# Generating and Checking DNN Verification Proofs

Hai Duong     ThanhVu Nguyen     Matthew Dwyer
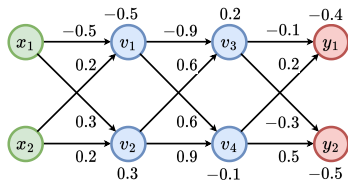
**GEORGE MASON UNIVERSITY**

**UNIVERSITY of VIRGINIA**

# DNN Problems

# DNN Verification

## DNN Verification Example



$$(x_1, x_2) \in [-2.0, 2.0] \times [-1.0, 1.0] \Rightarrow (y_1 > y_2)$$

DNN verifiers: return unsat for valid property and sat otherwise

# DNN Verification

## DNN Verification Example
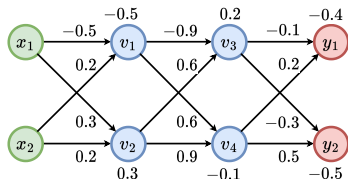


$(x_1, x_2) \in [-2.0, 2.0] \times [-1.0, 1.0] \Rightarrow (y_1 > y_2)$

DNN verifiers: return unsat for valid property and sat otherwise

## But, Could We Trust Verifiers?

- Currently, sat comes with counterexample; but unsat has no evidence!
- Bugs are inevitable (verification tools 20K+ lines of code)
- Multiple reports of soundness bugs in verifiers and literature
  - Claiming something is safe (unsat) when it is not

# Needs Proof of "Proved" Results!

# Challenges In Proof Generation and Checking

**Challenge 1: Compatibility**

Different verifiers use different algorithms and optimizations

Need a proof generation approach compatible with many verifiers

# Challenges In Proof Generation and Checking

**Challenge 1: Compatibility**

Different verifiers use different algorithms and optimizations
Need a proof generation approach compatible with many verifiers

**Challenge 2: Human-Readability**

Proofs needs to be represented and stored efficiently
Need a standard, compact, and human-readable format

# Challenges In Proof Generation and Checking

**Challenge 1: Compatibility**

Different verifiers use different algorithms and optimizations
Need a proof generation approach compatible with many verifiers

**Challenge 2: Human-Readability**

Proofs needs to be represented and stored efficiently
Need a standard, compact, and human-readable format

**Challenge 3: Efficiency**

Proofs can be *very* large!!
Need an efficient proof checking algorithm

# Challenge 1: Compatibility

**Initialization**
- Input: a DNN $N$ and a property $\phi$
- An `ActPatterns` $\leftarrow \{\emptyset\}$ to record branches (patterns)

**Branch-and-Bound Loop**
- **Select** (pop) a branch $\sigma$ from `ActPatterns`.
- **Deduce (bound)** checks feasibility of $\sigma$ wrt $(N, \phi)$
  - If *infeasible*:
    - Prune this $\sigma$ branch (verified)

  - If *feasible*:
    - **Decide (branch)** selects a neuron $v$
    - Splitting $\sigma \wedge (v = \text{on})$ and $\sigma \wedge (v = \text{off})$
    - Add new branches to `ActPatterns`

**Termination**
- If a counterexample is found during search, return (**sat,** cex)
- If no counterexample exists, return (**unsat,** )

# Challenge 1: Compatibility

**Initialization**
- Input: a DNN $N$ and a property $\phi$
- An ActPatterns $\leftarrow \{\emptyset\}$ to record branches (patterns)
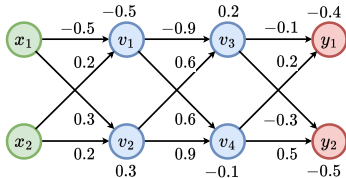- A proof tree prooftree $\leftarrow \{\ \}$ to record conflicts

**Branch-and-Bound Loop**
- **Select** (pop) a branch $\sigma$ from ActPatterns.
- **Deduce (bound)** checks feasibility of $\sigma$ wrt $(N, \phi)$
    - If *infeasible*:
        - Prune this $\sigma$ branch (verified)
        - Record $\sigma$ to prooftree $\rightarrow$ prooftree $\cup \{\sigma\}$
    - If feasible:
        - **Decide (branch)** selects a neuron $v$
        - Splitting $\sigma \wedge (v = \text{on})$ and $\sigma \wedge (v = \text{off})$
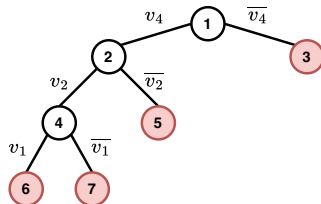        - Add new branches to ActPatterns

**Termination**
- If a counterexample is found during search, return (**sat**, cex)
- If no counterexample exists, return (**unsat**, prooftree)

## Example: BaB + Proof Generation



A simple DNN

A proof tree

Verifying $(x_1, x_2) \in [-2.0, 2.0] \times [-1.0, 1.0] \Rightarrow (y_1 > y_2)$

## Example: BaB + Proof Generation



A simple DNN



A proof tree

Verifying $(x_1, x_2) \in [-2.0, 2.0] \times [-1.0, 1.0] \Rightarrow (y_1 > y_2)$

## Proof Tree

- Binary tree structure
  - each neuron decision branches into two children (representing on and off cases)
- Captures unsatisfiability reasoning of BaB
  - each leaf represents an unsatisfiable branch
  - $l_3 : \overline{v_4}, \qquad l_5 : v_4 \wedge \overline{v_2}, \qquad l_6 : v_4 \wedge v_2 \wedge v_1, \qquad l_7 : v_4 \wedge v_2 \wedge \overline{v_1}$

# Challenge 2: Human-Readability

## APTP Proof Language

$$
\begin{array}{rcl}
\langle\textit{proof}\rangle & ::= & \langle\textit{declarations}\rangle\ \langle\textit{assertions}\rangle \\
\langle\textit{declarations}\rangle & ::= & \langle\textit{declaration}\rangle\ |\ \langle\textit{declaration}\rangle\ \langle\textit{declarations}\rangle \\
\langle\textit{declaration}\rangle & ::= & (\textbf{declare-const}\ \langle\textit{input-vars}\rangle\ \textbf{Real}) \\
& & |\ (\textbf{declare-const}\ \langle\textit{output-vars}\rangle\ \textbf{Real}) \\
& & |\ (\textbf{declare-pwl}\ \langle\textit{hidden-vars}\rangle\ \langle\textit{activation}\rangle) \\
\langle\textit{input-vars}\rangle & ::= & \langle\textit{input-var}\rangle\ |\ \langle\textit{input-var}\rangle\ \langle\textit{input-vars}\rangle \\
\langle\textit{output-vars}\rangle & ::= & \langle\textit{output-var}\rangle\ |\ \langle\textit{output-var}\rangle\ \langle\textit{output-vars}\rangle \\
\langle\textit{hidden-vars}\rangle & ::= & \langle\textit{hidden-var}\rangle\ |\ \langle\textit{hidden-var}\rangle\ \langle\textit{hidden-vars}\rangle \\
\langle\textit{activation}\rangle & ::= & \text{ReLU}\ |\ \text{Leaky ReLU}\ |\ \ldots \\
\langle\textit{assertions}\rangle & ::= & \langle\textit{assertion}\rangle\ |\ \langle\textit{assertion}\rangle\ \langle\textit{assertions}\rangle \\
\langle\textit{assertion}\rangle & ::= & (\textbf{assert}\ \langle\textit{formula}\rangle) \\
\langle\textit{formula}\rangle & ::= & (\langle\textit{operator}\rangle\ \langle\textit{term}\rangle\ \langle\textit{term}\rangle) \\
& & |\ (\textbf{and}\ \langle\textit{formula}\rangle\,+)\ |\ (\textbf{or}\ \langle\textit{formula}\rangle\,+) \\
\langle\textit{term}\rangle & ::= & \langle\textit{input-var}\rangle\ |\ \langle\textit{output-var}\rangle \\
& & |\ \langle\textit{hidden-var}\rangle\ |\ \langle\textit{constant}\rangle \\
\langle\textit{operator}\rangle & ::= & <\ |\ \leq\ |\ >\ |\ \geq \\
\langle\textit{input-var}\rangle & ::= & \text{X\_}\langle\textit{constant}\rangle \\
\langle\textit{output-var}\rangle & ::= & \text{Y\_}\langle\textit{constant}\rangle \\
\langle\textit{hidden-var}\rangle & ::= & \text{N\_}\langle\textit{constant}\rangle \\
\langle\textit{constant}\rangle & ::= & \textbf{Int}\ |\ \textbf{Real}
\end{array}
$$

# Challenge 2: Human-Readability

## APTP Proof Language

$\langle proof \rangle ::= \langle declarations \rangle \langle assertions \rangle$

$\langle declarations \rangle ::= \langle declaration \rangle \mid \langle declaration \rangle \langle declarations \rangle$

$\langle declaration \rangle ::=$ (**declare-const** $\langle input\text{-}vars \rangle$ **Real**)

$\mid$ (**declare-const** $\langle output\text{-}vars \rangle$ **Real**)

$\mid$ (**declare-pwl** $\langle hidden\text{-}vars \rangle \langle activation \rangle$)

$\langle input\text{-}vars \rangle ::= \langle input\text{-}var \rangle \mid \langle input\text{-}var \rangle \langle input\text{-}vars \rangle$

$\langle output\text{-}vars \rangle ::= \langle output\text{-}var \rangle \mid \langle output\text{-}var \rangle \langle output\text{-}vars \rangle$

$\langle hidden\text{-}vars \rangle ::= \langle hidden\text{-}var \rangle \mid \langle hidden\text{-}var \rangle \langle hidden\text{-}vars \rangle$

$\langle activation \rangle ::=$ ReLU | Leaky ReLU | ...

$\langle assertions \rangle ::= \langle assertion \rangle \mid \langle assertion \rangle \langle assertions \rangle$

$\langle assertion \rangle ::=$ (**assert** $\langle formula \rangle$)

$\langle formula \rangle ::= (\langle operator \rangle \langle term \rangle \langle term \rangle)$

$\mid$ (**and** $\langle formula \rangle +$) | (**or** $\langle formula \rangle +$)

$\langle term \rangle ::= \langle input\text{-}var \rangle \mid \langle output\text{-}var \rangle$

$\mid \langle hidden\text{-}var \rangle \mid \langle constant \rangle$

$\langle operator \rangle ::= < \mid \leq \mid > \mid \geq$

$\langle input\text{-}var \rangle ::=$ X_$\langle constant \rangle$

$\langle output\text{-}var \rangle ::=$ Y_$\langle constant \rangle$

$\langle hidden\text{-}var \rangle ::=$ N_$\langle constant \rangle$

$\langle constant \rangle ::=$ **Int** | **Real**



## Example

```
1  ; Declare variables
2  (declare-const X_0 Real)
3  (declare-const X_1 Real)
4  (declare-const Y_0 Real)
5  (declare-const Y_1 Real)
6  (declare-pwl N_1 ReLU)
7  (declare-pwl N_2 ReLU)
8  (declare-pwl N_3 ReLU)
9  (declare-pwl N_4 ReLU)
10 ; Input constraints
11 (assert (>= X_0 -2.0))
12 (assert (<= X_0  2.0))
13 (assert (>= X_1 -1.0))
14 (assert (<= X_1  1.0))
15 ; Output constraints
16 (assert (<= Y_0 Y_1))
17 ; Hidden constraints
18 (assert (or
19    (and (<  N_4 0))              ; l_3
20    (and (<  N_2 0) (>= N_4 0))   ; l_5
21    (and (>= N_2 0) (>= N_1 0) (>= N_4 0)) ; l_6
22    (and (>= N_2 0) (<  N_1 0) (>= N_4 0)) ; l_7
23 ))
```

# Challenge 3: Efficiency

## APTP Checker Algorithm

**input**    : DNN $\mathcal{N}$, property $\phi_{in} \Rightarrow \phi_{out}$, *proof*
**output**   : *certified* if proof is valid, otherwise *uncertified*

1 **if** $\neg$ RepOK *(proof)* **then** RaiseError(*Invalid proof tree*) ;
2 $model \leftarrow$ CreateMILP$(\mathcal{N}, \phi_{in}, \phi_{out})$ `// initialize MILP model with inputs`
3 **while** *proof* **do**
4      $node \leftarrow$ Select(*proof*) `// get node to check`
5      $model \leftarrow$ AddConstrs(*model*, *node*) `// add corresponding constraints`
6      **if** CheckFeasibility(*model*) **then**
7          **return** *uncertified* `// cannot certify`

8 **return** *certified*

# Challenge 3: Efficiency

## APTP Checker Algorithm

**input** : DNN $\mathcal{N}$, property $\phi_{in} \Rightarrow \phi_{out}$, proof
**output** : certified if proof is valid, otherwise uncertified

1 **if** $\neg$ RepOK (proof) **then** RaiseError(Invalid proof tree) ;
2 model $\leftarrow$ CreateMILP($\mathcal{N}, \phi_{in}, \phi_{out}$) // initialize MILP model with inputs
3 **while** proof **do**
4      node $\leftarrow$ Select(proof) // get node to check
5      model $\leftarrow$ AddConstrs(model, node) // add corresponding constraints
6      **if** CheckFeasibility(model) **then**
7          **return** uncertified // cannot certify

8 **return** certified

## Optimizations

❶ **Stabilization:** fix stable ReLUs to reduce MILP constraints

❷ **Pruning:** check parent nodes and if unsat then skip children

❸ **Parallelization:** check multiple leaves in parallel

# Evaluation

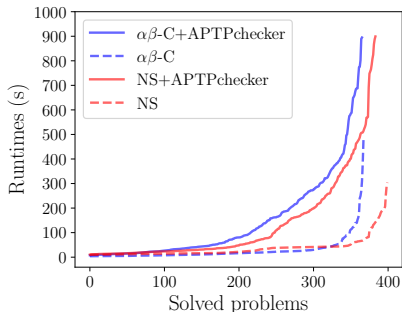| Name | Networks | | | | Instances |
|------|------|--------|---------|--------|-----------|
| | Num. | Layers | Neurons | Param. | Num. |
| CNN | 8 | 1-2C;1F | 320-3920 | 41K-180K | 200 |
| FNN | 8 | 2-6F | 64-3072 | 27K-1.7M | 200 |

**Benchmarks**

8 CNNs, 8 FNNs; 25 properties; total 400 problem instances

**Verifiers**

Add APTP proof generation to $\alpha\beta$-CROWN, `NeuralSAT`, and Marabou

# APTP Checker Performance



- Performance of verifiers (dashed) differ across configurations
  - They are able to verify between 337 and 400 problems
- APTP checker (solid) certified 93.7% and 99.4% of generated proofs
- This demonstrates that:
  - APTP proof format can encode proofs generated by different DNN verifiers
  - APTP checker can check them efficiently

# Trade-offs Between Verifiers

| Verifier | Num. Sub-Proofs | | MILP Complexity | |
|---|---|---|---|---|
| | Mean | Median | Mean | Median |
| NeuralSAT | 95 | 36 | 601 | 545 |
| $\alpha\beta$-CROWN | 230 | 180 | 414 | 179 |

- NeuralSAT: smaller proof trees, but with more complex MILP problems.
- $\alpha\beta$-CROWN: significantly larger proof trees, but with simpler MILP problems
- Strategies(?):
  - Generate larger proof trees with simpler MILPs for better parallelization
  - Adopting fast verification during development
  - Switching to proof-friendly strategies

# APTP Checker vs. Marabou Checker

| Checker | Proof checking time (seconds) | | |
|---------|:----:|:----:|:---:|
| | Mean | Median | Max |
| Marabou checker | 4 | 204 | 785 |
| APTP checker | 3 | 9 | 38 |

- Extract APTP proofs from 54 problems that Marabou verified
- For simple problems, both checkers handle quickly (similar mean times)
- For challenging problems, APTP checker is over $20\times$ speedup

# Key Takeaways

1. Verifiers producing unsound results: UNACCEPTABLE!
2. Proof generation technique applicable to a wide-range of verifiers
3. Standard, compact, and human-readable proof format
4. Efficient, lightweight, and verifier-independent proof checker

## Try APTP Checker:
https://github.com/dynaroars/aptpchecker

**Thank You!**