cross-execution linear constraints allow us to capture linear dependencies between the outputs of different DNN executions at each layer, making RaVeN more precise than existing state-of-the-art verifiers [68, 69, 88] which do not track linear dependencies at all layers. We use the linear constraints from DiffPoly analysis to formulate a mixed-integer linear program (MILP) (Section 4). We formally prove the soundness of RaVeN in Section 4.7.
- A complete implementation of RaVeN, including DiffPoly and MILP formulations capable of handling diverse relational properties defined over the same DNNs with the popular feedforward architectures and common activation functions like ReLU, Sigmoid, Tanh, etc.
- An extensive evaluation of RaVeN on a range of popular datasets, challenging fully-connected and convolutional networks, and diverse relational properties (e.g., UAP verification, monotonicity). Our results demonstrate that RaVeN achieves notably higher precision compared to prior approaches and can verify relational properties that are beyond the capabilities of current state-of-the-art verifiers (Section 5).

Our research can serve as a foundation for advancing relational verification in DNNs. Notably, our results indicate that DNNs exhibit improved provable robustness against universal attacks (UAPs), which are more realistic, compared to individual attacks. Recent studies [49, 85] demonstrate that defending against UAPs enhances accuracy and empirical robustness more effectively than defending against individual attacks [50]. In the future, integrating RaVeN into the training loop [52, 55, 90] can lead to DNNs with superior accuracy and provable robustness against UAPs. The supplementary materials¹ and code² are publicly available.

2 BACKGROUND

In this section, we present the essential background and notation used in this paper. Throughout the subsequent sections, lowercase letters (a, b , etc.) denote scalars, while uppercase letters (A, B , etc.) and the over barred lowercase letters (\bar{a}, \bar{b} , etc.) represent vectors and matrices.

Neural Networks: We primarily focus on feed-forward neural networks. However, since we use linear bound propagation techniques, similar to [86], our method can be extended to other architectures that can be expressed as DAGs (directed acyclic graphs). We use "DNN" to refer specifically to feed-forward neural networks. These DNNs, denoted as $N : \mathbb{R}^{n_0} \rightarrow \mathbb{R}^{n_l}$, are composed of l sequential layers N_1, \dots, N_l , where each $N_i : \mathbb{R}^{n_{i-1}} \rightarrow \mathbb{R}^{n_i}$ is a function. Each layer N_i applies either an affine function (convolution or linear function) or a non-linear activation function, such as ReLU, Sigmoid, or Tanh. Affine layers, represented as $N_i : \mathbb{R}^{n_{i-1}} \rightarrow \mathbb{R}^{n_i}$, are defined by $N_i(x) = A_i \cdot X + B_i$, where A_i is the weight matrix, and B_i is the bias vector.

2.1 Relational Verification of DNN

For a network $N : \mathbb{R}^{n_0} \rightarrow \mathbb{R}^{n_l}$ and a relational property defined over DNN inferences on k inputs, the input specification $\Phi : \mathbb{R}^{n_0 \times k} \rightarrow \{\text{true}, \text{false}\}$ is a boolean predicate. It encodes the input region

¹The latest version of the paper with appendix can be found at <https://focallab.org/files/raven.pdf>

²The code for RaVeN can be found at <https://github.com/uiuc-focal-lab/RaVeN>.

$\Phi_t \subseteq \mathbb{R}^{n_0 \times k}$ encompassing all potential inputs corresponding to each of the k DNN inferences. For any $X \in \mathbb{R}^{n_0 \times k}$ satisfying Φ , $X = (X_1, \dots, X_k)$ is a tuple of k points where $\forall i \in [k]. X_i \in \mathbb{R}^{n_0}$ and X_i is the input of the i -th DNN inference. Common DNN relational properties e.g. UAP verification [88], monotonicity [74], etc. can be encoded as the conjunction of k individual input specifications $\phi_{in}^i : \mathbb{R}^{n_0} \rightarrow \{true, false\}$ and cross-execution input specification $\Phi^\delta : \mathbb{R}^{n_0 \times k} \rightarrow \{true, false\}$. Each $\phi_{in}^i : \mathbb{R}^{n_0} \rightarrow \{true, false\}$ defines the input region $\phi_t^i \subseteq \mathbb{R}^{n_0}$ for i -th execution. Meanwhile, Φ^δ captures relationships between inputs used in distinct executions. Commonly Φ^δ bounds the difference between any pair of inputs $X_i, X_j \in \mathbb{R}^{n_0}$ used in different executions such as $L_{i,j} \leq X_i - X_j \leq U_{i,j}$ where $L_{i,j}, U_{i,j} \in \mathbb{R}^{n_0}$ are constant real vectors. Individual input regions ϕ_t^i are in general L_∞ regions [16] i.e. all $X_i \in \mathbb{R}^{n_0}$ such that $\|X_i - X_i^*\|_\infty \leq \epsilon$ around a concrete point $X_i^* \in \mathbb{R}^{n_0}$ with $\epsilon \in \mathbb{R}^+$. For any pair of inputs $X_i, X_j \in \mathbb{R}^{n_0}$, the cross-execution input specification between them $\phi_{i,j}^\delta$ are given by - $\phi_{i,j}^\delta(X_i, X_j) = (L_{i,j} \leq X_i - X_j) \wedge (X_i - X_j \leq U_{i,j})$. The output specification for relational properties is a boolean predicate $\Psi : \mathbb{R}^{n_l \times k} \rightarrow \{true, false\}$ defined over the outputs of all k DNN inferences. In this work, we consider output specifications Ψ that can be expressed as a logical formula in CNF (conjunctive normal form) with m clauses where each clause ψ_i is of the form below $C_{i,j,i'} \in \mathbb{R}^{n_l}$:

$$\psi_i(Y_1, \dots, Y_k) = \bigvee_{j=1}^n \psi_{i,j}(Y_1, \dots, Y_k) \quad \text{where } \psi_{i,j}(Y_1, \dots, Y_k) = \left(\sum_{i'=1}^k C_{i,j,i'}^T Y_{i'} \geq 0 \right)$$

Definition 2.1 (DNN Relational Verification Problem). The **relational verification** problem for a DNN $N : \mathbb{R}^{n_0} \rightarrow \mathbb{R}^{n_l}$, an input specification $\Phi : \mathbb{R}^{n_0 \times k} \rightarrow \{true, false\}$ and an output specification $\Psi : \mathbb{R}^{n_l \times k} \rightarrow \{true, false\}$ is to prove whether $\forall X_1, \dots, X_k \in \mathbb{R}^{n_0}. \Phi(X_1, \dots, X_k) \implies \Psi(N(X_1), \dots, N(X_k))$ or provide a counterexample otherwise.

2.2 Interesting Relational Properties of DNNs

UAP Verification. UAP verification problem verifies whether there exists a single perturbation that can be added to k DNN inputs to make it misclassify all of them. The UAP verification problem is fundamentally different from the commonly considered local L_∞ robustness verification where the adversary can perturb each input independently. However, as shown in recent studies [46, 47, 49] generating input-specific adversarial perturbation is unrealistic, and practical attacks require finding adversarial perturbation that works for a set of inputs instead of a single input. These works suggest that considering robustness against input-specific adversarial attacks is too conservative and presents a pessimistic view of practical DNN robustness. Since the adversarial perturbation is common across a set of inputs, the UAP verification problem requires a relational verifier that can exploit the dependency between perturbed inputs. We provide the input specification Φ and the output specification Ψ of the UAP verification problem in Appendix A.3. We describe another variation of UAP: targeted UAP in Appendix A.4.

Worst-case UAP accuracy: In general, for a given N , finding an adversarial perturbation that works for all inputs in a set is hard. However, an adversarial perturbation affecting a significant proportion of inputs also poses a threat to the DNN. Hence, most of the existing works compute the worst-case accuracy [88] of the DNN on an input set in the presence of a UAP adversary. The formal definition of worst-case UAP accuracy is as follows.

Definition 2.2 (Worst-case UAP accuracy). Given a DNN N , a set of inputs $I = \{X_1, \dots, X_k\}$, target outputs $O = \{Y_1, \dots, Y_k\}$ and perturbation norm bound $\epsilon \in \mathbb{R}$ the worst case UAP of N is $a^* = 1/k \min_{\|V\|_\infty \leq \epsilon} \sum_{i=1}^k (N(X_i + V) = Y_i)$ where V is the added perturbation.

Monotonicity Verification. Recent works have shown that local monotonicity of DNNs is interesting and verification for monotonic properties is desirable [21, 61]. This property asserts a

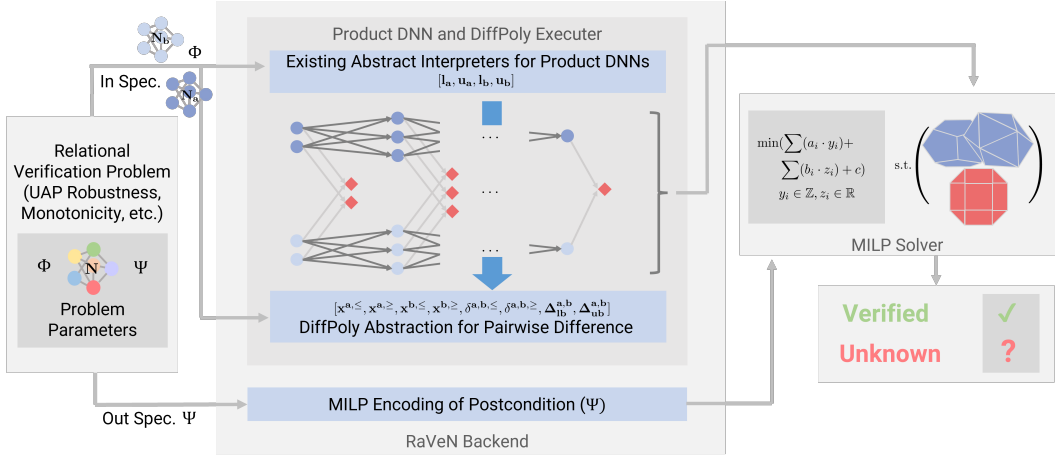


Fig. 1. The overview of the proposed sound and incomplete RaVeN verifier. Given a network N and a relational property (Φ, Ψ) relating k DNN inferences we show the flow of RaVeN along with the key steps - (i) constructing the product DNN by duplicating N k times and analyzing the product DNN with an existing DNN abstract interpreter, (ii) computing pairwise differences of outputs of all k inferences at each layer with DiffPoly analysis that uses concrete lower and upper bounds of each variable in the product DNN, (iii) combining DiffPoly analysis and product DNN analysis with an existing DNN abstract interpreter to infer layerwise linear constraints over outputs of all k DNN executions that preserves dependencies between different DNN executions, (iv) encoding the postcondition as a MILP objective and formulate MILP with layerwise linear constraints computed in step (iii). Finally, we use an off-the-shelf MILP solver [35] to verify the relational property by solving the corresponding MILP.

monotonic relationship between an input feature and the output. For instance, in predicting housing prices, a monotonic property could stipulate that a house with more rooms is consistently more expensive than a house with fewer rooms. We encode monotonicity as a relational property over a pair of DNN executions in Appendix A.6.

Hamming Distance. The Hamming distance between two strings is the number of substitutions needed to turn one string into the other [36]. Given a binary string (a list of images of binary digits), we want to formally verify the worst-case bounds on the hamming distance between the original binary string and classified binary string where each image of the binary digits can be perturbed by a common perturbation (formal definition in Appendix A.5). Hamming Distance serves as a valuable metric for tasks involving input string processing [62], like text comprehension or CAPTCHA solving.

Further Relational Verification Problems. Other than the properties described above, another interesting DNN property is fairness verification [40]. In fairness verification, we want to show a change in a sensitive feature does not change the output (i.e. the model is fair and unbiased towards that feature). We can encode the problem similarly to the monotonicity verification problem presented in the paper and verify it using RaVeN.

3 OVERVIEW

Fig. 1 illustrates the high-level idea behind the workings of RaVeN. It takes, as input, the DNN N and a relational property (Φ, Ψ) defined over k inferences of N . RaVeN computes a product DNN with k copies (one for each inference) of network N and runs existing DNN abstract interpreters [68, 69, 87] on each copy of N to obtain concrete lower and upper bounds of each variable in the

product DNN. However, the existing abstract interpreters analyze each DNN execution in isolation and as a result, fail to preserve the dependencies between outputs of different DNN executions. One of our key contributions is the design of a new abstract domain DiffPoly that can efficiently compute precise lower and upper bounds on differences between the outputs of a neuron corresponding to two DNN executions. While DiffPoly can be extended to track bounds on any linear combination of the layerwise outputs of any k DNN executions (Appendix G.5), we specifically focus on a pair of executions and track differences, not alternatives (e.g., sum), between them. This choice is motivated by the fact that for existing DNN relational properties (UAP verification, monotonicity, etc.), the difference between inputs used in multiple executions is bounded. Therefore, we naturally opt to track differences between the DNN’s outputs across multiple executions at subsequent hidden layers and the output layer. RaVeN combines the analysis of existing abstract interpreters on the product DNN and DiffPoly analysis on all $\binom{k}{2}$ pair of executions to infer linear constraints over the outputs of all k executions at each layer. The linear constraints computed by RaVeN capture the dependencies between different DNN executions at each layer making RaVeN more precise than the state-of-the-art relational verifier [88] that only tracks dependencies at the input layer but not at the hidden layers and loses precision as a result. At the final layer of N , we encode the output specification Ψ as a set of mixed-integer linear programming (MILP) constraints over the outputs of all k executions. Note that we use integer variables only to encode the output specification Ψ to limit the number of integer variables in the MILP formulation and subsequently avoid exponential blowup in MILP optimization time. Next, we elaborate on the workings of RaVeN with an illustrative example.

3.1 Illustrative Example

3.1.1 Network: For this example, we consider the network, N_{ex} , with three layers: two affine layers and one ReLU layer with two neurons each (Fig. 2). The weights on the edges represent the coefficients of the weight matrix used by the affine transformations applied at each layer and the learned bias for each neuron is shown above or below it. N_{ex} can be viewed as a loop-free straight-line program composed of a sequence of assignment statements - ReLU assignments $x_i \leftarrow \max(0, x_j)$ and affine assignments $x_i \leftarrow v + \sum_{j=1}^n w_j \cdot x_j$ where $v \in \mathbb{R}$ and $W = [w_1, \dots, w_n]^T \in \mathbb{R}^n$. In the example, N_{ex} is a program with 12 variables: 2 input variables - $\{i_1, i_2\}$, two output variables - $\{o_1, o_2\}$, 8 intermediate variables $\{x_1, \dots, x_8\}$ and a sequence assignment statements shown below:

$$\begin{array}{llllll} x_1 \leftarrow i_1 & x_3 \leftarrow x_1 - x_2 & x_5 \leftarrow \max(0, x_3) & x_7 \leftarrow x_5 - x_6 & o_1 \leftarrow x_7 \\ x_2 \leftarrow i_2 & x_4 \leftarrow -2 \cdot x_1 + x_2 & x_6 \leftarrow \max(0, x_4) & x_8 \leftarrow -x_5 + x_6 & o_2 \leftarrow x_8 \end{array} \quad (1)$$

3.1.2 Relational property: We verify the UAP verification problem described in Section 2.2 on N_{ex} where the relational property is defined over 2 separate executions of N_{ex} . Here the input specification $\forall X_1, X_2 \in \mathbb{R}^2. \Phi(X_1, X_2)$ is defined as follows where $X_1^* = [14, 11]^T$, $X_2^* = [11, 14]^T$, and $\epsilon = 6$.

$$\Phi(X_1, X_2) = (\|X_1 - X_1^*\|_\infty \leq \epsilon) \wedge (\|X_2 - X_2^*\|_\infty \leq \epsilon) \wedge (X_1 - X_2 = X_1^* - X_2^*) \quad (2)$$

In UAP verification, an adversary can select to attack the DNN with any perturbation δ such that $\|\delta\|_\infty \leq \epsilon$ but the same perturbation δ must be applied to both inputs - X_1^*, X_2^* . Therefore the two executions are related and tracking this relationship improves precision. In contrast, in the common local robustness problem, an adversary can choose different perturbations for the two inputs and therefore the two executions are unrelated and can be verified independently. Any input $X_1 \in \mathbb{R}^2$ inside the L_∞ ball defined by $\|X_1 - X_1^*\|_\infty \leq \epsilon$ is not misclassified if $(N_{ex}(X_1) = [o_1, o_2]^T) \wedge (o_1 - o_2 \geq 0)$ holds. Conversely, any input $X_2 \in \mathbb{R}^2$ lying inside the L_∞ ball - $\|X_2 - X_2^*\|_\infty \leq \epsilon$ is not misclassified if $(N_{ex}(X_2) = [o_1, o_2]^T) \wedge (o_2 - o_1 \geq 0)$ holds. We want to formally verify that there does not

exist an adversarial perturbation $\delta \in \mathbb{R}^2$ with $\|\delta\|_\infty \leq \epsilon$ such that both the inferences on inputs $X_1 = X_1^* + \delta$ and $X_2 = X_2^* + \delta$ produces incorrect classification results. In this case, the output specification Ψ can be encoded such that $\forall \delta \in \mathbb{R}^2$ and $\|\delta\|_\infty \leq \epsilon$ the network N_{ex} correctly classifies at least one of the two perturbed inputs $X_1 = X_1^* + \delta$ and $X_2 = X_2^* + \delta$.

$$\Psi(N_{ex}(X_1), N_{ex}(X_2)) = (C_1^T N_{ex}(X_1) \geq 0) \vee (C_2^T N_{ex}(X_2) \geq 0) \quad \text{where } C_1 = [1, -1]^T \wedge C_2 = [-1, 1]^T$$

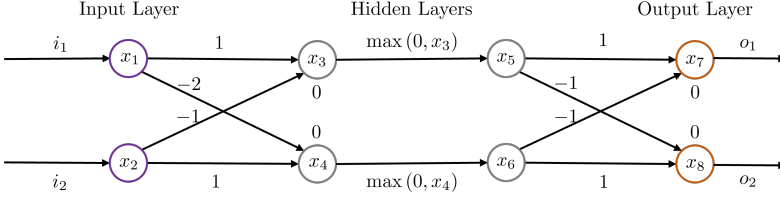


Fig. 2. Representation of N_{ex} used in the illustrative example

3.1.3 Product DNN construction & analysis: The input specification Φ (Eq. 2) relates two DNN executions on inputs from two input regions ϕ_1^1, ϕ_1^2 (not necessarily disjoint) defined by $\forall X_1 \in \mathbb{R}^2. \|X_1 - X_1^*\|_\infty \leq \epsilon$ and $\forall X_2 \in \mathbb{R}^2. \|X_2 - X_2^*\|_\infty \leq \epsilon$ respectively. So we construct the product DNN with two separate copies of the DNN - N_{ex}^1 and N_{ex}^2 where N_{ex}^1 and N_{ex}^2 track execution of N_{ex} on inputs from ϕ_1^1 and ϕ_1^2 respectively. The product DNN construction involves maintaining two separate copies of all 12 variables and all 10 assignment statements used in N_{ex} . In the product DNN, for each network N_{ex}^j where $j \in \{1, 2\}$, we rename input variables as $\{i_1^j, i_2^j\}$, output variables as $\{o_1^j, o_2^j\}$ and intermediate variables as $\{x_1^j, \dots, x_8^j\}$. N_{ex}^1 and N_{ex}^2 can be analyzed with any existing complete [24, 39] or incomplete DNN verifiers [69, 87]. However, for scalability, we use sound but incomplete abstract interpretation-based DNN verification techniques. We use the existing DeepZ [68] abstract interpreter to compute an overapproximated range of the possible values of each variable in N_{ex}^1 and N_{ex}^2 w.r.t. input regions ϕ_1^1 and ϕ_1^2 respectively. Fig. 12 in the appendix shows the range of values for each variable in the product DNN obtained by DeepZ analysis. The detailed execution of DeepZ for this example is in Appendix A.7.

3.1.4 Capturing dependencies between DNN executions: DeepZ (or, any other existing non-relational DNN verifier) analyze N_{ex}^1, N_{ex}^2 in isolation and do not track the relation captured in the cross-execution input constraint such as in Eq. 2 $\forall X_1, X_2. (X_1 - X_2 = X_1^* - X_2^*)$ that bounds the difference between the inputs used in different executions of the network. In contrast, the proposed DiffPoly can efficiently compute the bounds on the difference between two copies of the same variable corresponding to two different executions and as a result, can capture the dependencies between multiple executions. For example, given any variable x_i in N_{ex} DiffPoly computes lower and upper bound of $(x_i^1 - x_i^2)$ that holds for all possible inputs satisfying Φ . Overall, for any relational property defined over k DNN executions, we run $\binom{k}{2}$ DiffPoly for each pair of DNN executions. Note that since for any variable x_i , $(x_i^a - x_i^b) = -(x_i^b - x_i^a)$, for any pair of execution over inputs from ϕ_1^a , and ϕ_1^b , we only run DiffPoly analysis if $a < b$ to avoid redundant computations. For the rest of the paper, given a pair of variables $\langle x_i^a, x_i^b \rangle$ we use $\delta_{x_i}^{a,b}$ to denote their difference $(x_i^a - x_i^b)$.

3.1.5 DiffPoly domain: For two copies of the same variable from two separate executions e.g. x_i^a, x_i^b , the DiffPoly domain (formally described in Section 4.1), associates six linear constraints with $\langle x_i^a, x_i^b \rangle$: three upper linear constraints (symbolic upper bounds) $\delta_{x_i}^{a,b,\geq}, x_i^{a,\geq}, x_i^{b,\geq}$ and three lower linear constraints (symbolic lower bounds) $\delta_{x_i}^{a,b,\leq}, x_i^{a,\leq}, x_i^{b,\leq}$. The δ -constraints are the symbolic

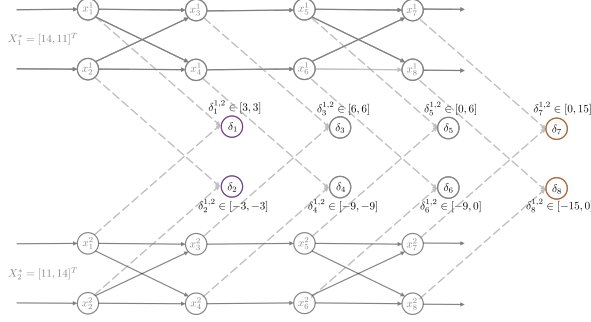


Fig. 3. Concrete bounds of difference as computed by DiffPoly analysis on the example network.

lower and upper bound on the difference $(x_i^a - x_i^b)$ satisfying $\delta_{x_i}^{a,b,\leq} \leq (x_i^a - x_i^b) \leq \delta_{x_i}^{a,b,\geq}$ while the other four constraints represent symbolic bounds on the variables x_i^a, x_i^b respectively. Additionally, the domain tracks concrete bounds - concrete lower bounds for each variable $(x_i^a - x_i^b)$, x_i^a , and x_i^b i.e. Δ_{lb}^{a,b,x_i} , l_{a,x_i} , and l_{b,x_i} and concrete upper bounds Δ_{ub}^{a,b,x_i} , u_{a,x_i} , and u_{b,x_i} . Note that as depicted in Fig 1, the concrete bounds - l_{a,x_i} , and l_{b,x_i} , u_{a,x_i} , and u_{b,x_i} are obtained from the analysis of the product DNN. At a high level, DiffPoly combines the ideas from the Zone domain [51], used for classical program analysis, that tracks concrete lower and upper bound on the difference of a pair of variables e.g. $l_{xy} \leq (x - y) \leq u_{xy}$ and the DeepPoly domain [69] that tracks symbolic lower and upper bound on variables of the DNN. However, DiffPoly is more precise than both the Zone domain which does not track symbolic bounds on the difference, and the DeepPoly domain which does not explicitly track any difference constraints making DiffPoly well suited for computing difference bounds across multiple DNN executions. Next, we show the format of symbolic bounds associated with DiffPoly below where $\delta_{x_j}^{a,b} = (x_j^a - x_j^b)$.

$$\delta_{x_i}^{a,b,\geq} = v + \sum_{j=1}^n \left(w_j^\delta \cdot \delta_{x_j}^{a,b} + w_j^a \cdot x_j^a + w_j^b \cdot x_j^b \right) \quad x_i^{a,\geq} = v_a^x + \sum_{j=1}^n w_j^{a,x} \cdot x_j^a \quad x_i^{b,\geq} = v_b^x + \sum_{j=1}^n w_j^{b,x} \cdot x_j^b \quad (3)$$

In Eq. 3, $v, v_a^x, v_b^x \in \mathbb{R}$, $W^\delta, W^a, W^b, W^{a,x}, W^{b,x} \in \mathbb{R}^n$ are the coefficients of the variables with w_i denoting the i -th coefficient for any vector $W \in \mathbb{R}^n$, n is the number of neurons in N_{ex} . We restrict the format of symbolic bounds and enforce $\forall j \geq i \quad w_j^\delta = w_j^a = w_j^b = w_j^{a,x} = w_j^{b,x} = 0$ so that symbolic bounds of any pair of variables $\langle x_i^a, x_i^b \rangle$ involve only variables that come before x_i^a, x_i^b (having smaller index) and their difference. These restrictions ensure that there are no cyclic dependencies between the symbolic bounds of the variables. Moreover, similar to the DeepPoly domain, we only allow a single symbolic lower, and upper bound to reduce the computation cost required to evaluate the concrete bounds for each variable. Otherwise, the unrestricted Polyhedra domain [20] though more precise, does not scale to the large DNNs considered in this work.

3.1.6 DiffPoly analysis: The analysis start with computing the symbolic and concrete bounds corresponding to $\langle x_1^1, x_1^2 \rangle$ and $\langle x_2^1, x_2^2 \rangle$. All pair of inputs X_1, X_2 satisfying input specification Φ satisfy $X_1 - X_2 = X_1^* - X_2^* = [3, -3]^T$. The linear constraints and concrete lower and upper bounds defining the range of the difference are as follows.

$$\delta_{x_1}^{1,2,\leq} = \delta_{x_1}^{1,2,\geq} = 3 \quad \delta_{x_2}^{1,2,\leq} = \delta_{x_2}^{1,2,\geq} = -3 \quad (x_1^1 - x_1^2) \in [3, 3] \quad (x_2^1 - x_2^2) \in [-3, -3]$$

At the input layer, the abstract elements also track linear constraints and concrete bounds for variables x_1^1, x_1^2, x_2^1 , and x_2^2 . However, for this example, we primarily focus on constraints $\delta_{x_i}^{1,2,\geq}$ and $\delta_{x_i}^{1,2,\leq}$ and show the rest of the constraints in the Appendix A.8. Next, we apply the affine

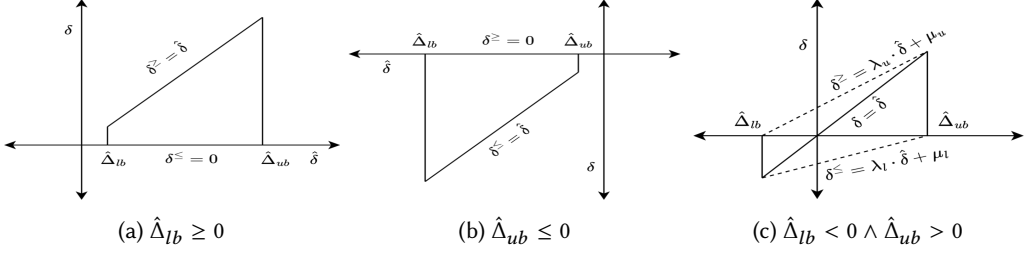


Fig. 4. The optimal (in terms of area) convex approximations for $\delta = \text{ReLU}(x) - \text{ReLU}(y)$ where $\hat{\delta} = (x - y)$, δ^{\geq} , and δ^{\leq} are symbolic upper bound and lower bound of δ respectively.

transformer (defined in Section 4.1) to calculate bounds corresponding to $\langle x_3^1, x_3^2 \rangle$ and $\langle x_4^1, x_4^2 \rangle$. We show the derivation of linear constraints $\delta_{x_3}^{1,2,\geq}$ and $\delta_{x_3}^{1,2,\leq}$ below where $\delta_{x_1}^{1,2} = (x_1^1 - x_1^2)$ and $\delta_{x_2}^{1,2} = (x_2^1 - x_2^2)$. The symbolic bounds $\delta_{x_4}^{1,2,\geq}$ and $\delta_{x_4}^{1,2,\leq}$ are obtained similarly.

$$\delta_{x_3}^{1,2} = (x_1^1 - x_2^1) - (x_1^2 - x_2^2) \implies \delta_{x_3}^{1,2,\geq} = \delta_{x_3}^{1,2,\leq} = (x_1^1 - x_1^2) - (x_2^1 - x_2^2) = \delta_{x_1}^{1,2} - \delta_{x_2}^{1,2} \quad (4)$$

To compute the concrete lower bound $\Delta_{lb}^{1,2,x_3}$ (or, upper bound) of $(x_3^1 - x_3^2)$ we substitute the concrete bounds of $\delta_{x_1}^{1,2}$ and $\delta_{x_2}^{1,2}$ in lower (upper) symbolic bounds of Eq. 4 for example:

$$\delta_{x_3}^{1,2,\leq} = \delta_{x_1}^{1,2} - \delta_{x_2}^{1,2} \implies \Delta_{lb}^{1,2,x_3} = \Delta_{lb}^{1,2,x_1} - \Delta_{ub}^{1,2,x_2} = 6$$

Next, we compute bounds corresponding to $\langle x_5^1, x_5^2 \rangle$ by using the ReLU abstract transformer (formally introduced in Section 4.2) for the assignments $x_5^1 \leftarrow \text{ReLU}(x_3^1)$ and $x_5^2 \leftarrow \text{ReLU}(x_3^2)$. In this case, choices for the symbolic bounds are non-unique. Fig. 4a shows one of two possible choices for linear constraints $\delta_{x_5}^{1,2,\geq} = \delta_{x_3}^{1,2}$ and $\delta_{x_5}^{1,2,\leq} = 0$. $\delta_{x_5}^{1,2,\geq} = x_5^{1,\geq} - x_5^{2,\leq}$ and $\delta_{x_5}^{1,2,\leq} = x_5^{1,\leq} - x_5^{2,\geq}$ are alternative candidates. However, in the abstract domain, we only allow only one choice for $\delta_{x_5}^{1,2,\geq}$ and one choice for $\delta_{x_5}^{1,2,\leq}$ so we greedily select one of two possible candidates for both $\delta_{x_5}^{1,2,\geq}$ and $\delta_{x_5}^{1,2,\leq}$. For both choices, we first evaluate the concrete bounds of $(x_5^1 - x_5^2)$ by substituting all variables in the symbolic lower (or upper) bound with their respective concrete bounds and then pick the candidate with the more precise concrete bound. For example, the choice $\delta_{x_5}^{1,2,\geq} = \delta_{x_3}^{1,2}$ yields concrete bound $\Delta_{ub}^{1,2,x_5} = 6.0$ which is more precise than $\Delta_{ub}^{1,2,x_5} = 20.625$ calculated from $\delta_{x_5}^{1,2,\geq} = x_5^{1,\geq} - x_5^{2,\leq}$. Thus, we select $\delta_{x_5}^{1,2,\geq} = \delta_{x_3}^{1,2}$. Finally, we obtain bounds corresponding to $\langle x_7^1, x_7^2 \rangle$ and $\langle x_8^1, x_8^2 \rangle$ by applying the affine abstract transformer. We show concrete bounds for the difference of each pair of variables $(x_i^1 - x_i^2)$ in Fig. 3 and detailed analysis in Appendix A.8.

3.1.7 Back-substitution for concrete bounds: We obtain the concrete bounds of each $(x_i^1 - x_i^2)$ by the back-substitution strategy used in most of the popular non-relational DNN verifiers e.g. CROWN [92], DeepPoly [69], α -CROWN [86], etc. In back-substitution, we start with the symbolic bounds $\delta_{x_i}^{a,b,\geq}$ (or, $\delta_{x_i}^{a,b,\leq}$) of $(x_i^1 - x_i^2)$ and then obtain concrete bounds Δ_{ub}^{a,b,x_i} (or, Δ_{lb}^{a,b,x_i}) of $(x_i^1 - x_i^2)$ by substituting concrete bounds of all the variables in $\delta_{x_i}^{a,b,\geq}$ (or, $\delta_{x_i}^{a,b,\leq}$). Commonly, back-substitution does not stop after a single concrete substitution step rather it refines Δ_{ub}^{a,b,x_i} (or, Δ_{lb}^{a,b,x_i}) by a sequence of steps with each step including a symbolic substitution, where all the variables in $\delta_{x_i}^{a,b,\geq}$ (or, $\delta_{x_i}^{a,b,\leq}$) are replaced by the corresponding symbolic bounds, followed by a concrete substitution. Although back-substitution is computationally more expensive than a single concrete substitution step, it obtains more precise concrete bounds Δ_{ub}^{a,b,x_i} (or, Δ_{lb}^{a,b,x_i}) which in turn improves the precision of RaVeN.

3.2 Using Analysis Bounds to Solve the UAP Verification Problem

We will now explain how RaVeN combines DiffPoly analysis with product DNN analysis to create the MILP formulation. Additionally, through our illustrative example, we will compare RaVeN’s approach to state-of-the-art baseline methods like [40] and [88]. This comparison will demonstrate that while the baseline methods fall short in confirming the absence of a UAP in our example, our approach successfully verifies the non-existence of a UAP.

3.2.1 State-of-the-art DNN relational verifiers. [40] only analyzes the product DNN and uses the concrete bounds obtained independently for each execution to verify UAP robustness. This approach does not track any dependencies across executions and just leverages standard DNN local robustness verification of individual inferences. However, DeepZ analysis on the product DNN computes for input region ϕ_t^1 the lower bound of $C_1^T N_{ex}(X_1)$ is -13.25 and for ϕ_t^2 the lower bound of $C_2^T N_{ex}(X_2)$ is -31.44 . Since the lower bounds of both $C_1^T N_{ex}(X_1)$ and $C_2^T N_{ex}(X_2)$ are less than 0 this method can not prove that UAP does not exist. Next, we focus on the state-of-the-art approach (referred to as I/O formulation in the rest of the paper) for UAP verification introduced by [88]. The I/O formulation initially applies non-relational DNN verifiers (e.g., DeepZ) to the product DNN. Based on DeepZ analysis, for each execution, it extracts linear constraints connecting output variables to input variables specific to that execution. Lastly, it translates the cross-execution input constraints into linear constraints, represents the output specification Ψ as a MILP objective, and employs standard MILP solvers to find the optimal solution (detailed formulation in Appendix B.1). For our illustrative example, the I/O formulation can only prove the absence of a UAP when the MILP solution is non-negative. However, the optimal MILP solution in this case is $-5.306 < 0$, highlighting that the I/O formulation lacks the precision to verify the relational property. This imprecision arises because the I/O formulation, while tracking dependencies at the input layers, neglects subsequent hidden layers, leading to a loss of precision.

3.2.2 RaVeN MILP formulation. We introduce a two-step enhancement to the MILP encoding in comparison to I/O formulation (same MILP objective) using our tool, RaVeN. To begin with, we relate the output of each layer to the output of the preceding layer by employing a set of linear constraints, commencing from the input layer. We replace non-linear activation layers (e.g., ReLU, Sigmoid, etc.) with convex overapproximations using concrete bounds obtained from DeepZ analysis, such as triangle relaxation [70] for ReLU. RaVeN’s layerwise approach effectively captures linear dependencies across executions at the hidden layers, yielding an improved optimal solution of -1.564 compared to the I/O formulation (details behind this improvement in Appendix B.2). Nonetheless, it remains insufficient for verifying the absence of UAP. In this case, the issue lies in the isolated computation of convex overapproximations for non-linear activation functions, which disregards the inter-dependencies between executions. To address this limitation, RaVeN utilizes the DiffPoly analysis and incorporates DiffPoly’s custom abstract transformers for non-linear activation functions defined over pairs of executions. This approach computes convex overapproximations that consider inter-dependencies between execution pairs. Figure 5 illustrates this enhancement, showing how constraints derived from the DiffPoly analysis enhance the precision of the convex region at the hidden layers. The addition of the difference constraints from the DiffPoly analysis to the layerwise formulation of RaVeN improves the optimal value to 0 thereby proving the absence of UAP in the illustrative example. It is important to note that RaVeN employs the same MILP encoding for Ψ

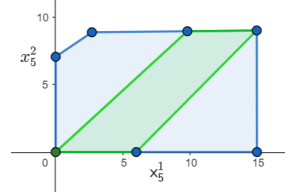


Fig. 5. For the variables x_5^1 and x_5^2 the convex region (green) obtained with constraints from DiffPoly analysis is more precise than the convex region (blue) formed without the difference constraints.

as utilized in the I/O formulation. The observed improvement is the result of RaVeN's enhanced capability in capturing the linear dependencies between outputs from multiple executions. The detailed MILP formulation for RaVeN is in Appendix B.3.

4 RAVEN ALGORITHM

In this section, we present RaVeN's pseudocode, discuss its key components, and assess its asymptotic runtime. We provide a sketch of the soundness proofs of RaVeN in Section 4.7 with detailed proofs in Appendix F. We first formally introduce the product DNN.

Definition 4.1 (Product DNN). Given any l layer DNN $N : \mathbb{R}^{n_0} \rightarrow \mathbb{R}^{n_l}$ and input specification Φ defined over k executions of N the product DNN $\mathcal{N}^k : \mathbb{R}^{n_0 \times k} \rightarrow \mathbb{R}^{n_l \times k}$ defined as sequential composition of l functions $\mathcal{N}_i^k : \mathbb{R}^{n_{i-1} \times k} \rightarrow \mathbb{R}^{n_i \times k}$ where $\mathcal{N}_i^k((X_1^i, \dots, X_k^i)) = [N_i(X_1^i), \dots, N_i(X_k^i)]^T$, for all $j \in [k]$. $X_j^i \in \mathbb{R}^{n_{i-1}}$ and $N_i : \mathbb{R}^{n_{i-1}} \rightarrow \mathbb{R}^{n_i}$ is the i -th layer of N .

Algorithm 1 shows the pseudocode for RaVeN. For the product DNN, an existing non-relational verifier (e.g. DeepZ) is used to obtain the concrete bounds for the outputs of all k executions at all layers, say \mathcal{A}^k (line 5). We use the concrete bounds from product DNN analysis (line 7) to initialize DiffPoly analysis for all $\kappa = \binom{k}{2}$ pair of executions (line 8). Next, DiffPoly computes the symbolic and concrete bounds (denoted as $\mathcal{A}_\delta^{a,b}$) of the outputs and their differences w.r.t each pair of executions (line 8). Note that aside from handling differences, DiffPoly also maintains symbolic bounds on the variables from the product DNN that are relevant to the pair of executions it is analyzing. This allows DiffPoly to calculate the concrete bounds of these product DNN variables using back-substitution although DiffPoly can also be run independently from product DNN analysis. However, we decide to utilize the concrete bounds from the product DNN analysis, as they can be more precise compared to the bounds obtained by DiffPoly. Furthermore, this approach enables DiffPoly to benefit from any improvements made in the product DNN analysis. We produce linear constraints for all layers by utilizing the symbolic and concrete bounds obtained from DiffPoly analysis on all κ pairs of executions. (line 10). After layerwise linear constraints computation, we encode Ψ , as a MILP objective (line 11). Finally, we invoke a MILP solver on the MILP formulated using the linear constraints and MILP objective function to verify the relational verification problem (line 12). Note, Algorithm 1 shows a sequential implementation of RaVeN. However, we can parallelly run existing DNN abstract interpreters on each of k copies of N and parallelly execute DiffPoly interpreter on all $\binom{k}{2}$ difference networks. Next, we formally define the building blocks of RaVeN algorithm: DiffPoly domain and layerwise MILP formulation.

4.1 DiffPoly Abstract Domain

Next, we formally introduce the DiffPoly domain and the corresponding abstract transformers for the affine and activation (ReLU, Sigmoid, Tanh, etc.) assignments. For a list of $2n$ variables $[x_1^a, \dots, x_n^a], [x_1^b, \dots, x_n^b]$ corresponding to a pair of execution of N the corresponding element in the DiffPoly domain \mathcal{A}_{2n} is defined as $\bar{a} = [a_1, \dots, a_n]$. Here each a_i is associated with a pair of variables $\langle x_i^a, x_i^b \rangle$. a_i associates (i) six symbolic bounds: symbolic lower and upper bounds for x_i^a, x_i^b and $(x_i^a - x_i^b)$ and (ii) six concrete bounds: concrete lower and upper bounds for x_i^a, x_i^b and $(x_i^a - x_i^b)$. We represent each a_i as a tuple $a_i = \langle C_{sym}^i, C_{con}^i \rangle$ with C_{sym}^i and C_{con}^i denoting the symbolic and concrete bounds respectively:

$$C_{sym}^i = \{x_i^{a,\leq}, x_i^{b,\leq}, \delta_{x_i}^{a,b,\leq}, x_i^{a,\geq}, x_i^{b,\geq}, \delta_{x_i}^{a,b,\geq}\} \quad C_{con}^i = \{l_{a,x_i}, l_{b,x_i}, \Delta_{lb}^{a,b,x_i}, u_{a,x_i}, u_{b,x_i}, \Delta_{ub}^{a,b,x_i}\}$$

The monotonic concretization function $\gamma_{2n} : \mathcal{A}_{2n} \rightarrow \wp(\mathbb{R}^{2n})$ mapping each abstract element \bar{a} to the corresponding element in the concrete domain $\wp(\mathbb{R}^{2n})$ (powerset of \mathbb{R}^{2n}), is shown in Eq. 5

Algorithm 1 RaVeN Algorithm

```

1: procedure RAVeN( $\Phi, \Psi, N$ )
2:   Input:  $\Phi : \mathbb{R}^{n_0 \times k} \rightarrow \{\text{true}, \text{false}\}, \Psi : \mathbb{R}^{n_1 \times k} \rightarrow \{\text{true}, \text{false}\}, N : \mathbb{R}^{n_0} \rightarrow \mathbb{R}^{n_1}$ .
3:   Verify:  $\forall X_1, \dots, X_k \in \mathbb{R}^{n_0}. \Phi(X_1, \dots, X_k) \implies \Psi(N(X_1), \dots, N(X_k))$ .
4:    $\mathcal{N}^k \leftarrow \text{ConstructProductDNN}(N, \Phi)$ 
5:    $\mathcal{A}^k \leftarrow \text{ProdDNNAnalyzer}(\mathcal{N}^k, \Phi, \mathcal{V})$  ▷  $\mathcal{V}$  is existing non-relational DNN verifier
6:   for  $a, b \in [k] \wedge a < b$  do
7:      $\mathcal{L}^a, \mathcal{U}^a, \mathcal{L}^b, \mathcal{U}^b \leftarrow \text{ExtractConcreteBounds}(\mathcal{A}_i^k, a, b)$ 
8:      $\mathcal{A}_{\delta}^{a,b} \leftarrow \text{DiffPolyExecutor}(N^a, N^b, \Phi, \mathcal{L}^a, \mathcal{U}^a, \mathcal{L}^b, \mathcal{U}^b)$ 
9:   end for
10:   $\mathcal{M} \leftarrow [\text{LayerwiseConstraints}(\mathcal{A}_{\delta}^{a,b}, N, \Phi) \mid a, b \in [k] \wedge b < a]$  ▷ Constraints
11:   $\mathcal{M}^{\Psi} \leftarrow \text{RaVeNObjectiveFunction}(\Psi)$  ▷ Objective Function Formulation
12:  return MILPSolver( $\mathcal{M}, \mathcal{M}^{\Psi}$ ) ▷ MILP Solver Invocation
13: end procedure

```

where for any $X \in \mathbb{R}^n$ we represent i -th coordinate of X as x_i .

$$\begin{aligned}
\varphi_{2n}^{\delta}(X^a, X^b) &= (X^a, X^b \in \mathbb{R}^n) \wedge (\forall i \in [n]. (\delta_{x_i}^{a,b,\leq} \leq (x_i^a - x_i^b) \leq \delta_{x_i}^{a,b,\geq} \wedge \Delta_{lb}^{a,b,x_i} \leq (x_i^a - x_i^b) \leq \Delta_{ub}^{a,b,x_i})) \\
\varphi_n(X^a) &= (X^a \in \mathbb{R}^n) \wedge (\forall i \in [n]. (x_i^{a,\leq} \leq x_i^a \leq x_i^{a,\geq} \wedge l_{a,x_i} \leq x_i^a \leq u_{a,x_i})) \\
\gamma_{2n}(\bar{a}) &= \{(X^a, X^b) \mid X^a, X^b \in \mathbb{R}^n \wedge \varphi_n(X^a) \wedge \varphi_n(X^b) \wedge \varphi_{2n}^{\delta}(X^a, X^b)\}
\end{aligned} \tag{5}$$

In the DiffPoly domain, for any deterministic function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ the abstract transformer $T_f^{\sharp} : \mathcal{A}_{2n} \rightarrow \mathcal{A}_{2m}$ is required to satisfy the following soundness condition for all abstract elements $\bar{a} \in \mathcal{A}_{2n}$ where $T_f : \wp(\mathbb{R}^{2n}) \rightarrow \wp(\mathbb{R}^{2m})$ defines the corresponding concrete transformer

$$T_f(\gamma_{2n}(\bar{a})) \subseteq \gamma_{2m}(T_f^{\sharp}(\bar{a})) \quad \text{where } \forall X \in \wp(\mathbb{R}^{2n}). T_f(X) = \{(f(X), f(Y)) \mid (X, Y) \in X\}$$

Next, we define abstract transformers for the DiffPoly domain.

4.2 DiffPoly ReLU Abstract Transformer

$\text{ReLU} : \mathbb{R} \rightarrow \mathbb{R}$ is defined as $\text{ReLU}(x) = \max(0, x)$. Let, $T_R^{\sharp} : \mathcal{A}_{2i} \rightarrow \mathcal{A}_{2i+2}$ be the abstract transformer that executes assignment statements $y_i^a \leftarrow \text{ReLU}(x_i^a), y_i^b \leftarrow \text{ReLU}(x_i^b)$. For $\bar{a} = [a_1, \dots, a_i] \in \mathcal{A}_{2i}$, let $\bar{a}' = T_R^{\sharp}(\bar{a})$ represent the output of the transformer. First, for $\bar{a}' = [a'_1, \dots, a'_{i+1}]$, we compute the symbolic bounds $C_{sym}^{j,j}$ for each a'_j where $j \in [i+1]$. In this case, for all $j \in [i]. a'_j = a_j$ and a'_{i+1} is associated with the variable pair $< y_i^a, y_i^b >$. Since ReLU is piecewise linear, we separately analyze cases where ReLU acts as a linear function and cases where it demonstrates non-linear behavior. Table 1 summarizes the separate cases we consider while designing the abstract transformer for ReLU . In Table 1, for any variable v , v_+ (or, v_-) denotes the case when values taken by v are always positive (or negative) and v_{\pm} denotes the case when v can be both positive and negative. **Symbolic bounds for $(y_i^a - y_i^b)$.** We first consider cases where at least one of $\text{ReLU}(x_i^a)$ or $\text{ReLU}(x_i^b)$ behaves as a linear function and separately consider the case where both of them are non-linear. Similarly, we consider 3 scenarios based on the concrete bounds of $\delta_{x_i}^{a,b} = x_i^a - x_i^b$ (shown in Fig. 4) where we characterize the convex region having a minimum area that captures all possible values of $(y_i^a - y_i^b)$. In Table 2, we show the computation of the symbolic bounds for $(y_i^a - y_i^b)$ based on the cases for x_i^a and x_i^b . The first column shows the case, the second column shows the symbolic expression for $(y_i^a - y_i^b)$, and the last column shows its symbolic bounds. For the first four

Table 1. DiffPoly ReLU Cases

Cases from x_i^a	$x_-^{a,i} = (u_{a,x_i} \leq 0)$	$x_+^{a,i} = (l_{a,x_i} \geq 0)$	$x_{\pm}^{a,i} = \neg x_-^{a,i} \wedge \neg x_+^{a,i}$
Cases from x_i^b	$x_-^{b,i} = (u_{b,x_i} \leq 0)$	$x_+^{b,i} = (l_{b,x_i} \geq 0)$	$x_{\pm}^{b,i} = \neg x_-^{b,i} \wedge \neg x_+^{b,i}$
Cases from $\delta_{x_i}^{a,b}$	$\delta_-^i = (\Delta_{ub}^{a,b,x_i} \leq 0)$	$\delta_+^i = (\Delta_{lb}^{a,b,x_i} \geq 0)$	$\delta_{\pm}^i = \neg \delta_-^i \wedge \neg \delta_+^i$

Table 2. Computation of the symbolic bounds for $\delta_{y_i}^{a,b}$ based on cases for x_i^a and x_i^b .

Case	$\delta_{y_i}^{a,b}$	Symbolic bounds $\delta_{y_i}^{a,b,\leq}$ and $\delta_{y_i}^{a,b,\geq}$
$x_-^{a,i} \wedge x_-^{b,i}$	0	$(\delta_{y_i}^{a,b,\leq} = 0) \wedge (\delta_{y_i}^{a,b,\geq} = 0)$
$x_+^{a,i} \wedge x_+^{b,i}$	$x_i^a - x_i^b$	$(\delta_{y_i}^{a,b,\leq} = \delta_{x_i}^{a,b}) \wedge (\delta_{y_i}^{a,b,\geq} = \delta_{x_i}^{a,b})$
$x_+^{a,i} \wedge x_-^{b,i}$	x_i^a	$(\delta_{y_i}^{a,b,\leq} = x_i^a) \wedge (\delta_{y_i}^{a,b,\geq} = x_i^a)$
$x_-^{a,i} \wedge x_+^{b,i}$	$-x_i^b$	$(\delta_{y_i}^{a,b,\leq} = -x_i^b) \wedge (\delta_{y_i}^{a,b,\geq} = -x_i^b)$
$x_{\pm}^{a,i} \wedge x_{\pm}^{b,i}$	$ReLU(x_i^a)$	$(\delta_{y_i}^{a,b,\leq} = y_i^{a,\leq}) \wedge (\delta_{y_i}^{a,b,\geq} = y_i^{a,\geq})$
$x_-^{a,i} \wedge x_{\pm}^{b,i}$	$-ReLU(x_i^b)$	$(\delta_{y_i}^{a,b,\leq} = -y_i^{b,\geq}) \wedge (\delta_{y_i}^{a,b,\geq} = -y_i^{b,\leq})$
$x_{\pm}^{a,i} \wedge x_+^{b,i}$	$ReLU(x_i^a) - x_i^b$	$(\delta_{y_i}^{a,b,\leq} = y_i^{a,\leq} - x_i^b) \wedge (\delta_{y_i}^{a,b,\geq} = y_i^{a,\geq} - x_i^b)$
$x_+^{a,i} \wedge x_{\pm}^{b,i}$	$x_i^a - ReLU(x_i^b)$	$(\delta_{y_i}^{a,b,\leq} = x_i^a - y_i^{b,\geq}) \wedge (\delta_{y_i}^{a,b,\geq} = x_i^a - y_i^{b,\leq})$
$x_{\pm}^{a,i} \wedge x_{\pm}^{b,i}$	$ReLU(x_i^a) - ReLU(x_i^b)$	$(\delta_{y_i}^{a,b,\leq} = y_i^{a,\leq} - y_i^{b,\geq}) \wedge (\delta_{y_i}^{a,b,\geq} = y_i^{a,\geq} - y_i^{b,\leq})$

Table 3. Computation of the symbolic bounds for $\delta_{y_i}^{a,b}$ based on cases for $(x_i^a - x_i^b)$.

Case	Symbolic bounds $\delta_{y_i}^{a,b,\leq}$ and $\delta_{y_i}^{a,b,\geq}$ for ReLU activation
δ_+^i	$(\delta_{y_i}^{a,b,\leq} = 0) \wedge (\delta_{y_i}^{a,b,\geq} = \delta_{x_i}^{a,b})$
δ_-^i	$(\delta_{y_i}^{a,b,\leq} = \delta_{x_i}^{a,b}) \wedge (\delta_{y_i}^{a,b,\geq} = 0)$
δ_{\pm}^i	$(\delta_{y_i}^{a,b,\leq} = \lambda_{lb}^{\delta} \cdot \delta_{x_i}^{a,b} + \mu_{lb}^{\delta}) \wedge (\delta_{y_i}^{a,b,\geq} = \lambda_{ub}^{\delta} \cdot \delta_{x_i}^{a,b} + \mu_{ub}^{\delta})$ with $\lambda_{ub}^{\delta} = \frac{\Delta_{ub}^{a,b,x_i} - \Delta_{lb}^{a,b,x_i}}{\Delta_{ub}^{a,b,x_i} - \Delta_{lb}^{a,b,x_i}}, \lambda_{lb}^{\delta} = -\frac{\Delta_{lb}^{a,b,x_i}}{\Delta_{ub}^{a,b,x_i} - \Delta_{lb}^{a,b,x_i}}, -\mu_{ub}^{\delta} = \mu_{lb}^{\delta} = \frac{\Delta_{lb}^{a,b,x_i} \times \Delta_{ub}^{a,b,x_i}}{\Delta_{ub}^{a,b,x_i} - \Delta_{lb}^{a,b,x_i}}$

cases, $ReLU(x_i^a) - ReLU(x_i^b)$ behaves as a linear function and therefore our symbolic bounds are exact. For the remaining 5 cases, we compute symbolic bounds for $(y_i^a - y_i^b)$ overapproximating the exact values based on the symbolic bounds of y_i^a, y_i^b, x_i^a and x_i^b . We also consider 3 separate cases depicted in Table 3 (and in Fig 4) based on concrete bounds of $(x_i^a - x_i^b)$ where $\delta_{y_i}^{a,b,\leq}$ and $\delta_{y_i}^{a,b,\geq}$ are linear function of $\delta_{x_i}^{a,b} = (x_i^a - x_i^b)$. The cases described above are not mutually exclusive, resulting in multiple symbolic bound choices for $(y_i^a - y_i^b)$. However, in DiffPoly, we only allow a single symbolic upper bound and a lower bound for $(y_i^a - y_i^b)$. To resolve this, as described in Section 3, we greedily select the symbolic bounds that yield more precise concrete bounds based on concrete substitution (see Eq. 7). For example, consider the case specified by $(x_{\pm}^{a,i} \wedge x_{\pm}^{b,i} \wedge \delta_+)$ there are two choices for $\delta_{y_i}^{a,b,\geq} = y_i^{a,\geq} - y_i^{b,\leq}$ and $\delta_{y_i}^{a,b,\geq} = \delta_{x_i}^{a,b}$. Let, $S_c(y_i^{a,\geq} - y_i^{b,\leq})$ and $S_c(\delta_{x_i}^{a,b})$ be their respective concrete upper bounds. Then we pick $\delta_{y_i}^{a,b,\geq} = y_i^{a,\geq} - y_i^{b,\leq}$ if $S_c(y_i^{a,\geq} - y_i^{b,\leq}) < S_c(\delta_{x_i}^{a,b})$ otherwise select $\delta_{y_i}^{a,b,\geq} = \delta_{x_i}^{a,b}$. Next, we discuss symbolic bound computation for y_i^a and y_i^b .

Symbolic bounds for y_i^a and y_i^b . For cases $x_{\pm}^{a,i}$ and $x_{\pm}^{b,i}$ where the $ReLU$ behaves like a linear function, the symbolic bounds for y_i^a can be directly expressed as a linear function of x_i^a . However, for the case, $x_{\pm}^{a,i}$ the $ReLU$ function is no longer linear and we apply the linear relaxation [69, 92] to obtain the symbolic bounds of y_i^a using the concrete bounds l_{a,x_i} and u_{a,x_i} . The details are in the Appendix (Fig. 13). Bounds for y_i^b are derived similarly.

Concrete bounds for y_i^a, y_i^b . We get concrete bounds for y_i^a, y_i^b from the product DNN execution.

Concrete bounds for $(y_i^a - y_i^b)$. For $(y_i^a - y_i^b)$, we find concrete bounds using *back-substitution*. Each *back-substitution* step recursively applies symbolic substitution (Eq. 6) followed by concrete substitution (Eq. 7) to generate a set of possible candidates for concrete bounds and picks the most precise one. We provide a pseudo-code of the back-substitution algorithm in Appendix E. For any variable δ , its symbolic upper bound $\delta^{\geq} = \overline{v_0} + \sum_i \overline{w_i} \cdot \overline{v_i}$ and symbolic lower bound $\delta^{\leq} = \underline{v_0} + \sum_i \underline{w_i} \cdot \underline{v_i}$, the symbolic substitutions $S_s(\delta^{\geq})$, $S_s(\delta^{\leq})$ and concrete substitutions $S_c(\delta^{\geq})$, $S_c(\delta^{\leq})$ are shown below. Here, $\overline{v_0}, \underline{v_0} \in \mathbb{R}$ and $\overline{v_i}^{\geq}, \overline{v_i}^{\leq}, \underline{v_i}^{\geq}, \underline{v_i}^{\leq}$ are symbolic bounds of variables, $\overline{v_i}^{lb}, \overline{v_i}^{ub}, \underline{v_i}^{lb}, \underline{v_i}^{ub}$ are the respective concrete bounds and $\overline{w_i}^+ = \max(0, \overline{w_i})$, $\overline{w_i}^- = \min(0, \overline{w_i})$, $\underline{w_i}^+ = \max(0, \underline{w_i})$, $\underline{w_i}^- = \min(0, \underline{w_i})$. Note, both symbolic and concrete substitutions for upper and lower bounds satisfy that $(S_s(\delta^{\geq}) \geq \delta) \wedge (S_c(\delta^{\geq}) \geq \delta)$ and $(S_s(\delta^{\leq}) \leq \delta) \wedge (S_c(\delta^{\leq}) \leq \delta)$.

$$S_s(\delta^{\geq}) = \overline{v_0} + \sum_i \overline{w_i}^+ \cdot \overline{v_i}^{\geq} + \sum_i \overline{w_i}^- \cdot \overline{v_i}^{\leq} \quad S_s(\delta^{\leq}) = \underline{v_0} + \sum_i \underline{w_i}^+ \cdot \underline{v_i}^{\leq} + \sum_i \underline{w_i}^- \cdot \underline{v_i}^{\geq} \quad (6)$$

$$S_c(\delta^{\geq}) = \overline{v_0} + \sum_i \overline{w_i}^+ \cdot \overline{v_i}^{ub} + \sum_i \overline{w_i}^- \cdot \overline{v_i}^{lb} \quad S_c(\delta^{\leq}) = \underline{v_0} + \sum_i \underline{w_i}^+ \cdot \underline{v_i}^{lb} + \sum_i \underline{w_i}^- \cdot \underline{v_i}^{ub} \quad (7)$$

4.3 DiffPoly Abstract Transformer For Differentiable Activations

For any differentiable function $g: \mathbb{R} \rightarrow \mathbb{R}$, we define $T_g^{\#}: \mathcal{A}_{2i} \rightarrow \mathcal{A}_{2i+2}$ as the abstract transformer for the assignments $y_i^a \leftarrow g(x_i^a)$ and $y_i^b \leftarrow g(x_i^b)$. Both Sigmoid and Tanh, being differentiable everywhere, can be modeled via g . We use the lower bound and the upper bound on the derivative of g to compute the symbolic bounds of $(y_i^a - y_i^b)$. The concrete bounds of y_i^a and y_i^b are obtained from product DNN analysis while concrete bounds of $(y_i^a - y_i^b)$ are calculated by back-substitution.

Symbolic bounds computation: Let, $l_{g'}$ and $u_{g'}$ be the lower and upper bound of $g'(x)$ over the range $x \in [l, u]$ where $l = \min(l_{a,x_i}, l_{b,x_i})$ and $u = \max(u_{a,x_i}, u_{b,x_i})$. We consider three cases from the 3rd row of Table 1 and show the symbolic bounds of $(y_i^a - y_i^b)$ for all three cases in Table 4 (also depicted in Appendix Fig. 14). This formulation holds for any differentiable function g provided $l_{g'}$ and $u_{g'}$ are easy to compute. For Sigmoid and Tanh, the derivative $g'(x)$ has a closed form, and $g'(x)$ is maximum at $x = 0$ and decreases as x increases (or, decreases). So, $l_{g'}$ and $u_{g'}$ computation only takes constant time given values of l and u . For y_i^a and y_i^b , we use concrete bounds $-l_{a,x_i}, u_{a,x_i}, l_{b,x_i}, u_{b,x_i}$ and apply the linear relaxation from [92], which also extends to differentiable functions with a closed form of the differential.

Table 4. Computation of the symbolic bounds for $(y_i^a - y_i^b)$ based on cases for $(x_i^a - x_i^b)$.

Case	Symbolic bounds $\delta_{y_i}^{a,b,\leq}$ and $\delta_{y_i}^{a,b,\geq}$ for any differentiable activation g
δ_+^i	$(\delta_{y_i}^{a,b,\leq} = l_{g'} \cdot \delta_{x_i}^{a,b}) \wedge (\delta_{y_i}^{a,b,\geq} = u_{g'} \cdot \delta_{x_i}^{a,b})$
δ_-^i	$(\delta_{y_i}^{a,b,\leq} = u_{g'} \cdot \delta_{x_i}^{a,b}) \wedge (\delta_{y_i}^{a,b,\geq} = l_{g'} \cdot \delta_{x_i}^{a,b})$
δ_{\pm}^i	$(\delta_{y_i}^{a,b,\leq} = \lambda_{lb}^{\delta} \cdot \delta_{x_i}^{a,b} + \mu_{lb}^{\delta}) \wedge (\delta_{y_i}^{a,b,\geq} = \lambda_{ub}^{\delta} \cdot \delta_{x_i}^{a,b} + \mu_{ub}^{\delta})$ with $l_{g'} = \min_{x \in [l,u]} g'(x)$ and $u_{g'} = \max_{x \in [l,u]} g'(x)$ $\lambda_{ub}^{\delta} = \frac{u_{g'} \times \Delta_{ub}^{a,b,x_i} - l_{g'} \times \Delta_{lb}^{a,b,x_i}}{\Delta_{ub}^{a,b,x_i} - \Delta_{lb}^{a,b,x_i}}, \lambda_{lb}^{\delta} = \frac{l_{g'} \times \Delta_{ub}^{a,b,x_i} - u_{g'} \times \Delta_{lb}^{a,b,x_i}}{\Delta_{ub}^{a,b,x_i} - \Delta_{lb}^{a,b,x_i}}, -\mu_{ub}^{\delta} = \mu_{lb}^{\delta} = \frac{(u_{g'} - l_{g'}) \times \Delta_{lb}^{a,b,x_i} \times \Delta_{ub}^{a,b,x_i}}{\Delta_{ub}^{a,b,x_i} - \Delta_{lb}^{a,b,x_i}}$

4.4 DiffPoly Affine Abstract Transformer

We describe the affine abstract transformer $T_A^{\#}: \mathcal{A}_{2i} \rightarrow \mathcal{A}_{2i+2}$ corresponding to the assignment statements $x_{i+1}^a \leftarrow v + \sum_{j=1}^i w_j \cdot x_j^a$ and $x_{i+1}^b \leftarrow v + \sum_{j=1}^i w_j \cdot x_j^b$ where v and all w_j are real numbers. In this case, the difference $(x_{i+1}^a - x_{i+1}^b)$ can be represented as $(x_{i+1}^a - x_{i+1}^b) = \sum_{j=1}^i w_j \cdot (x_j^a - x_j^b)$. Since for affine assignments, x_{i+1}^a (and x_{i+1}^b) is a linear function over x_j^a s (and x_j^b s), we can directly compute the

linear constraints that represent the symbolic bounds. For $\bar{a} \in \mathcal{A}_{2i}$, let $\bar{a}' = [a'_1, \dots, a'_{i+1}] = T_A^\#(\bar{a})$ where $\bar{a}' \in \mathcal{A}_{2i+2}$ and $\forall j \in [i]$. ($a_j = a'_j$). We show the symbolic bounds corresponding to a'_{i+1} in Eq. 8. The product DNN analysis provides the concrete bounds of x_{i+1}^a and x_{i+1}^b while $\Delta_{lb}^{a,b,x_{i+1}}$ and $\Delta_{ub}^{a,b,x_{i+1}}$ are calculated by performing back-substitution on $\delta_{x_{i+1}}^{a,b,\leq}$ and $\delta_{x_{i+1}}^{a,b,\geq}$ respectively.

$$x_{i+1}^{a,\leq} = x_{i+1}^{a,\geq} = v + \sum_{j=1}^i w_j \cdot x_j^a \quad x_{i+1}^{b,\leq} = x_{i+1}^{b,\geq} = v + \sum_{j=1}^i w_j \cdot x_j^b \quad \delta_{x_{i+1}}^{a,b,\leq} = \delta_{x_{i+1}}^{a,b,\geq} = \sum_{j=1}^i w_j \cdot \delta_{x_j}^{a,b} \quad (8)$$

DiffPoly vs DeepPoly with transformer for the difference of activations: In Section 3.1.5, we explain why the existing DeepPoly domain is not suited for difference-bound computation between the outputs of a pair of DNN executions. It is natural to ask whether the precision improvement in difference tracking achieved by DiffPoly can be replicated by just designing a new abstract transformer for the DeepPoly domain handling the following assignments $y_i^a \leftarrow \sigma(x_i^a)$, $y_i^b \leftarrow \sigma(x_i^b)$ and $(y_i^a - y_i^b) \leftarrow \sigma(x_i^a) - \sigma(x_i^b)$ where $\sigma: \mathbb{R} \rightarrow \mathbb{R}$ is the non-linear activation function. In this case, the DeepPoly domain lacks concrete, symbolic bounds on the difference $(x_i^a - x_i^b)$ and can only use the concrete, symbolic bounds of the individual variables x_i^a, x_i^b . This results in imprecise concrete bounds Δ_{lb}^{a,b,y_i} and Δ_{ub}^{a,b,y_i} of $(y_i^a - y_i^b)$ which in turn results in imprecise symbolic bounds (Table 3 and 4 uses the sign of the concrete bounds of difference for selecting the symbolic bounds). For instance, in the illustrative example, the symbolic upper bound of $(\delta_{x_5}^{1,2})$ with DeepPoly bounds results in concrete upper bound $\Delta_{ub}^{1,2,x_5} = 20.625$ while DiffPoly produces more precise concrete upper bound $\Delta_{ub}^{1,2,x_5} = 6.0$. Overall DiffPoly is more general and can precisely handle bivariate non-linear functions such as $\sigma(x) - \sigma(y)$ with inputs x, y coming from two distinct copies of the network. Furthermore, we demonstrate in Appendix G.5 that DiffPoly can be expanded to encompass any linear combination of variables from k executions. This makes DiffPoly the first domain capable of computing precise bounds (both concrete and symbolic) of any linear combination of DNN outputs at each layer coming from different related executions.

4.5 RaVeN's Layerwise Constraint Formulation

In this section, we formally introduce RaVeN's layerwise constraint formulation. Consider $\mathcal{A}_\Delta = [\mathcal{A}_\delta^1, \dots, \mathcal{A}_\delta^\kappa]^T$, that stores the symbolic and concrete bounds computed by all κ DiffPoly analyses, with \mathcal{A}_δ^j representing the bounds computed by the j -th analysis. RaVeN's constraint formulation algorithm takes as input \mathcal{A}_Δ , network $N: \mathbb{R}^{n_0} \rightarrow \mathbb{R}^{n_l}$, and the input specification Φ and generates a set of linear constraints for each layer. Let, \mathcal{L}^i represent the set of linear constraints over the outputs of the i -th layer, defining the convex region $\mathcal{L}_t^i \subseteq \mathbb{R}^{n_i \times k}$. In this case, \mathcal{L}_t^i contains all possible outputs at i -th layer for all k executions. We compute \mathcal{L}^i by adding linear constraints for all n_i variables at the i -th layer for each pair of executions using the concrete and symbolic bounds from the DiffPoly analysis for that pair. For instance, consider $a \in [k] \wedge b \in [k] \wedge (a < b)$, which defines a pair of executions. Here, $[x_1^a, \dots, x_{n_i}^a]$ and $[x_1^b, \dots, x_{n_i}^b]$ represent variables at the i -th layer for the pair of executions (a, b) . Then the linear constraints added for this pair of executions are as follows where $j \in [n_i]$ and the concrete and symbolic bounds are from the DiffPoly analysis which in turn inherits the concrete bounds $l_{a,x_j}, u_{a,x_j}, l_{b,x_j}, u_{b,x_j}$ from product DNN analysis:

$$\begin{array}{lll} x_j^{a,\leq} \leq x_j^a \leq x_j^{a,\geq} & x_j^{b,\leq} \leq x_j^b \leq x_j^{b,\geq} & \delta_{x_j}^{a,b,\leq} \leq (x_j^a - x_j^b) \leq \delta_{x_j}^{a,b,\geq} \\ l_{a,x_j} \leq x_j^a \leq u_{a,x_j} & l_{b,x_j} \leq x_j^b \leq u_{b,x_j} & \Delta_{lb}^{a,b,x_j} \leq (x_j^a - x_j^b) \leq \Delta_{ub}^{a,b,x_j} \end{array} \quad (9)$$

In Eq. 9, the third column illustrates the additional difference constraints added for a variable pair, while the remaining constraints constitute RaVeN's layerwise formulation, as elaborated in

Section 3.2.2. Note that, as discussed earlier, in DiffPoly analysis, up to two valid symbolic lower or upper bounds can be generated for each variable and their difference. For efficiency in concrete bounds computation with back-substitution, DiffPoly restricts to a single symbolic lower and upper bound. However, in the MILP formulation, all valid bounds are incorporated. The input specification Φ , defined as a conjunction of linear constraints over the inputs, is directly encoded as a set of linear constraints \mathcal{L}^0 at the input layer. The linear constraints for all l layers are then generated by aggregating layerwise constraints \mathcal{L}^i with input linear constraints \mathcal{L}^0 .

4.6 RaVeN MILP encoding

We provide the general encoding of Ψ as MILP objective for relational DNN specifications described in Section 2.1. We add the MILP encoding of Ψ to the layerwise constraints from Section 4.5 to formulate the MILP instance. Let Y_1, \dots, Y_k be the DNN's output for k executions, for all $i \in [m]$ and $j \in [n]$, $x_{i,j}$ and z_i be integer variables and for all $i' \in [k]$, $C_{i,j,i'} \in \mathbb{R}^{n_l}$ where m is the number of clauses in Ψ and n is number of literals in each clause (see Section 2.1). Then the MILP objective is as follows

$$\min_{(Y_1, \dots, Y_k)} \sum_{i=1}^m z_i \quad \text{s.t.} \quad x_{i,j} = \psi_{i,j}(Y_1, \dots, Y_k) = \left(\sum_{i'=1}^k C_{i,j,i'}^T Y_{i'} \geq 0 \right); z_i = \left(\sum_{j=1}^n x_{i,j} \geq 0 \right) \quad (10)$$

The proof of the correctness of the MILP formulation is in Appendix F.6. For the common properties (e.g. UAP, targeted-UAP, worst-case hamming distance, etc.) $m = k$, $n = n_l$ and the MILP objective introduces only $k \times (n_l + 1)$ integer variables where n_l is the output dimension of the DNN (Appendix G.4). Hence irrespective of the size of the network, the number of integer variables only depends on the number of executions k and n_l which is in general a small constant (i.e. 10 for commonly used MNIST and CIFAR10 networks). Since the number of integer variables is the primary bottleneck of MILP optimization, RaVeN scales to large DNNs by only introducing a small number of integer variables ($n_l + 1$) per execution. This differs from the naive MILP which introduces an integer variable at each activation and does not scale past even small networks containing a few hundred neurons. Besides decreasing the count of integer variables, RaVeN efficiently infers linear constraints for the MILP encoding that are sound while improving the precision of the over-approximated convex region (illustrated in Figure 5 of the paper). This requires - (i) recognizing that tracking the difference between the outputs of a pair of DNN executions helps in improving precision while maintaining scalability, and (ii) designing and leveraging DiffPoly analysis on $\binom{k}{2}$ pairs of executions while computing provably correct constraints across multiple executions.

4.7 Soundness Proof Sketch of RaVeN

In this section, we outline the soundness proof for various components of RaVeN. Detailed proofs are in Appendix F. We start with the soundness proofs of all DiffPoly transformers.

4.7.1 Soundness of DiffPoly ReLU transformer. We first state the lemmas required to prove the soundness of $T_R^\#$. Proofs of all cases shown in Fig. 4, Lemma 4.2, and 4.3 are in Appendix G.1.

LEMMA 4.2. (Correctness of symbolic bounds in Table 2 and 3) If $x_i^a \in [l_{a,x_i}, u_{a,x_i}]$, $x_i^b \in [l_{b,x_i}, u_{b,x_i}]$ and $\delta_{x_i}^{a,b} = (x_i^a - x_i^b) \in [\Delta_{lb}^{a,b,x_i}, \Delta_{ub}^{a,b,x_i}]$ and $\delta_{y_i}^{a,b} = \text{ReLU}(x_i^a) - \text{ReLU}(x_i^b)$ then $\delta_{y_i}^{a,b,\leq} \leq \delta_{y_i}^{a,b} \leq \delta_{y_i}^{a,b,\geq}$ where $\delta_{y_i}^{a,b,\leq}$ and $\delta_{y_i}^{a,b,\geq}$ defined in Table 2 and 3.

LEMMA 4.3. (Correctness of concrete bounds computed by the ReLU transformer) If $x_i^a \in [l_{a,x_i}, u_{a,x_i}]$, $x_i^b \in [l_{b,x_i}, u_{b,x_i}]$ and $\delta_{x_i}^{a,b} = (x_i^a - x_i^b) \in [\Delta_{lb}^{a,b,x_i}, \Delta_{ub}^{a,b,x_i}]$, $y_i^a = \text{ReLU}(x_i^a)$, $y_i^b = \text{ReLU}(x_i^b)$, $\delta_{y_i}^{a,b} =$

$y_i^a - y_i^b$ then $l_{a,y_i} \leq y_i^a \leq u_{a,y_i}$, $l_{b,y_i} \leq y_i^b \leq u_{b,y_i}$, and $\Delta_{lb}^{a,b,y_i} \leq \delta_{y_i}^{a,b} \leq \Delta_{ub}^{a,b,y_i}$ where Δ_{lb}^{a,b,y_i} and Δ_{ub}^{a,b,y_i} computed by applying back-substitution on $\delta_{y_i}^{a,b,\leq}$ and $\delta_{y_i}^{a,b,\geq}$ respectively.

The concrete transformer $T_R : \wp(\mathbb{R}^{2i}) \rightarrow \wp(\mathbb{R}^{2i+2})$ for the ReLU assignments $y_i^a \leftarrow \text{ReLU}(x_i^a)$, $y_i^b \leftarrow \text{ReLU}(x_i^b)$ is defined as $T_R(X) = \{([x_1^a, \dots, x_i^a, y_i^a]^T, [x_1^b, \dots, x_i^b, y_i^b]^T) \mid (X^a, X^b) \in X\}$ where $y_i^a = \text{ReLU}(x_i^a)$, $y_i^b = \text{ReLU}(x_i^b)$, $X \subseteq \mathbb{R}^{2i}$ and $X^a = [x_1^a, \dots, x_i^a]^T \in \mathbb{R}^i$, $X^b = [x_1^b, \dots, x_i^b]^T \in \mathbb{R}^i$.

THEOREM 4.4. (*Soundness of DiffPoly Relu Transformer*) For any abstract element $\bar{a} \in \mathcal{A}_{2i}$ $T_R(\gamma_{2i}(\bar{a})) \subseteq \gamma_{2i+2}(T_R^\#(\bar{a}))$.

PROOF. The proof is in Appendix F.1. □

4.7.2 Soundness of DiffPoly differentiable activation transformer. Proof of all the cases from Table. 4 are in Appendix G.2. Lemma F.1 proves the soundness of the symbolic bounds, while Lemma F.2 proves the soundness of concrete bounds. The comprehensive soundness proof for the DiffPoly's transformer for differentiable activations is in Appendix F.2.

4.7.3 Soundness of DiffPoly Affine transformer. Lemma F.4 proves the soundness of the symbolic bounds corresponding to the DiffPoly affine transformer, while Lemma F.5 proves the soundness of the corresponding concrete bounds. A comprehensive soundness proof for the DiffPoly affine transformer is in Appendix F.3.

4.7.4 Soundness of product DNN analysis. We prove that the output region $\mathbb{P} \subseteq \mathbb{R}^{n_l \times k}$ obtained by running existing DNN abstract interpreters e.g. [68] on each of k copies of N contains all possible output w.r.t all k executions on inputs satisfying Φ . Let, $\forall i \in [k]$ $\phi_{in}^i : \mathbb{R}^{n_0} \rightarrow \{\text{true}, \text{false}\}$ defines the L_∞ input region $\phi_{in}^i = \|X - X_i^*\|_\infty \leq \epsilon$ for each of k executions. Existing DNN abstract interpreters operate on these individual input regions ϕ_{in}^i and compute the overapproximated output region $\mathcal{P}_i \subseteq \mathbb{R}^{n_l}$ that satisfies $\forall X \in \mathbb{R}^{n_0}. \phi_{in}^i(X) \implies (N(X) \in \mathcal{P}_i)$. The output region $\mathbb{P} \subseteq \mathbb{R}^{n_l \times k}$ is the cross-product of all k output regions $\mathbb{P} = \times_{i=1}^k \mathcal{P}_i$. Now, we show that \mathbb{P} contains all possible outputs of $N^k(X)$ provided $X \in \mathbb{R}^{n_0} \times k$ satisfies Φ .

THEOREM 4.5. (*Soundness of Product DNN analysis*) $\forall (X_1, \dots, X_k) \in \mathbb{R}^{n_0 \times k}. \Phi((X_1, \dots, X_k)) \implies (N^k((X_1, \dots, X_k)) \in \mathbb{P})$.

PROOF. The proof is in Appendix F.4. □

4.7.5 Soundness of RaVeN MILP formulation. We prove that for all layer $i \in [L]$ the convex region $\mathcal{L}_i^t \subseteq \mathbb{R}^{n_i \times k}$ defined by the linear constraints \mathcal{L}^i contain all possible outputs at i -th layer for all k executions. For the input region, we show $\Phi_t \subseteq \mathcal{L}_1^0$.

THEOREM 4.6. (*Soundness of Linear constraints*) $\Phi_t \subseteq \mathcal{L}_t^0$ and $\forall i \in [L]. \forall X_1, \dots, X_k \in \mathbb{R}^{n_0}. \Phi(X_1, \dots, X_k) \implies (N^i(X_1), \dots, N^i(X_k)) \in \mathcal{L}_i^t$ where $N^i : \mathbb{R}^{n_0} \rightarrow \mathbb{R}^{n_i}$ is the composition of first i layers of the network N , $N^i = N_1 \circ \dots \circ N_i$.

PROOF. The proof is in Appendix F.5. □

4.8 Asymptotic Runtime Analysis

First, we describe the runtime analysis of DiffPoly. Let the original DNN have n neurons. Symbolic bound computations for each variable pair $\langle x_i^a, x_i^b \rangle$ at worst take $O(n)$ time. Overall, the worst-case complexity for symbolic bound computation for all variable pairs is $O(n^2)$. The back-substitution algorithm used for computing concrete bounds in the worst case explores $O(n)$ symbolic bounds before terminating. Obtaining the concrete bounds by substituting concrete values for all variables in each symbolic bound takes $O(n)$ time. The worst-case runtime for obtaining concrete bounds for

each variable pair is $O(n^2)$ and the asymptotic runtime of a single DiffPoly analysis is $O(n^3)$. Since we consider $\binom{k}{2}$ pairs of executions the total cost of DiffPoly analysis is $O(k^2 \cdot n^3)$. For product DNN analysis we use an existing DNN abstract interpreter for each of k copies of the original network N . We assume analyzing each copy of N takes C_N time. So analyzing the product DNN takes $k \cdot C_N$ time. For the MILP formulation, we add in the worst-case $O(k)$ of constraints per variable and the product DNN contains $O(k \cdot n)$ variables. Then the total size of the MILP in terms of the number of linear constraints is $O(k^2 \cdot n)$. Since we formulate the MILP using the constraints obtained from the DiffPoly analysis, in the worst case, MILP formulation takes $O(k^2 \cdot n^3)$. Suppose it takes C_M worst case time to optimize the MILP, then worst case time complexity of RaVeN is $O(k^2 \cdot n^3) + k \cdot C_N + C_M$. Note, C_M depends on the MILP encoding of Ψ which is the only source of integer variables in RaVeN's formulation.

5 EVALUATION

We evaluate the effectiveness of RaVeN on a wide range of relational properties and a diverse set of neural networks and datasets. We consider the following relational properties: UAP, targeted UAP, hamming distance, and monotonicity as formally defined in Appendix A.3. For UAP and Hamming Distance properties, we compare our method to the existing baselines highlighted above in Section 3. The first baseline we consider is individual verification (see Section 3.2.1) which is work by Khedr and Shoukry [40]. The second baseline is an instantiation of the work done by Zeng et al. [88] with state-of-the-art non-relational verifiers DeepZ [68] and DeepPoly [69] which we call I/O Formulation (see Section 3.2.1). For these properties, our experimental results indicate that RaVeN is always more precise than existing methods and can verify significantly more properties. For monotonicity, we compare our methods to two existing baselines Liu et al. [48] and Pasado [44].

5.1 Experimental Setup

Datasets. For UAP based experiments, we use the popular MNIST [45] and CIFAR10 [42] image datasets. We also use MNIST for the Hamming distance experiments. For our monotonicity experiments, we use the Boston Housing (BH) dataset [37] and the Adult dataset [8]. The BH dataset contains 12 housing attributes such as age, tax, rooms, etc. and the target is housing price. The Adult dataset contains 87 features such as age, education, marital status, etc.

Neural Networks. Table 5 shows the MNIST, CIFAR10, BH, and Adult neural network architectures used in our experiments. We use standard network architectures (Convolutional and Fully-connected) commonly seen in other neural network verification works [68, 69]. We consider networks trained with standard training, DiffAI [53], CROWN-IBP [90], projected gradient descent (PGD) [50], and a monotonicity training scheme [34].

Non-relational verifier. We instantiate both RaVeN and I/O Formulation with either DeepPoly or DeepZ. Although RaVeN works with other non-relational verifiers including SOTA "Branch and Bound" based verifiers like α , β -CROWN [79] and MNBaB [28]. We use DeepPoly or DeepZ because they are fast and widely used for initializing complete verifiers. For example, α , β -CROWN uses CROWN (equivalent to DeepPoly). We also compare RaVeN's performance with α , β -CROWN and MNBaB in Section 5.6.

Implementation Details. We implemented our method in Python with Pytorch V1.11 and Gurobi V10.0.3 as an off-the-shelf MILP solver. Our MNIST experiments were performed on an Intel(R) Core(TM) i7-12800HX @ 4.80 GHz with 16 GB of memory and the remainder of our experiments on an Intel(R) Core(TM) i9-9900KS CPU @ 4.00GHz with 64 GB of memory. Unless otherwise specified, we use DeepZ [68] to perform bound analysis on the product DNN and use the same verifier for the baselines. We use Gurobi with a timeout of 5 minutes to solve MILP problems.

Table 5. Network Information and Runtime (s) averaged over ϵ values considered in this paper

DATASET	MODEL	TYPE	TRAIN	# LAYERS	# PARAMS	IND. VERI.	I/O FORM.	RaVeN	MILP TIME
MNIST	IBP-SMALL	CONV	IBP	7	60K	0.04	0.12	1.98	1.01
	CONVSMALL	CONV	DIFFAI	7	80K	0.30	0.39	7.40	4.06
	IBP	CONV	IBP	9	400K	0.42	0.46	19.33	7.79
	CONVBIG	CONV	DIFFAI	13	1.8M	6.46	6.50	23.19	16.61
	HAMMING	FC	PGD	3	39K	0.04	0.14	2.21	2.02
CIFAR10	IBP-SMALL	CONV	IBP	7	60K	0.29	0.47	8.39	5.03
	CONVSMALL	CONV	DIFFAI	7	80K	0.44	0.57	12.59	6.61
	IBP	CONV	IBP	9	2.2M	36.44	36.56	200.16	161.66
	CONVBIG	CONV	DIFFAI	13	2.5 M	16.19	16.29	185.05	161.63
DATASET	MODEL	TYPE	TRAIN	# LAYERS	# PARAMS	LIU ET AL.	PASADO	RaVeN	DIFFPOLY
BH	12x1	FC	MONO	3	312	0.25	×	-	0.02
ADULT	10 x 10	FC	STANDARD	5	980	×	36.70	4.23	0.87

5.2 Relational Properties

The formal definitions for UAP, targeted UAP, and hamming distance given in Appendix A.3 involve verifying that there does not exist an attack that can change all DNN predictions on a given input set by perturbing all the inputs with a single perturbation. While RaVeN can handle this problem, it is pessimistic and perturbations of this nature, although dangerous, rarely occur in reality. Instead, we bound the worst-case accuracy of the neural network under a UAP attack. Formally, we report a the verified worst-case accuracy which is a lower bound (as RaVeN is incomplete) on a^* , the true worst-case accuracy. For network N and inputs X_1, \dots, X_k where $\forall v \in \mathbb{R}^{n_0}$ s.t. $\|v\|_p \leq \epsilon$, $\frac{1}{k} \sum_{i=1}^k (N(X_i + v) = Y_i) \geq a$ and Y_i is the correct label of X_i . Note that a result is better if it more tightly approximates a^* in this case since all presented methods are sound the best result is the one with the greatest value. For hamming distance, we perform a similar relaxation upper bounding the true worst case hamming distance. Thus, for hamming distance, smaller is better. For monotonicity, we are given a set of monotonic features and report the percentage of those features we can verify. For monotonicity, larger is better.

5.3 Universal Adversarial Perturbation Verification

We compare the performance of RaVeN vs the two baselines for worst-case accuracy under UAP attack on the MNIST and CIFAR10 networks. For each experiment, we verify a batch of 5 images. We repeat 20 times on randomly selected images, reporting the average worst-case accuracy. We use the standard ϵ values used in the literature [68, 69]. We additionally analyze RaVeN vs. baselines on the targeted UAP verification problem in Appendix H.1.

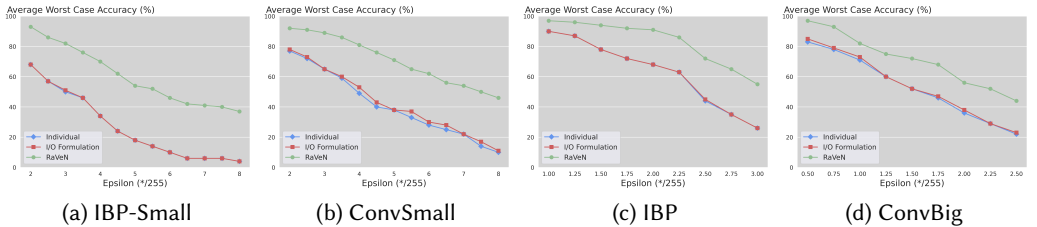


Fig. 6. Average worst case UAP accuracy for convolutional networks trained on CIFAR10

5.3.1 Comparison on CIFAR10 networks. Figure 6 compares the worst-case accuracy (%) on the CIFAR10 dataset with a variety of training methods (Crown-IBP, DiffAI) and network architectures

(IBP-Small, ConvSmall, IBP, ConvBig). We observe that RaVeN outperforms all baselines significantly for all networks, training methods, and ϵ s. For example, we see that for IBP-Small trained with Crown-IBP that RaVeN obtains at least 25% higher average worst case accuracy verified when compared to baselines on all ϵ s and a maximum of 38% higher accuracy at $\epsilon = 4.5$. On the same network, I/O Formulation, the SOTA UAP verification method, obtains at most 1% higher than the Individual baseline. For the IBP-Small network, even when the baselines achieve close to 0% at $\epsilon = 8/255$ RaVeN still obtains 37% accuracy. We observe similar results on the other networks.

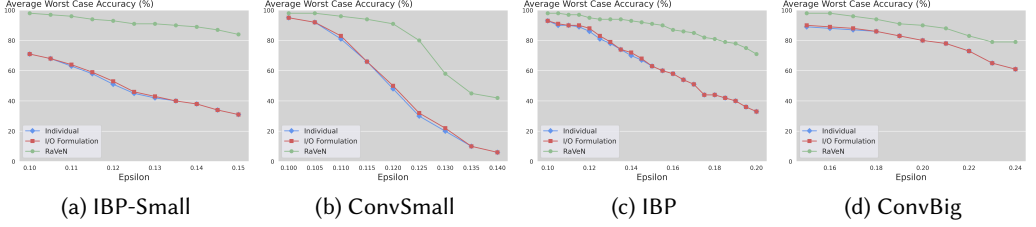


Fig. 7. Average Worst case UAP accuracy for convolutional networks trained on MNIST

5.3.2 Comparison on MNIST Networks. Figure 7 shows similar results to CIFAR10 with the same diverse range of networks and training methods. Particularly, we observe that for IBP-Small RaVeN verifies an additional 53% accuracy when compared to baselines at $\epsilon = 0.15$. We observe that as ϵ grows RaVeN’s relative benefit is greater, this is especially clear when for IBP (Figure 7 c).

5.3.3 Runtime Analysis. Table 5 shows the average runtime in seconds for each method. We observe that RaVeN time > I/O Formulation time > Independent Verification time. We note that even with more time the baseline approaches would not achieve any better results as they are limited and can not get more precise. Note that a majority of the time for RaVeN is taken by the MILP solver as seen in Table 5. As RaVeN is the first tool to show that cross-execution information aids in relational verification we believe runtime can be improved with future research. We also note that our timings are comparable to the timeouts given in the SOTA competition for verification of NNs (VNN-Comp [12]) (216 seconds per instance) even though we are verifying sets of 5 images.

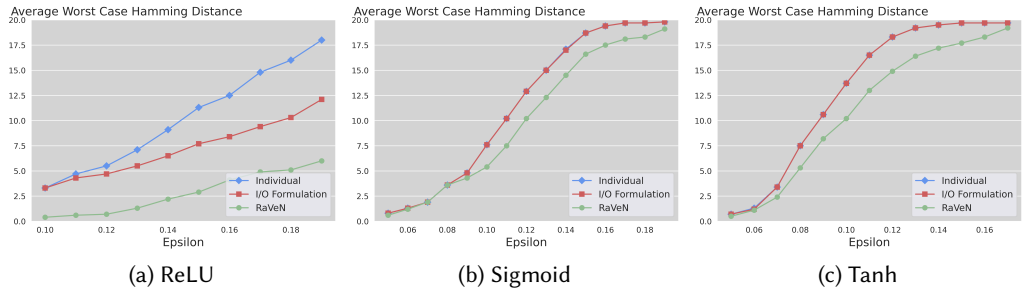


Fig. 8. Average Worst Case Hamming Distance with different activation functions (smaller is better)

5.4 Hamming Distance Verification

We use MNIST as the base dataset and train a 3-layer fully connected network with 200 neurons in the hidden layers. We use a range of activation functions (ReLU, Tanh, Sigmoid). The network is adversarially trained with PGD to identify between classes 0 and 1. In this experiment, DeepPoly is used to instantiate both the baselines and RaVeN. Figure 8 shows the worst case hamming distance for strings of length 20 for different activation functions and ϵ values. For all ϵ values and string lengths, RaVeN outperforms both baselines, e.g. at $\epsilon = 0.3$ for Tanh the baselines obtain 20 and 19.85 while RaVeN obtains 15. We especially see that for Sigmoid and Tanh activations the baselines perform identically while RaVeN significantly outperforms both of them.

5.5 Monotonicity Verification

We verify the monotonicity of networks with both Tanh and ReLU activations trained on the Adult [8] and BH [37] datasets respectively. We compare our methods against the SOTA monotonicity verifier for Tanh networks, Pasado [44] using the Adult dataset with 5 monotonic features (same features as previous works [44, 66]). Monotonicity can be verified directly by DiffPoly without the need for any MILP formulation. For incomplete verifiers such as RaVeN, imprecisions accumulate during the analysis. By splitting the input region and verifying each region separately we can get a sound analysis which is sometimes more precise than the original analysis with some additional computation cost. Input splitting is a common tool used in other verification papers as a way to increase precision [38]. We use input splitting for monotonicity for two reasons: 1. the monotonic input specification only has one dimension of variation and is thus easy to split, and 2. DiffPoly/RaVeN verifies monotonicity very quickly in comparison to SOTA methods so we can split to gain precision while still having faster runtime. For both RaVeN and DiffPoly we split the input region 10 times before verifying. Figure 9 shows the results of RaVeN and DiffPoly compared to Pasado and its baselines (Zonotope, Interval). For small ϵ Pasado slightly outperforms RaVeN (92% vs 94%); however, as ϵ grows the benefit of RaVeN becomes clear (66% vs 2% at $\epsilon = 4$). We observe that DiffPoly alone can perform on par with Pasado while running significantly faster (0.87s vs 36.7s, while RaVeN sits in the middle at 4.23s). For ReLU networks we compare against Liu et al. [48] as Pasado is unable to handle ReLU (Liu et al. [48] only handles ReLU). We verify a single feature on the Boston Housing dataset over the 98 test images. Liu et al. [48] can verify all 98 inputs for monotonicity for each $\epsilon = [10, 20, 30]$. On the other hand, DiffPoly is able to verify [96, 95, 95] inputs for $\epsilon = [10, 20, 30]$, but we note that DiffPoly is significantly faster (0.02s vs 0.25s). We observe that DiffPoly and RaVeN are powerful monotonicity verifiers that can handle a wider range of networks/activation functions than both baselines achieving good results in significantly less time.

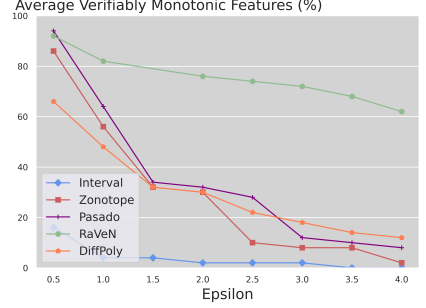


Fig. 9. Average % of Verified Monotonic Features on Adult Dataset

5.6 Ablation Studies

In this section, we show an ablation study comparing RaVeN to stronger individual verifiers: MNBaB [28] and α, β -CROWN [79]. We further show an ablation study on the benefits of adding difference constraints compared to only adding the layerwise formulation. In Appendix H.2, we show RaVeN performs well compared to baselines when all of them use DeepPoly [69] instead of DeepZ [68].

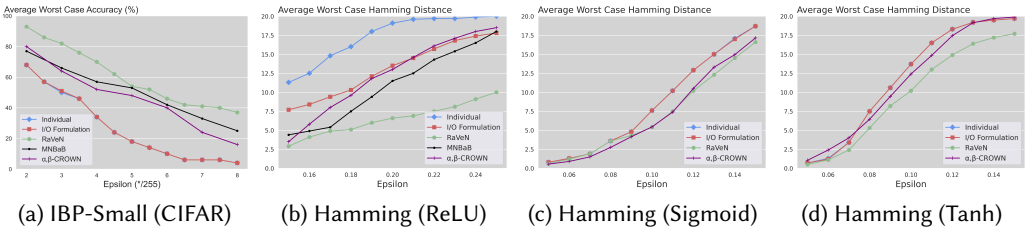


Fig. 10. Comparison of RaVeN against MNBaB and α, β -CROWN

5.6.1 Comparison to MNBaB and α, β -CROWN. MNBaB [28] and α, β -CROWN [79] use branching to obtain better precision at the cost of runtime. Although both MNBaB and α, β -CROWN are

complete for non-relational properties for DNNs with piece-wise linear activations such as ReLU, they are imprecise for relational verification as they do not take the cross-execution constraints into account. Furthermore, both MNBaB and α, β -CROWN cannot verify monotonicity, whereas both DiffPoly and RaVeN can handle monotonicity. We instantiate MNBaB and α, β -CROWN with a 2-minute timeout per individual input. Note that although RaVeN is given a timeout of 5 minutes for MILP solving, for individual verifiers to perform UAP verification they must individually verify each input in the batch giving MNBaB and α, β -CROWN a total of 10 and 40 minutes for UAP and hamming distance verification respectively. Figure 10 compares RaVeN to MNBaB and α, β -CROWN on UAP verification for IBP-Small on CIFAR10 and for hamming distance verification on MNIST with different activations. Note that MNBaB does not currently support Sigmoid or Tanh activations. Similar to the above experiments, we instantiate RaVeN with DeepZ for IBP-Small and DeepPoly for hamming distance networks. We observe that RaVeN consistently performs better than MNBaB and α, β -CROWN (except for the hamming distance network with sigmoid activations for small ϵ s). For example, for hamming distance with ReLU activations at $\epsilon = 0.25$, RaVeN can verify an average worst-case hamming distance of 10 while MNBaB and α, β -CROWN only obtain 18 and 18.5 respectively. For IBP-Small on CIFAR10 at $\epsilon = 8/255$, RaVeN can verify a worst-case UAP accuracy of 37% while MNBaB and α, β -CROWN only obtain 25% and 16% respectively.

In Table 6, we show a runtime comparison between RaVeN, MNBaB, and α, β -CROWN on the same networks as Figure 10. We observe that RaVeN takes less time than MNBaB and α, β -CROWN in all instances. Note that for Sigmoid and Tanh activations, α, β -CROWN is equivalent to α -CROWN [87] which does not support branching resulting in lower runtimes. In all instances, MNBaB and α, β -CROWN take significantly more time ($> 37.7\times$ more time for hamming distance with ReLU activations).

Table 6. Runtime Comparison (in secs) between RaVeN, MNBaB, and α, β -CROWN

DATASET	MODEL	ACTIVATION	RaVeN	MNBaB	α, β -CROWN
MNIST	HAMMING	ReLU	4.92	209.38	185.91
	HAMMING	SIGMOID	1.15	×	3.05
	HAMMING	TANH	2.37	×	5.77
CIFAR10	IBP-SMALL	ReLU	8.39	23.13	39.92
ADULT	10 × 10	TANH	4.23	×	×

5.6.2 Benefits of Difference Constraints. Figure 11 shows the benefits of adding difference constraints. In each example, RaVeN with difference constraints outperforms RaVeN layerwise without difference constraints. For example, for IBP-Small on CIFAR10 we see at $\epsilon = 8$ adding difference constraints increases the accuracy bound from 15% to 37%. The benefit of difference constraints is especially highlighted in the hamming distance example (d) as only by adding difference constraints is RaVeN able to outperform the baseline methods. A runtime comparison between RaVeN layerwise and RaVeN with difference constraints can be found in Appendix H.3.

6 RELATED WORK

DNN verifiers. Prior works in DNN verification [1] primarily focus on proving whether a DNN satisfies L_∞ robustness [69, 80] property. In this case, existing DNN verifiers show that all inputs inside a given L_∞ region [16] are properly classified. The DNN verifiers are broadly categorized into three main categories - (i) sound but incomplete verifiers which may not always prove property even if it holds [31, 63, 67–69, 86, 87], (ii) complete verifiers that can always prove the property if it holds [5, 13, 14, 25, 28, 30, 64, 71, 78, 79, 91] and (iii) verifiers with probabilistic guarantees

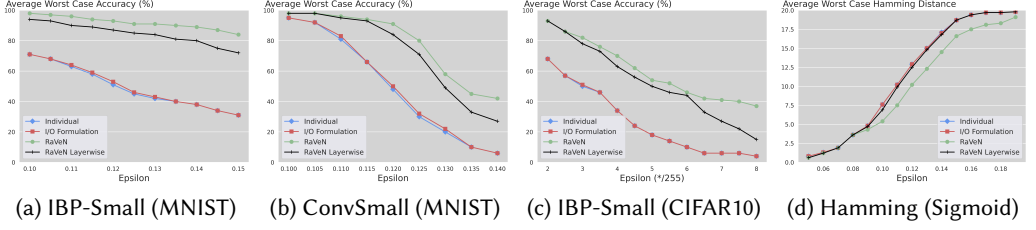


Fig. 11. Comparison of RaVeN with difference constraints with RaVeN with only layerwise formulation.

[19]. However, all of these verifiers verify properties defined over single DNN execution and are ineffective for verifying interesting relational properties [17] such as UAP verification [88] and monotonicity [74] defined over multiple DNN executions.

DNN relational verifiers. Existing DNN relational verifiers can be grouped into two main categories - (i) verifiers for relational properties (UAP, monotonicity, etc.) defined over multiple executions of the same DNN, [40, 88], (ii) verifiers for relational properties (local DNN equivalence [58]) defined over multiple executions of different DNN on the same input [58, 59]. For relational properties defined over multiple executions of the same DNN the existing verifiers [40] reduce the verification problem into L_∞ robustness problem by constructing product DNN with multiple copies of the same DNN. However, the relational verifier in [40] treats all k executions of the DNN as independent and loses precision. The state-of-the-art DNN relational verifier [88] although tracks the relationship between inputs used in multiple executions at the input layer, does not track the relationship between the inputs fed to the subsequent hidden layers and can only achieve a marginal improvement over the baseline verifiers that treat all executions independently. ITNE [81] is a verifier for global robustness based on difference tracking. Global robustness measures the largest change to the output of a single class over the entire dataset (local robustness lifted to the dataset) whereas the UAP property considered in this work focuses on the number of points a single perturbation can cause to misclassify over a set of inputs which can be from different classes. Furthermore, RaVeN is more precise (Eq. 6 in [81] is covered by Table 2, RaVeN gains precision by also considering the constraints in Table 3) and handles more activations than ITNE.

Relational verification of programs. Compared to DNNs, significantly more work exist for verifying different relational properties, such as information flow security, determinism, etc. on programs [7, 9, 11, 15, 18, 26, 27, 29, 41, 65, 73, 77]. Standard programs and DNNs have different computational structure. For example, programs have loops while DNNs have a large number of non-linear activations. These structural differences create specific challenges for the relational verification of DNNs not seen for programs and vice-versa.

7 CONCLUSION

In this work, we developed a new framework called RaVeN to verify the relational properties of DNNs based on our novel approach of difference tracking with the DiffPoly abstract domain. We run extensive experiments on multiple relational properties including UAP verification, monotonicity, etc., and show that RaVeN outperforms the state-of-the-art relational verifier [88] on all of them. We have primarily considered relational properties defined over multiple executions of the same DNN, however, RaVeN can be extended to relational properties involving two or more different DNNs - local equivalence of pair of DNNs [58], properties defined over an ensemble of DNNs, etc. RaVeN can also be integrated inside the training loop to obtain more trustworthy and safe neural networks. We leave this as future work. Also, the current implementation of RaVeN is sequential but as stated above certain steps like the product DNN analysis and pairwise difference computation with DiffPoly can be parallelized to reduce the verification cost.

ACKNOWLEDGEMENT

We thank the anonymous reviewers for their insightful comments. This work was supported in part by NSF Grants No. CCF-2238079, CCF-2316233, CNS-2148583, Google Research Scholar award, and Qualcomm Innovation Fellowship.

ARTIFACT STATEMENT

The artifact, on which the evaluation was done, is available at: <https://zenodo.org/records/10807316> with DOI [10.5281/zenodo.10807316](https://doi.org/10.5281/zenodo.10807316) [23]. The artifact includes instructions to reproduce the claimed results of the paper.

REFERENCES

- [1] Aws Albarghouthi. 2021. Introduction to Neural Network Verification. *Found. Trends Program. Lang.* 7, 1-2 (2021), 1–157. <https://doi.org/10.1561/25000000051>
- [2] Filippo Amato, Alberto López, Eladia María Peña-Méndez, Petr Vaňhara, Aleš Hampl, and Josef Havel. 2013. Artificial neural networks in medical diagnosis. *Journal of Applied Biomedicine* 11, 2 (2013).
- [3] Greg Anderson, Shankara Pailoor, Isil Dillig, and Swarat Chaudhuri. 2019. Optimization and Abstraction: A Synergistic Approach for Analyzing Neural Network Robustness. In *Proc. Programming Language Design and Implementation (PLDI)*, 731–744.
- [4] Stanley Bak, Taylor Dohmen, K. Subramani, Ashutosh Trivedi, Alvaro Velasquez, and Piotr Wojciechowski. 2023. The Octatope Abstract Domain for Verification of Neural Networks. In *Formal Methods - 25th International Symposium, FM 2023, Lübeck, Germany, March 6-10, 2023, Proceedings (Lecture Notes in Computer Science, Vol. 14000)*, Marsha Chechik, Joost-Pieter Katoen, and Martin Leucker (Eds.). Springer, 454–472. https://doi.org/10.1007/978-3-031-27481-7_26
- [5] Stanley Bak, Hoang-Dung Tran, Kerianne Hobbs, and Taylor T. Johnson. 2020. Improved Geometric Path Enumeration for Verifying ReLU Neural Networks. In *Computer Aided Verification - 32nd International Conference, CAV 2020, Los Angeles, CA, USA, July 21-24, 2020, Proceedings, Part I (Lecture Notes in Computer Science, Vol. 12224)*, Shuvendu K. Lahiri and Chao Wang (Eds.). Springer, 66–96. https://doi.org/10.1007/978-3-030-53288-8_4
- [6] Mislav Balunovic, Maximilian Baader, Gagandeep Singh, Timon Gehr, and Martin Vechev. 2019. Certifying Geometric Robustness of Neural Networks. In *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (Eds.), Vol. 32. Curran Associates, Inc. https://proceedings.neurips.cc/paper_files/paper/2019/file/f7fa6aca028e7ff4ef62d75ed025fe76-Paper.pdf
- [7] G. Barthe, P.R. D'Argenio, and T. Rezk. 2004. Secure information flow by self-composition. In *Proceedings. 17th IEEE Computer Security Foundations Workshop, 2004*. 100–114. <https://doi.org/10.1109/CSFW.2004.1310735>
- [8] Barry Becker and Ronny Kohavi. 1996. Adult. UCI Machine Learning Repository. DOI: <https://doi.org/10.24432/C5XW20>.
- [9] Raven Beutner and Bernd Finkbeiner. 2022. Software Verification of Hyperproperties Beyond k-Safety. In *Computer Aided Verification - 34th International Conference, CAV 2022, Haifa, Israel, August 7-10, 2022, Proceedings, Part I (Lecture Notes in Computer Science, Vol. 13371)*, Sharon Shoham and Yakir Vizel (Eds.). Springer, 341–362.
- [10] Mariusz Bojarski, Davide Del Testa, Daniel Dworakowski, Bernhard Firner, Beat Flepp, Prasoon Goyal, Lawrence D Jackel, Mathew Monfort, Urs Muller, Jiakai Zhang, et al. 2016. End to end learning for self-driving cars. *arXiv preprint arXiv:1604.07316* (2016).
- [11] Laura Bozzelli, Adriano Peron, and César Sánchez. 2021. Asynchronous Extensions of HyperLTL. In *36th Annual ACM/IEEE Symposium on Logic in Computer Science, LICS 2021, Rome, Italy, June 29 - July 2, 2021*. IEEE, 1–13. <https://doi.org/10.1109/LICS52264.2021.9470583>
- [12] Christopher Brix, Mark Niklas Müller, Stanley Bak, Taylor T Johnson, and Changliu Liu. 2023. First three years of the international verification of neural networks competition (VNN-COMP). *International Journal on Software Tools for Technology Transfer* (2023), 1–11.
- [13] Rudy Bunel, Jingyue Lu, Ilker Turkaslan, Pushmeet Kohli, P Torr, and P Mudigonda. 2020. Branch and bound for piecewise linear neural network verification. *Journal of Machine Learning Research* 21, 2020 (2020).
- [14] Rudy R Bunel, Oliver Hinder, Srinadh Bhojanapalli, and Krishnamurthy Dvijotham. 2020. An efficient nonconvex reformulation of stagewise convex optimization problems. *Advances in Neural Information Processing Systems* 33 (2020).
- [15] Jacob Burnim and Koushik Sen. 2009. Asserting and Checking Determinism for Multithreaded Programs. In *Proceedings of the 7th Joint Meeting of the European Software Engineering Conference and the ACM SIGSOFT Symposium on The Foundations of Software Engineering (Amsterdam, The Netherlands) (ESEC/FSE '09)*. Association for Computing Machinery, New York, NY, USA, 3–12. <https://doi.org/10.1145/1595696.1595700>

- [16] Nicholas Carlini and David Wagner. 2017. Towards evaluating the robustness of neural networks. In *2017 IEEE Symposium on Security and Privacy (SP)*. Ieee, 39–57.
- [17] Maria Christakis, Hasan Ferit Eniser, Jörg Hoffmann, Adish Singla, and Valentin Wüstholtz. 2022. Specifying and Testing k -Safety Properties for Machine-Learning Models. *arXiv preprint arXiv:2206.06054* (2022).
- [18] Berkeley R. Churchill, Oded Padon, Rahul Sharma, and Alex Aiken. 2019. Semantic program alignment for equivalence checking. In *Proceedings of the 40th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI 2019, Phoenix, AZ, USA, June 22-26, 2019*, Kathryn S. McKinley and Kathleen Fisher (Eds.). ACM, 1027–1040. <https://doi.org/10.1145/3314221.3314596>
- [19] Jeremy Cohen, Elan Rosenfeld, and Zico Kolter. 2019. Certified Adversarial Robustness via Randomized Smoothing. In *Proceedings of the 36th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 97)*, Kamalika Chaudhuri and Ruslan Salakhutdinov (Eds.). PMLR, 1310–1320. <https://proceedings.mlr.press/v97/cohen19c.html>
- [20] Patrick Cousot and Nicolas Halbwachs. 1978. Automatic Discovery of Linear Restraints among Variables of a Program. In *Proceedings of the 5th ACM SIGACT-SIGPLAN Symposium on Principles of Programming Languages* (Tucson, Arizona) (POPL '78). Association for Computing Machinery, New York, NY, USA, 84–96. <https://doi.org/10.1145/512760.512770>
- [21] Hennie Daniels and Marina Velikova. 2010. Monotone and partially monotone neural networks. *IEEE Transactions on Neural Networks* 21, 6 (2010), 906–917.
- [22] Sumanth Dathathri, Krishnamurthy Dvijotham, Alexey Kurakin, Aditi Raghunathan, Jonathan Uesato, Rudy Bunel, Shreya Shankar, Jacob Steinhardt, Ian J. Goodfellow, Percy Liang, and Pushmeet Kohli. 2020. Enabling certification of verification-agnostic networks via memory-efficient semidefinite programming. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (Eds.). <https://proceedings.neurips.cc/paper/2020/hash/397d6b4c83c91021fe928a8c4220386b-Abstract.html>
- [23] D.Banerjee. 2024. RaVeN: v1.1. <https://doi.org/10.5281/zenodo.10807316>
- [24] Hai Duong, Linhan Li, ThanhVu Nguyen, and Matthew Dwyer. 2023. A DPLL (T) Framework for Verifying Deep Neural Networks. *arXiv preprint arXiv:2307.10266* (2023).
- [25] Ruediger Ehlers. 2017. Formal verification of piece-wise linear feed-forward neural networks. In *International Symposium on Automated Technology for Verification and Analysis*.
- [26] Marco Eilers, Peter Müller, and Samuel Hitz. 2018. Modular Product Programs. In *Programming Languages and Systems - 27th European Symposium on Programming, ESOP 2018, Held as Part of the European Joint Conferences on Theory and Practice of Software, ETAPS 2018, Thessaloniki, Greece, April 14-20, 2018, Proceedings (Lecture Notes in Computer Science, Vol. 10801)*, Amal Ahmed (Ed.). Springer, 502–529. https://doi.org/10.1007/978-3-319-89884-1_18
- [27] Azadeh Farzan and Anthony Vandikas. 2019. Automated Hypersafety Verification. In *Computer Aided Verification - 31st International Conference, CAV 2019, New York City, NY, USA, July 15-18, 2019, Proceedings, Part I (Lecture Notes in Computer Science, Vol. 11561)*, Isil Dillig and Serdar Tasiran (Eds.). Springer, 200–218. https://doi.org/10.1007/978-3-030-25540-4_11
- [28] Claudio Ferrari, Mark Niklas Mueller, Nikola Jovanović, and Martin Vechev. 2022. Complete Verification via Multi-Neuron Relaxation Guided Branch-and-Bound. In *International Conference on Learning Representations*. https://openreview.net/forum?id=l_amHf1oaK
- [29] Bernd Finkbeiner, Markus N. Rabe, and César Sánchez. 2015. Algorithms for Model Checking HyperLTL and HyperCTL *. In *Computer Aided Verification - 27th International Conference, CAV 2015, San Francisco, CA, USA, July 18-24, 2015, Proceedings, Part I (Lecture Notes in Computer Science, Vol. 9206)*, Daniel Kroening and Corina S. Pasareanu (Eds.). Springer, 30–48. https://doi.org/10.1007/978-3-319-21690-4_3
- [30] Aymeric Fromherz, Klas Leino, Matt Fredrikson, Bryan Parno, and Corina Pasareanu. 2021. Fast Geometric Projections for Local Robustness Certification. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=zWy1uxjDdZJ>
- [31] Timon Gehr, Matthew Mirman, Dana Drachler-Cohen, Petar Tsankov, Swarat Chaudhuri, and Martin Vechev. 2018. Ai2: Safety and robustness certification of neural networks with abstract interpretation. In *2018 IEEE Symposium on Security and Privacy (SP)*.
- [32] Chuqin Geng, Nham Le, Xiaojie Xu, Zhaoyue Wang, Arie Gurfinkel, and Xujie Si. 2023. Towards reliable neural specifications. In *International Conference on Machine Learning*. PMLR, 11196–11212.
- [33] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. 2014. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572* (2014).
- [34] Akhil Gupta, Naman Shukla, Lavanya Marla, Arinbjörn Kolbeinsson, and Kartik Yellepeddi. 2019. How to incorporate monotonicity in deep networks while preserving flexibility? *arXiv preprint arXiv:1909.10662* (2019).
- [35] Gurobi Optimization, LLC. 2018. Gurobi Optimizer Reference Manual.

- [36] Richard W Hamming. 1950. Error detecting and error correcting codes. *The Bell system technical journal* 29, 2 (1950), 147–160.
- [37] David Harrison Jr and Daniel L Rubinfeld. 1978. Hedonic housing prices and the demand for clean air. *Journal of environmental economics and management* 5, 1 (1978), 81–102.
- [38] Anan Kabaha and Dana Drachler-Cohen. 2022. Boosting Robustness Verification of Semantic Feature Neighborhoods. In *Static Analysis - 29th International Symposium, SAS 2022, Auckland, New Zealand, December 5-7, 2022, Proceedings (Lecture Notes in Computer Science, Vol. 13790)*, Gagandeep Singh and Caterina Urban (Eds.). Springer, 299–324. https://doi.org/10.1007/978-3-031-22308-2_14
- [39] Guy Katz, Derek Huang, Duligur Ibeling, Kyle Julian, Christopher Lazarus, Rachel Lim, Parth Shah, Shantanu Thakoor, Haoze Wu, Aleksandar Zeljić, David Dill, Mykel Kochenderfer, and Clark Barrett. 2019. *The Marabou Framework for Verification and Analysis of Deep Neural Networks*. 443–452.
- [40] Haitham Khedr and Yasser Shoukry. 2023. CertiFair: A Framework for Certified Global Fairness of Neural Networks. *Proceedings of the AAAI Conference on Artificial Intelligence* 37, 7 (Jun. 2023), 8237–8245.
- [41] Máté Kovács, Helmut Seidl, and Bernd Finkbeiner. 2013. Relational abstract interpretation for the verification of 2-hypersafety properties. In *2013 ACM SIGSAC Conference on Computer and Communications Security, CCS’13, Berlin, Germany, November 4-8, 2013*, Ahmad-Reza Sadeghi, Virgil D. Gligor, and Moti Yung (Eds.). ACM, 211–222. <https://doi.org/10.1145/2508859.2516721>
- [42] Alex Krizhevsky. 2009. Learning Multiple Layers of Features from Tiny Images. (2009).
- [43] Jianglin Lan, Yang Zheng, and Alessio Lomuscio. 2022. Tight Neural Network Verification via Semidefinite Relaxations and Linear Reformulations. In *Thirty-Sixth AAAI Conference on Artificial Intelligence, AAAI 2022, Thirty-Fourth Conference on Innovative Applications of Artificial Intelligence, IAAI 2022, The Twelveth Symposium on Educational Advances in Artificial Intelligence, EAAI 2022 Virtual Event, February 22 - March 1, 2022*. AAAI Press, 7272–7280. <https://ojs.aaai.org/index.php/AAAI/article/view/20689>
- [44] Jacob Laurel, Siyuan Brant Qian, Gagandeep Singh, and Sasa Misailovic. 2023. Synthesizing precise static analyzers for automatic differentiation. *Proceedings of the ACM on Programming Languages* 7, OOPSLA2 (2023), 1964–1992.
- [45] Yann LeCun, Bernhard E. Boser, John S. Denker, Donnie Henderson, Richard E. Howard, Wayne E. Hubbard, and Lawrence D. Jackel. 1989. Handwritten Digit Recognition with a Back-Propagation Network. In *NIPS*. 396–404.
- [46] Juncheng Li, Shuhui Qu, Xinjian Li, Joseph Szurley, J. Zico Kolter, and Florian Metze. 2019. Adversarial Music: Real world Audio Adversary against Wake-word Detection System. In *Proc. Neural Information Processing Systems (NeurIPS)*. 11908–11918.
- [47] Juncheng Li, Frank R. Schmidt, and J. Zico Kolter. 2019. Adversarial camera stickers: A physical camera-based attack on deep learning systems. In *Proc. International Conference on Machine Learning, ICML*, Vol. 97. 3896–3904.
- [48] Xingchao Liu, Xing Han, Na Zhang, and Qiang Liu. 2020. Certified monotonic neural networks. *Advances in Neural Information Processing Systems* 33 (2020), 15427–15438.
- [49] Zikun Liu, Changming Xu, Emerson Sie, Gagandeep Singh, and Deepak Vasisht. 2023. Exploring Practical Vulnerabilities of Machine Learning-based Wireless Systems. In *20th USENIX Symposium on Networked Systems Design and Implementation, NSDI 2023, Boston, MA, April 17-19, 2023*. USENIX Association, 1801–1817.
- [50] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. 2018. Towards Deep Learning Models Resistant to Adversarial Attacks. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=rJzIBfZAb>
- [51] Antoine Miné. 2001. A new numerical abstract domain based on difference-bound matrices. In *Programs as Data Objects: Second Symposium, PADO2001 Aarhus, Denmark, May 21–23, 2001 Proceedings*. Springer, 155–172.
- [52] Matthew Mirman, Timon Gehr, and Martin Vechev. 2018. Differentiable abstract interpretation for provably robust neural networks. In *Proc. International Conference on Machine Learning (ICML)*. 3578–3586.
- [53] Matthew Mirman, Timon Gehr, and Martin Vechev. 2018. Differentiable Abstract Interpretation for Provably Robust Neural Networks. In *Proceedings of the 35th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 80)*, Jennifer Dy and Andreas Krause (Eds.). PMLR, 3578–3586. <https://proceedings.mlr.press/v80/mirman18b.html>
- [54] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, Omar Fawzi, and Pascal Frossard. 2017. Universal adversarial perturbations. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 1765–1773.
- [55] Mark Niklas Müller, Franziska Eckert, Marc Fischer, and Martin T. Vechev. 2023. Certified Training: Small Boxes are All You Need. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net. <https://openreview.net/pdf?id=7oFuxtJtUMH>
- [56] Satoshi Munakata, Caterina Urban, Haruki Yokoyama, Koji Yamamoto, and Kazuki Munakata. 2023. Verifying Attention Robustness of Deep Neural Networks Against Semantic Perturbations. In *NASA Formal Methods - 15th International Symposium, NFM 2023, Houston, TX, USA, May 16-18, 2023, Proceedings (Lecture Notes in Computer Science, Vol. 13903)*, Kristin Yvonne Rozier and Swarat Chaudhuri (Eds.). Springer, 37–61. https://doi.org/10.1007/978-3-031-33170-1_3

- [57] Alessandro De Palma, Harkirat S. Behl, Rudy Bunel, Philip H. S. Torr, and M. Pawan Kumar. 2021. Scaling the Convex Barrier with Active Sets. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net. <https://openreview.net/forum?id=uQfOy7LrITR>
- [58] Brandon Paulsen, Jingbo Wang, and Chao Wang. 2020. ReluDiff: Differential Verification of Deep Neural Networks. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering (Seoul, South Korea) (ICSE '20)*. Association for Computing Machinery, New York, NY, USA, 714–726. <https://doi.org/10.1145/3377811.3380337>
- [59] Brandon Paulsen, Jingbo Wang, Jiawei Wang, and Chao Wang. 2021. NeuroDiff: Scalable Differential Verification of Neural Networks Using Fine-Grained Approximation. In *Proceedings of the 35th IEEE/ACM International Conference on Automated Software Engineering (Virtual Event, Australia) (ASE '20)*. Association for Computing Machinery, New York, NY, USA, 784–796. <https://doi.org/10.1145/3324884.3416560>
- [60] Yannik Potdevin, Dirk Nowotka, and Vijay Ganesh. 2019. An empirical investigation of randomized defenses against adversarial attacks. *arXiv preprint arXiv:1909.05580* (2019).
- [61] Rob Potharst and Adrianus Johannes Feelders. 2002. Classification trees for problems with monotonicity constraints. *ACM SIGKDD Explorations Newsletter* 4, 1 (2002), 1–10.
- [62] Chongli Qin, Krishnamurthy (Dj) Dvijotham, Brendan O'Donoghue, Rudy Bunel, Robert Stanforth, Sven Gowal, Jonathan Uesato, Grzegorz Swirszcz, and Pushmeet Kohli. 2019. Verification of Non-Linear Specifications for Neural Networks. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=HyeFasRctQ>
- [63] Hadi Salman, Greg Yang, Huan Zhang, Cho-Jui Hsieh, and Pengchuan Zhang. 2019. A Convex Relaxation Barrier to Tight Robustness Verification of Neural Networks. In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*.
- [64] Joseph Scott, Guanting Pan, Elias B. Khalil, and Vijay Ganesh. 2022. Goose: A Meta-Solver for Deep Neural Network Verification. In *Proceedings of the 20th Internal Workshop on Satisfiability Modulo Theories co-located with the 11th International Joint Conference on Automated Reasoning (IJCAR 2022) part of the 8th Federated Logic Conference (FLoC 2022), Haifa, Israel, August 11-12, 2022 (CEUR Workshop Proceedings, Vol. 3185)*, David Déharbe and Antti E. J. Hyvärinen (Eds.). CEUR-WS.org, 99–113. <https://ceur-ws.org/Vol-3185/extended678.pdf>
- [65] Ron Shemer, Arie Gurfinkel, Sharon Shoham, and Yakir Vizel. 2019. Property Directed Self Composition. In *Computer Aided Verification - 31st International Conference, CAV 2019, New York City, NY, USA, July 15-18, 2019, Proceedings, Part I (Lecture Notes in Computer Science, Vol. 11561)*, Isil Dillig and Serdar Tasiran (Eds.). Springer, 161–179. https://doi.org/10.1007/978-3-030-25540-4_9
- [66] Zhouxing Shi, Yihan Wang, Huan Zhang, J Zico Kolter, and Cho-Jui Hsieh. 2022. Efficiently computing local lipschitz constants of neural networks via bound propagation. *Advances in Neural Information Processing Systems* 35 (2022), 2350–2364.
- [67] Gagandeep Singh, Rupanshu Ganvir, Markus Püschel, and Martin Vechev. 2019. Beyond the single neuron convex barrier for neural network certification. In *Advances in Neural Information Processing Systems*.
- [68] Gagandeep Singh, Timon Gehr, Matthew Mirman, Markus Püschel, and Martin Vechev. 2018. Fast and effective robustness certification. *Advances in Neural Information Processing Systems* 31 (2018).
- [69] Gagandeep Singh, Timon Gehr, Markus Püschel, and Martin Vechev. 2019. An abstract domain for certifying neural networks. *Proceedings of the ACM on Programming Languages* 3, POPL (2019).
- [70] Gagandeep Singh, Timon Gehr, Markus Püschel, and Martin Vechev. 2019. Robustness Certification with Refinement. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=HJgeEh09KQ>
- [71] Matthew Sotoudeh, Zhe Tao, and Aditya V Thakur. 2023. SyReNN: a tool for analyzing deep neural networks. *International Journal on Software Tools for Technology Transfer* 25, 2 (2023), 145–165.
- [72] Matthew Sotoudeh and Aditya V Thakur. 2020. Abstract neural networks. In *Static Analysis: 27th International Symposium, SAS 2020, Virtual Event, November 18–20, 2020, Proceedings 27*. Springer, 65–88.
- [73] Marcelo Sousa and Isil Dillig. 2016. Cartesian hoare logic for verifying k-safety properties. In *Proceedings of the 37th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI 2016, Santa Barbara, CA, USA, June 13-17, 2016*, Chandra Krantz and Emery D. Berger (Eds.). ACM, 57–69. <https://doi.org/10.1145/2908080.2908092>
- [74] Cheng Tan, Yibo Zhu, and Chuanxiong Guo. 2021. Building Verified Neural Networks with Specifications for Systems. In *Proceedings of the 12th ACM SIGOPS Asia-Pacific Workshop on Systems (Hong Kong, China) (APSys '21)*. 42–47.
- [75] Hoang-Dung Tran, Diago Manzananas Lopez, Patrick Musau, Xiaodong Yang, Luan Viet Nguyen, Weiming Xiang, and Taylor T. Johnson. 2019. Star-Based Reachability Analysis of Deep Neural Networks. In *Formal Methods – The Next 30 Years*, Maurice H. ter Beek, Annabelle McIver, and José N. Oliveira (Eds.). Springer International Publishing, Cham, 670–686.
- [76] Shubham Ugare, Gagandeep Singh, and Sasa Misailovic. 2022. Proof transfer for fast certification of multiple approximate neural networks. *Proc. ACM Program. Lang.* 6, OOPSLA1 (2022), 1–29. <https://doi.org/10.1145/3527319>
- [77] Hiroshi Unno, Tachio Terauchi, and Eric Koskinen. 2021. Constraint-Based Relational Verification. In *Computer Aided Verification - 33rd International Conference, CAV 2021, Virtual Event, July 20-23, 2021, Proceedings, Part I (Lecture*

- Notes in Computer Science, Vol. 12759), Alexandra Silva and K. Rustan M. Leino (Eds.). Springer, 742–766. https://doi.org/10.1007/978-3-030-81685-8_35
- [78] Shiqi Wang, Kexin Pei, Justin Whitehouse, Junfeng Yang, and Suman Jana. 2018. Efficient formal safety analysis of neural networks. In *Advances in Neural Information Processing Systems*.
- [79] Shiqi Wang, Huan Zhang, Kaidi Xu, Xue Lin, Suman Jana, Cho-Jui Hsieh, and J Zico Kolter. 2021. Beta-CROWN: Efficient Bound Propagation with Per-neuron Split Constraints for Complete and Incomplete Neural Network Verification. *arXiv preprint arXiv:2103.06624* (2021).
- [80] Shiqi Wang, Huan Zhang, Kaidi Xu, Xue Lin, Suman Jana, Cho-Jui Hsieh, and J Zico Kolter. 2021. Beta-CROWN: Efficient Bound Propagation with Per-neuron Split Constraints for Neural Network Robustness Verification. In *Advances in Neural Information Processing Systems*, A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan (Eds.). <https://openreview.net/forum?id=ahYIIRBeCFw>
- [81] Zhilu Wang, Chao Huang, and Qi Zhu. 2022. Efficient global robustness certification of neural networks via interleaving twin-network encoding. In *2022 Design, Automation & Test in Europe Conference & Exhibition (DATE)*. IEEE, 1087–1092.
- [82] Eric Wong and J. Zico Kolter. 2018. Provable Defenses against Adversarial Examples via the Convex Outer Adversarial Polytope. In *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholm, Sweden, July 10-15, 2018 (Proceedings of Machine Learning Research, Vol. 80)*, Jennifer G. Dy and Andreas Krause (Eds.). PMLR, 5283–5292. <http://proceedings.mlr.press/v80/wong18a.html>
- [83] Haoze Wu, Clark Barrett, Mahmood Sharif, Nina Narodytska, and Gagandeep Singh. 2022. Scalable Verification of GNN-Based Job Schedulers. *Proc. ACM Program. Lang.* 6, OOPSLA2, Article 162 (oct 2022), 30 pages. <https://doi.org/10.1145/3563325>
- [84] Haoze Wu, Teruhiro Tagomori, Alexander Robey, Fengjun Yang, Nikolai Matni, George Pappas, Hamed Hassani, Corina Pasareanu, and Clark Barrett. 2023. Toward certified robustness against real-world distribution shifts. In *2023 IEEE Conference on Secure and Trustworthy Machine Learning (SaTML)*. IEEE, 537–553.
- [85] Changming Xu and Gagandeep Singh. 2022. Robust Universal Adversarial Perturbations. *CoRR* abs/2206.10858 (2022). <https://doi.org/10.48550/arXiv.2206.10858> arXiv:2206.10858
- [86] Kaidi Xu, Zhouxing Shi, Huan Zhang, Yihan Wang, Kai-Wei Chang, Minlie Huang, Bhavya Kailkhura, Xue Lin, and Cho-Jui Hsieh. 2020. Automatic Perturbation Analysis for Scalable Certified Robustness and Beyond. (2020).
- [87] Kaidi Xu, Huan Zhang, Shiqi Wang, Yihan Wang, Suman Jana, Xue Lin, and Cho-Jui Hsieh. 2021. Fast and Complete: Enabling Complete Neural Network Verification with Rapid and Massively Parallel Incomplete Verifiers. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=nVZtXBI6LNn>
- [88] Yi Zeng, Zhouxing Shi, Ming Jin, Feiyang Kang, Lingjuan Lyu, Cho-Jui Hsieh, and Ruoxi Jia. 2023. Towards Robustness Certification Against Universal Perturbations. In *The Eleventh International Conference on Learning Representations*. <https://openreview.net/forum?id=7GEvPKxjtt>
- [89] Mustafa Zeqiri, Mark Niklas Müller, Marc Fischer, and Martin T. Vechev. 2023. Efficient Certified Training and Robustness Verification of Neural ODEs. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net. <https://openreview.net/pdf?id=KyoVpYvWWnK>
- [90] Huan Zhang, Hongge Chen, Chaowei Xiao, Sven Gowal, Robert Stanforth, Bo Li, Duane Boning, and Cho-Jui Hsieh. 2020. Towards stable and efficient training of verifiably robust neural networks. In *Proc. International Conference on Learning Representations (ICLR)*.
- [91] Huan Zhang, Shiqi Wang, Kaidi Xu, Linyi Li, Bo Li, Suman Jana, Cho-Jui Hsieh, and J Zico Kolter. 2022. General Cutting Planes for Bound-Propagation-Based Neural Network Verification. In *Advances in Neural Information Processing Systems*, Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (Eds.). <https://openreview.net/forum?id=5haAJAcofjc>
- [92] Huan Zhang, Tsui-Wei Weng, Pin-Yu Chen, Cho-Jui Hsieh, and Luca Daniel. 2018. Efficient neural network robustness certification with general activation functions. *Advances in neural information processing systems* 31 (2018).