

Appendix: DNNV

David Shriver, Matthew B. Dwyer, and Sebastian Elbaum

University of Virginia, Charlottesville, VA, USA
{dls2fc,matthewbdwyer,selbaum}@virginia.edu

A Verification Benchmarks

We examine the benchmarks used to evaluate each of the 13 verifiers supported by DNNV, and determine whether each verifier can run on the benchmark out of the box, and also whether they could be run on the benchmark when DNNV is applied. Here we provide a short description of each of the 19 verification benchmarks that we have identified. A short summary of some of the features of each verifier relevant to DNNV are shown in Table 1. These features include whether any properties cannot represent their input constraints using hyper-rectangles (\neg HR), whether any network in the benchmark contains convolution operations (C), whether any network contains residual structures (R), and whether any network uses any non-ReLU activation functions (\neg ReLU).

The ACAS Xu (AX) benchmark, introduced for *Reluplex* [7], is one of the most used verification benchmarks [2, 4, 8, 16]. The benchmark consists of 10 properties. Property ϕ_1 is a reachability property, specifying an upper bound on one of the 5 output variables. Properties ϕ_5 , ϕ_6 , ϕ_9 , and ϕ_{10} are all traditional class robustness properties, specifying the desired class for the given input region. Properties ϕ_3 , ϕ_4 , ϕ_7 and ϕ_8 are reachability properties, specifying a set of acceptable classes for the input region. Property ϕ_2 is also a reachability property, specifying that a given output value cannot be greater than all others. Each of the properties are applied to a subset of 45 networks trained on an aircraft collision avoidance dataset, with 5 inputs, 5 output classes and 6 layers of 50 neurons each. The original benchmark included networks in *Reluplex*-NNET format, and a custom version of *Reluplex* was written for each property. Later uses of the benchmark translated the verification problems into RLV format, which is used by *Planet*, *BaB*, and *BaBSB*, as well as translating the networks into ONNX. The benchmark in ONNX and DNNP format is fully supported by DNNV.

The Collision Detection (CD) benchmark [5], introduced for the evaluation of *Planet*, consists of 500 local robustness properties for an 80 neuron network with a fully connected layer and max pooling layer that classifies whether 2 simulated vehicles will collide, given their current state. The verification problems, in RLV format, are supported by *Planet*, *BaB*, and *BaBSB*. The problems have also been modified to convert max pooling operations to a sequence of fully-connected layers with ReLU activations, and then translated to *Reluplex*-NNET format, enabling off the shelf support by *Marabou*, and a generalized version of *Reluplex*. This benchmark is one of the few that is not supported by DNNV,

Table 1. Verifier benchmarks.

Key	Name	Uses	#P	# \mathcal{N}	Features		
					\neg HR	C	R \neg ReLU
AX	ACAS Xu	[2, 4, 7, 8, 16]	10	45			
CD	Collision Detection	[4, 5, 8]	500	1			
PM	<i>Planet</i> MNIST	[5]	7	1	✓	✓	
TS	TwinStream	[3]	1	81			
PCA	PCAMNIST	[4]	12	17			
MM	<i>MIPVerify</i> MNIST	[15]	10000	5		✓	
MC	<i>MIPVerify</i> CIFAR10	[15]	10000	2		✓	✓
NM	<i>Neurify</i> MNIST	[6, 16]	500	4		✓	
NDB	<i>Neurify</i> Drebin	[16]	500	3			
NDv	<i>Neurify</i> DAVE	[16]	200	1	✓	✓	
DZM	<i>DeepZono</i> MNIST	[12]	1700	10		✓	✓
DZC	<i>DeepZono</i> CIFAR10	[12]	1700	5		✓	✓
DPM	<i>DeepPoly</i> MNIST	[6, 13]	1500	8		✓	✓
DPC	<i>DeepPoly</i> CIFAR10	[13]	800	5		✓	
RZM	<i>RefineZono</i> MNIST	[14]	800	8		✓	
RZC	<i>RefineZono</i> CIFAR10	[14]	200	2		✓	
RPM	<i>RefinePoly</i> MNIST	[11]	600	6		✓	
RPC	<i>RefinePoly</i> CIFAR10	[11]	300	3		✓	✓
VC	<i>VeriNet</i> CIFAR10	[6]	250	1		✓	

since the network contains structures that are not easily supported by ONNX. In particular, the max-pooling operation in the original network, applied to a flat tensor, cannot be encoded by ONNX from their original format.

The *Planet* MNIST (PM) benchmark [5] is a set of 7 properties over a convolutional network trained on the MNIST dataset [10]. The first 4 of these are reachability properties with hyper-rectangle input constraints, while the next 2 are local robustness properties with hyper-rectangle input constraints, and the final property is an local robustness property with halfspace-polytope input constraints. The original benchmark was provided in RLV format. The first 6 of these properties are currently supported by DNNV, while the final property could be supported by DNNV with additional engineering effort.

The TwinStream (TS) benchmark [3] consists of 1 property applied to 81 networks that output a constant value. The property asserts that for all inputs, the output of the network is positive. The original benchmark was provided in RLV format. This benchmark is fully supported by DNNV for all verifiers.

The PCAMNIST (PCA) benchmark [4] consists of 12 reachability properties applied to 17 networks trained on modified versions of the MNIST dataset to predict the parity of the digit represented by the first k components of the PCA decomposition of an image. The original benchmark was provided in RLV format. This benchmark is fully supported by DNNV for all verifiers.

MIPVerify MNIST (MM) consists of 10000 local robustness properties applied to 5 networks trained on the MNIST dataset. The networks have varied

structures: 2 networks are fully connected and 3 are convolutional. We could not find an available version of the benchmark used by *MIPVerify* to evaluate its original input format. This benchmark is fully supported by DNNV for all verifiers except *Reluplex*, which does not support convolution operations.

MIPVerify CIFAR (MC) consists of 10000 local robustness properties applied to 2 networks trained on the CIFAR10 dataset [9]. One of these networks is a convolutional network and the other is a residual network. We could not find an available version of the benchmark used by *MIPVerify* to evaluate its original input format. This benchmark is supported by DNNV for verifiers that can support residual connections, including: *Planet*, *DeepZono*, *DeepPoly*, *RefineZono*, and *RefinePoly*. While the benchmark is supported by the version of *MIPVerify* used in its study, it is not supported through DNNV, since the publicly available version of *MIPVerify* does not support residual connections.

The *Neurify* MNIST (NM) benchmark [16] consists of 500 L_∞ local robustness properties across 4 MNIST networks, 3 fully connected networks with 58, 110, and 1034 neurons respectively, and a convolutional network with 4814 neurons. The original benchmark was provided in *Neurify*-NNET format, with properties hard-coded into the verifier. DNNV enables almost all verifiers to run on this benchmark. *Reluplex* cannot be run due to the presence of convolutional layers, which are not supported. *MIPVerify* cannot be run due to the presence of non-hypercube input constraints. While this limitation of the verifier can be satisfied with a property reduction for fully-connected networks, DNNV does not currently support such a reduction for convolutional networks.

The *Neurify* Drebin (NDB) benchmark [16] consists of 500 L_∞ local robustness properties across 3 fully connected Drebin [1] networks with 102, 212, and 402 neurons each. The original benchmark was provided in *Neurify*-NNET format, with properties hard-coded into the verifier. This benchmark is fully supported by DNNV for all verifiers.

The *Neurify* DAVE (NDV) benchmark [16] consists of 200 local reachability properties, with 4 different types of input constraints (50 properties of each type). The first type of input constraint is an L_∞ constraint, which is equivalent to a hyper-rectangle constraint. The second type of input constraint is an L_1 constraint, which can be written as a halfspace polytope constraint. The third and fourth type of input constraint are image brightening and contrast, which can be written as halfspace polytope constraints. The properties are applied to a convolutional network for an autonomous vehicle, with 10276 neurons. The original benchmark was provided in *Neurify*-NNET format, with properties hard-coded into the verifier. Similar to the *Neurify* MNIST benchmark, DNNV enables almost all verifiers to run on this benchmark. *Reluplex* cannot be run, due to the presence of convolutional layers, which are not supported, and *MIPVerify* cannot be run due to the presence of non-hypercube input constraints.

The *DeepZono* MNIST (DZM) benchmark [12] consists of 1700 local robustness properties, subsets of which are applied to 10 networks trained on the MNIST dataset. The networks have varied structures and activation functions: 3 networks are fully connected, 1 of which uses ReLU activations, 1 with Tanh

activations, and 1 with Sigmoid activations; 6 are convolutional, 4 of which have ReLU activations, 1 with Tanh activations, and 1 with Sigmoid activations; and 1 is a residual network. The networks in the original benchmark were provided in a custom human-readable text format, with properties hard-coded into the verifier. DNNV does not increase the support for this benchmark due to the presence of both a residual network and non-ReLU activation functions.

The *DeepZono* CIFAR10 (DZC) benchmark [12] consists of 1700 local robustness properties, subsets of which are applied to 5 networks trained on the CIFAR10 dataset. The networks have varied structures and activation functions: 3 networks are fully connected, 1 of which uses ReLU activations, 1 with Tanh activations, and 1 with Sigmoid activations; and 2 are convolutional with ReLU activations. The networks in the original benchmark were provided in a custom human-readable text format, with properties hard-coded into the verifier. DNNV enables *VeriNet* to run on this benchmark. Other verifiers are not supported due to the non-ReLU activation functions.

The *DeepPoly* MNIST (DPM) benchmark [13] consists of 1500 local robustness properties, subsets of which are applied to 8 networks trained on the MNIST dataset. The networks have varied structures and activation functions: 5 networks are fully connected, 3 of which uses ReLU activations, 1 with Tanh activations, and 1 with Sigmoid activations; and 3 are convolutional with ReLU activations. The networks in the original benchmark were provided in a custom human-readable text format, with properties hard-coded into the verifier. DNNV enables *VeriNet* to run on this benchmark. Other verifiers are not supported due to the non-ReLU activation functions.

The *DeepPoly* CIFAR10 (DPC) benchmark [13] consists of 800 local robustness properties, subsets of which are applied to 5 networks trained on the CIFAR10 dataset. The networks have varied structures: 3 networks are fully connected with ReLU activations; and 2 are convolutional with ReLU activations. The networks in the original benchmark were provided in a custom human-readable text format, with properties hard-coded into the verifier. DNNV enables several additional verifiers to support this benchmark. In particular, it enables most verifiers that can be applied to convolutional networks with relu activations.

The *RefineZono* MNIST (RZM) benchmark [14] consists of 800 local robustness properties, subsets of which are applied to 8 networks trained on the MNIST dataset. 5 networks are fully connected with ReLU activations and 3 are convolutional with ReLU activations. The networks in the original benchmark were provided in a custom human-readable text format, with properties hard-coded into the verifier. DNNV enables several additional verifiers to support this benchmark. In particular, it enables most verifiers that can be applied to convolutional networks with relu activations.

The *RefineZono* CIFAR10 (RZC) benchmark [14] consists of 200 local robustness properties, subsets of which are applied to 2 networks trained on the CIFAR10 dataset. One of the networks is fully connected with ReLU activations and the other is convolutional with ReLU activations. The networks in the orig-

inal benchmark were provided in a custom human-readable text format, with properties hard-coded into the verifier. DNNV enables several additional verifiers to support this benchmark. In particular, it enables most verifiers that can be applied to convolutional networks with relu activations.

The *RefinePoly* MNIST (RPM) benchmark [11] consists of 600 local robustness properties, subsets of which are applied to 6 networks trained on the MNIST dataset. 4 networks are fully connected with ReLU activations and 2 are convolutional with ReLU activations. The networks in the original benchmark were provided in a custom human-readable text format, with properties hard-coded into the verifier. DNNV enables several additional verifiers to support this benchmark. In particular, it enables most verifiers that can be applied to convolutional networks with relu activations.

The *RefinePoly* CIFAR10 (RPC) benchmark [11] consists of 300 local robustness properties, subsets of which are applied to 3 networks trained on the MNIST dataset. Two of the networks are convolutional with ReLU activations and the third is a residual network with ReLU activations. The networks in the original benchmark were provided in a custom human-readable text format, with properties hard-coded into the verifier. DNNV enables the *Planet* verifier to support this benchmark. In particular, it enables most verifiers that can be applied to convolutional networks with relu activations. Other verifiers do not support the residual structure of one of the networks.

The *VeriNet* CIFAR10 (VC) benchmark [6] consists of 250 local robustness properties applied to 1 convolutional network with ReLU activations. The networks were provided in ONNX format, with hard-coded properties. DNNV enables support of this benchmark by most of the integrated verifiers. *Reluplex* does not support convolutional networks, and *MIPVerify* does not support properties with input constraints that are not hyper-cubes.

A.1 Support

We summarize the support of each verifier for each of the benchmarks in Table 2. Each row of this table corresponds to one of the 13 verifiers supported by DNNV, and each column corresponds to one of the 19 benchmarks identified in Table 1. Each cell of the table may contain a circle that identifies the support of the verifier for the benchmark. The left half of the circle is filled black if the verifier can support the benchmark out of the box, and is filled white otherwise. The right half is filled black if the verifier supports the benchmark through DNNV, filled gray if support is planned, and is filled white otherwise. Planned support means that DNNV will support the benchmark after the implementation for halfspace polytope constraints in the input space is completed. We plan to implement support for this feature by the notification deadline on December 23, 2020. An absent circle indicates that the verifier can not be made to support some aspect of the benchmark. For the benchmarks shown here, this is always due to the presence of non-ReLU activation functions in some of the networks in the benchmarks.

Table 2. Benchmark support by each verifier. The left half of the circle is black if the verifier can support the benchmark out of the box, and is white otherwise. The right half is black if the verifier supports the benchmark through DNNV, gray if support is pending DNNV implementation, and is white otherwise. An absent circle indicates that the verifier can not be made to support some aspect of the benchmark.

Verifier	Benchmark																		
	AX	CD	PM	TS	PCA	MM	MC	NM	NDB	NDV	DZM	DZC	DPM	DPC	RZM	RZC	RPM	RPC	VC
Reluplex	●	◐	○	●	●	○	○	○	○	◐	○			○	○	○	○	○	○
Planet	●	◐	◐	●	●	◐	◐	◐	◐	◐				○	◐	◐	◐	◐	◐
BaB	●	◐	◐	●	●	○	○	◐	◐	◐				○	◐	◐	◐	◐	◐
BaBSB	●	◐	◐	●	●	○	○	◐	◐	◐				○	◐	◐	◐	◐	◐
MIPVerify	○	○	○	◐	◐	◐	◐	◐	◐	○				○	○	○	○	○	○
Neurify	●	○	◐	◐	◐	○	○	●	●	◐				○	◐	◐	◐	◐	◐
DeepZono	●	○	◐	◐	◐	◐	◐	◐	◐	◐	●	●	●	●	●	●	●	●	◐
DeepPoly	●	○	◐	◐	◐	◐	◐	◐	◐	◐	●	●	●	●	●	●	●	●	◐
RefineZono	●	○	◐	◐	◐	◐	◐	◐	◐	◐	●	●	●	●	●	●	●	●	◐
RefinePoly	●	○	◐	◐	◐	◐	◐	◐	◐	◐	●	●	●	●	●	●	●	●	◐
Marabou	●	◐	◐	●	●	○	○	◐	◐	◐				○	◐	◐	◐	○	◐
nenum	●	○	◐	◐	◐	○	○	◐	◐	◐				○	◐	◐	◐	○	◐
VeriNet	◐	○	◐	◐	◐	○	○	●	◐	◐	○	◐	◐	◐	◐	◐	◐	○	●

B Property Reduction

In this section, we provide the algorithm for reducing properties to reachability properties, as well as proofs for the equivalidity of the resulting set of reachability properties and original property. Algorithm 1 is the overall reduction algorithm, while Algorithm 2 and 3 are subprocedures used by the main algorithm. The algorithm and proofs for reduction to other property types (such as robustness) are very similar.

We assume that properties are of the form $\forall x \in \mathbb{R}^n : \phi_{\mathcal{X}}(x) \rightarrow \phi_{\mathcal{Y}}(\mathcal{N}(x))$, where $\phi_{\mathcal{X}}$ is a set of constraints over the inputs – the pre-condition, and $\phi_{\mathcal{Y}}$ is a set of constraints over the outputs – the post-condition. We also assume that constraints are represented as linear inequalities.

B.1 Proofs

In order to prove that the property reduction produces a set of correctness problems equivalid to the original problem, we first prove the following lemmas:

Lemma 1. *Let ϕ be a conjunction of linear inequalities over the variables x_i for i from 0 to $n - 1$. We can construct a halfspace polytope $H = (A, b)$ with Algorithm 2 such that $(Ax \leq b) \Leftrightarrow (x \models \phi)$.*

Algorithm 1: Property Reduction

Input: Correctness problem $\langle \mathcal{N}, \phi \rangle$
Output: A set of robustness problems $\{\langle \mathcal{N}_1, \phi_1 \rangle, \dots, \langle \mathcal{N}_i, \phi_i \rangle\}$

```

1 begin
2    $\phi' \leftarrow DNF(\neg\phi)$ 
3    $\Psi \leftarrow \emptyset$ 
4   for  $disjunct \in \phi'$  do
5      $\phi_x \leftarrow \text{extract\_input\_constraints}(disjunct)$ 
6      $\phi_y \leftarrow \text{extract\_output\_constraints}(disjunct)$ 
7      $hspoly \leftarrow \text{disjunct\_to\_hpolytope}(\phi_y)$ 
8      $suffix \leftarrow \text{construct\_suffix}(hspoly)$ 
9      $\mathcal{N}' \leftarrow suffix \circ \mathcal{N}$ 
10     $\phi' \leftarrow \forall x. (x \in \phi_x \implies \mathcal{N}'(x)_0 > \mathcal{N}'(x)_1)$ 
11     $\Psi \leftarrow \Psi \cup \langle \mathcal{N}', \phi' \rangle$ 
12  return  $\Psi$ 

```

Proof. We first show that every linear inequality in the conjunction can be reformulated to the form $a_0x_0 + a_1x_1 + \dots + a_{n-1}x_{n-1} \leq b$. It is trivial to show that inequalities with a \geq comparison can be manipulated to an equivalent form with \leq , and $>$ can be manipulated to become $<$. It is also trivial to show that the inequality can be manipulated to have variables on lhs and a constant value on rhs. This results in a conjunction of linear inequalities of the form $a_0x_0 + a_1x_1 + \dots + a_{n-1}x_{n-1} < b$ and $a_0x_0 + a_1x_1 + \dots + a_{n-1}x_{n-1} \leq b$. Finally, the $<$ comparison can be changed to a \leq comparison by decrementing the constant on the right-hand-side from b to b' where b' is the largest representable number less than b .

We prove that linear inequalities using the $<$ comparison can be reformulated to use a \leq comparison using a proof by contradiction. Assume that either $a_0x_0 + a_1x_1 + \dots + a_{n-1}x_{n-1} < b$ and $a_0x_0 + a_1x_1 + \dots + a_{n-1}x_{n-1} > b'$ or $a_0x_0 + a_1x_1 + \dots + a_{n-1}x_{n-1} \geq b$ and $a_0x_0 + a_1x_1 + \dots + a_{n-1}x_{n-1} \leq b'$. Then one of two cases must be true. Either $b' < a_0x_0 + a_1x_1 + \dots + a_{n-1}x_{n-1} < b$, a contradiction, since $a_0x_0 + a_1x_1 + \dots + a_{n-1}x_{n-1}$ cannot be both larger than the largest representable number less than b and also less than b .¹ Or $b \leq a_0x_0 + a_1x_1 + \dots + a_{n-1}x_{n-1} \leq b'$, a contradiction, since $b' < b$ by definition.

Given a conjunction of linear inequalities in the form $a_0x_0 + a_1x_1 + \dots + a_{n-1}x_{n-1} \leq b$, Algorithm 2 constructs A and b with a row in A and value in b corresponding to each conjunct. There are two cases to prove: $(Ax \leq b) \rightarrow (x \models \phi)$ and $(x \models \phi) \rightarrow (Ax \leq b)$.

We prove case 1 by contradiction. Assume $(Ax \leq b)$ and $(x \not\models \phi)$. By the construction of H in Algorithm 2, each conjunct of ϕ is exactly 1 constraint in H . If $Ax \leq b$, then all constraints in H must be satisfied, and thus all conjuncts in ϕ must be satisfied and $x \models \phi$, a contradiction.

¹ We further discuss the assumption that such a number exists in Section B.2.

Algorithm 2: disjunct_to_hpolytope

Input: Conjunction of linear inequalities ϕ_i
Output: Halfspace polytope H

```

1 begin
2    $H \leftarrow (A, b)$  where  $A$  is an  $(|\phi_i|) \times (m)$  matrix where columns correspond
   to the output variables  $N(x)_0$  to  $N(x)_{m-1}$ 
3   for  $ineq_j \in \phi_i$  do
4     if  $ineq_j$  uses  $\geq$  then
5       swap lhs and rhs; switch inequality to  $\leq$ 
6     else if  $ineq_j$  uses  $>$  then
7       swap lhs and rhs; switch inequality to  $<$ 
8       move variables to lhs; move constants to rhs
9     if  $ineq_j$  uses  $<$  then
10      decrement rhs; switch inequality to  $\leq$ 
11       $A_j \leftarrow$  coefficients of variables on lhs
12       $b_j \leftarrow$  rhs constant
13  return  $H$ 

```

Algorithm 3: construct_suffix

Input: Halfspace polytope $H = (A, b)$
Output: A DNN with 2 fully connected layers S

```

1 begin
2    $S_h \leftarrow \text{ReLU}(\text{FullyConnectedLayer}(A, -b))$ 
3    $W \leftarrow \begin{bmatrix} 1 & 1 & \dots & 1 \\ 0 & 0 & \dots & 0 \end{bmatrix}$ 
4    $S_o \leftarrow \text{FullyConnectedLayer}(W, \vec{0})$ 
5    $S \leftarrow S_o \circ S_h$ 
6  return  $S$ 

```

We prove case 2 by contradiction. Assume $(x \models \phi)$ and $(Ax \not\leq b)$. By the construction of H in Algorithm 2, each conjunct of ϕ is exactly 1 constraint in H . If $x \models \phi$, then all conjuncts in ϕ must be satisfied, and thus all constraints in H must be satisfied and $Ax \leq b$, a contradiction.

Lemma 2. *Let $H = (A, b)$ be a halfspace polytope such that $Ax \leq b$. Then, a DNN, \mathcal{N}_s , can be built with Algorithm 3 that classifies whether its outputs satisfy $A(\mathcal{N}(x)) \leq b$ or not. Formally, $\mathcal{N}(x) \in H \Leftrightarrow \mathcal{N}_s(x)_0 \leq \mathcal{N}_s(x)_1$.*

Proof. There are 2 cases:

1. $\mathcal{N}(x) \in H \rightarrow \mathcal{N}_s(x)_0 \leq \mathcal{N}_s(x)_1$
2. $\mathcal{N}_s(x)_0 \leq \mathcal{N}_s(x)_1 \rightarrow \mathcal{N}(x) \in H$

We prove case 1 by contradiction. Assume $\mathcal{N}(x) \in H$ and $\mathcal{N}_s(x)_0 > \mathcal{N}_s(x)_1$. From Algorithm 3, each neuron in the hidden layer of \mathcal{N}_s corresponds to one

constraint in H . The weights of each neuron are the values in the corresponding row of A , and the bias is the negation of the corresponding value of b . If the output $\mathcal{N}(x)$ satisfies the constraint, then the value of the neuron will be less than or equal to 0, otherwise it will be greater than 0. After application of the ReLU activation function, all neurons will be equal to 0 if their corresponding constraint is satisfied by $\mathcal{N}(x)$ and greater than 0 otherwise. The first neuron in the final layer sums all of the neurons in the hidden layer, while the second neuron has a constant value of 0. If $\mathcal{N}(x) \in H$, then all neurons in the hidden layer after activation must have a value of 0 since all constraints are satisfied. However, if all neurons have a value of 0, then their sum must also have a value of zero, and therefore $\mathcal{N}_s(x)_0 = \mathcal{N}_s(x)_1$, a contradiction.

We prove case 2 by contradiction. Assume $\mathcal{N}_s(x)_0 \leq \mathcal{N}_s(x)_1$ and $\mathcal{N}(x) \notin H$. From Algorithm 3, each neuron in the hidden layer of \mathcal{N}_s corresponds to one constraint in H . The weights of each neuron are the values in the corresponding row of A , and the bias is the negation of the corresponding value of b . If the output $\mathcal{N}(x)$ satisfies the constraint, then the value of the neuron will be less than or equal to 0, otherwise it will be greater than 0. After application of the ReLU activation function, all neurons will be equal to 0 if their corresponding constraint is satisfied by $\mathcal{N}(x)$ and greater than 0 otherwise. The first neuron in the final layer sums all of the neurons in the hidden layer, while the second neuron has a constant value of 0. If $\mathcal{N}(x) \notin H$, then at least one neurons in the hidden layer after activation must have a value greater than 0 since at least one constraint is not satisfied. However, if any neuron has a value greater than 0, then their sum must also have a value greater than zero, and therefore $\mathcal{N}_s(x)_0 > \mathcal{N}_s(x)_1$, a contradiction.

Theorem 1. *Let $\psi = \langle \mathcal{N}, \phi \rangle$ be an arbitrary correctness problem with a DNN correctness property defined as a formula of disjunctions and conjunctions of linear inequalities over the input and output variables of \mathcal{N} . Property Reduction (Algorithm 1) maps ψ to an equivalent set of correctness problems $reduce(\psi) = \{\langle \mathcal{N}_1, \phi_1 \rangle, \dots, \langle \mathcal{N}_k, \phi_k \rangle\}$.*

$$\mathcal{N} \models \psi \Leftrightarrow \forall \langle \mathcal{N}_i, \phi_i \rangle \in reduce(\psi). \mathcal{N}_i \models \phi_i$$

Proof. A model that satisfies any disjunct of $DNF(\neg\phi)$ falsifies ϕ . If ϕ is falsifiable, then at least one disjunct of $DNF(\neg\phi)$ is satisfiable.

Algorithm 1 constructs a correctness problem for each disjunct. For each disjunct, Algorithm 1 constructs a halfspace polytope, H , which is used to construct a suffix network, \mathcal{N}_s . The algorithm then constructs the network $\mathcal{N}'(x) = \mathcal{N}_s(\mathcal{N}(x))$. Algorithm 1 pairs each constructed network with the property $\phi = \forall x. x \in [0, 1]^n \rightarrow \mathcal{N}'(x)_0 > \mathcal{N}'(x)_1$. A violation occurs only when $\mathcal{N}'(x)_0 \leq \mathcal{N}'(x)_1$. By Lemmas 1 and 2, we get that $\mathcal{N}'(x)_0 \leq \mathcal{N}'(x)_1$ if and only if $\mathcal{N}'(x) \in H$. If $\mathcal{N}'(x) \in H$ then $\mathcal{N}'(x)$ satisfies the disjunct and is therefore a violation of the original property.

B.2 On the Existence of a Bounded Largest Representable Number

Our proof that property reduction generates a set of robustness problems equivalent to an arbitrary problem relies on the assumption that strict inequalities can be converted to non-strict inequalities. To do so we rely on the existence of a largest representable number that is less than some given value. While this is not necessarily true for all sets of numbers (e.g., \mathbb{R}), it is true for most numeric representations used in computation (e.g., IEEE 754 floating point arithmetic).

C DNN Simplifications

In this section, we describe the DNN simplifications currently performed by DNNV. This is not a full list of all possible simplifications, but have been useful for some networks we have encountered in practice.

C.1 BatchNormalization Simplification

BatchNormalization simplification removes BatchNormalization operations from a network by combining them with a preceding Conv operation or Gemm operation. If no applicable preceding layer exists, the batch normalization layer is converted into an equivalent Conv operation. This simplification can decrease the number of operations in the model and increase verifier support, since many verifiers do not support BatchNormalization operations.

C.2 Identity Removal

DNNV removes many types of identity operations from DNN models, including explicit Identity operations, Concat operations with a single input, and Flatten operations applied to flat tensors. Such operations can occur in DNN models due to user error, or through automated processes, and their removal does not affect model behavior.

C.3 Convert MatMul followed by Add to Gemm

DNNV converts instances of MatMul (matrix multiplication) operations, followed immediately by Add operations to an equivalent Gemm (generalized matrix multiplication) operation. The Gemm operation generalizes the matrix multiplication and addition, and can simplify subsequent processing and analysis of the DNN.

C.4 Combine Consecutive Gemm

DNNV combines two consecutive Gemm operations into a single equivalent Gemm operation, reducing the number of operations in the DNN.

C.5 Combine Consecutive Conv

In special cases, DNNV can combine consecutive Conv (convolution) operations into a single equivalent Conv operation, reducing the number of operations in the DNN. Currently, DNNV can combine Conv operations when the first Conv uses a diagonal 1 by 1 kernel with a stride of 1 and no zero padding, and the second Conv has no zero padding. This case can occur after converting a normalization layer (such as BatchNormalization) to a Conv operation.

C.6 Bundle Pad

DNNV can bundle explicit Pad operations with an immediately succeeding Conv or MaxPool operation. This both simplifies the DNN model, and increases support, since many verifiers do not support explicit Pad operations (but can support padding as part of a Conv or MaxPool operation).

C.7 Move Activations Backward

DNNV moves activation functions through reshaping operations to immediately succeed the most recent non-reshaping operation. This is possible since activation functions are element-wise operations. This transformation can simplify pattern matching in later analysis steps by reducing the number of possible patterns.

References

1. Arp, D., Spreitzenbarth, M., Hubner, M., Gascon, H., Rieck, K.: DREBIN: effective and explainable detection of android malware in your pocket. In: 21st Annual Network and Distributed System Security Symposium, NDSS 2014, San Diego, California, USA, February 23-26, 2014. The Internet Society (2014), <https://www.ndss-symposium.org/ndss2014/drebin-effective-and-explainable-detection-android-malware-your-pocket>
2. Bak, S., Tran, H.D., Hobbs, K., Johnson, T.T.: Improved geometric path enumeration for verifying relu neural networks. In: Lahiri, S.K., Wang, C. (eds.) Computer Aided Verification. pp. 66–96. Springer International Publishing, Cham (2020)
3. Bunel, R., Turkaslan, I., Torr, P.H.S., Kohli, P., Kumar, M.P.: Piecewise linear neural network verification: A comparative study. CoRR **abs/1711.00455v1** (2017), <http://arxiv.org/abs/1711.00455v1>
4. Bunel, R.R., Turkaslan, I., Torr, P.H.S., Kohli, P., Mudigonda, P.K.: A unified view of piecewise linear neural network verification. In: NeurIPS. pp. 4795–4804 (2018)
5. Ehlers, R.: Formal verification of piece-wise linear feed-forward neural networks. In: Automated Technology for Verification and Analysis - 15th International Symposium, ATVA 2017, Pune, India, October 3-6, 2017, Proceedings. pp. 269–286 (2017). https://doi.org/10.1007/978-3-319-68167-2_19, https://doi.org/10.1007/978-3-319-68167-2_19

6. Henriksen, P., Lomuscio, A.: Efficient neural network verification via adaptive refinement and adversarial search. In: Giacomo, G.D., Catalá, A., Dilkina, B., Milano, M., Barro, S., Bugarín, A., Lang, J. (eds.) ECAI 2020 - 24th European Conference on Artificial Intelligence, 29 August-8 September 2020, Santiago de Compostela, Spain, August 29 - September 8, 2020 - Including 10th Conference on Prestigious Applications of Artificial Intelligence (PAIS 2020). *Frontiers in Artificial Intelligence and Applications*, vol. 325, pp. 2513–2520. IOS Press (2020). <https://doi.org/10.3233/FAIA200385>, <https://doi.org/10.3233/FAIA200385>
7. Katz, G., Barrett, C.W., Dill, D.L., Julian, K., Kochenderfer, M.J.: Reluplex: An efficient SMT solver for verifying deep neural networks. In: *Computer Aided Verification - 29th International Conference, CAV 2017, Heidelberg, Germany, July 24-28, 2017, Proceedings, Part I*. pp. 97–117 (2017). https://doi.org/10.1007/978-3-319-63387-9_5, https://doi.org/10.1007/978-3-319-63387-9_5
8. Katz, G., Huang, D.A., Ibeling, D., Julian, K., Lazarus, C., Lim, R., Shah, P., Thakoor, S., Wu, H., Zeljić, A., et al.: The Marabou framework for verification and analysis of deep neural networks. In: *International Conference on Computer Aided Verification*. pp. 443–452. Springer (2019)
9. Krizhevsky, A., Hinton, G.: Learning multiple layers of features from tiny images (2009)
10. LeCun, Y., Cortes, C., Burges, C.J.: The mnist database of handwritten digits
11. Singh, G., Ganvir, R., Püschel, M., Vechev, M.T.: Beyond the single neuron convex barrier for neural network certification. In: Wallach, H.M., Larochelle, H., Beygelzimer, A., d’Alché-Buc, F., Fox, E.B., Garnett, R. (eds.) *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, 8-14 December 2019, Vancouver, BC, Canada*. pp. 15072–15083 (2019), <http://papers.nips.cc/paper/9646-beyond-the-single-neuron-convex-barrier-for-neural-network-certification>
12. Singh, G., Gehr, T., Mirman, M., Püschel, M., Vechev, M.: Fast and effective robustness certification. In: Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., Garnett, R. (eds.) *Advances in Neural Information Processing Systems 31*, pp. 10802–10813. Curran Associates, Inc. (2018), <http://papers.nips.cc/paper/8278-fast-and-effective-robustness-certification.pdf>
13. Singh, G., Gehr, T., Püschel, M., Vechev, M.T.: An abstract domain for certifying neural networks. *PACMPL* **3**(POPL), 41:1–41:30 (2019)
14. Singh, G., Gehr, T., Püschel, M., Vechev, M.T.: Boosting robustness certification of neural networks. In: *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net (2019), <https://openreview.net/forum?id=HJgeEh09KQ>
15. Tjeng, V., Xiao, K.Y., Tedrake, R.: Evaluating robustness of neural networks with mixed integer programming. In: *International Conference on Learning Representations* (2019), <https://openreview.net/forum?id=HyGIIdiRqtm>
16. Wang, S., Pei, K., Whitehouse, J., Yang, J., Jana, S.: Efficient formal safety analysis of neural networks. In: *NeurIPS*. pp. 6369–6379 (2018)