### Scalable DNN verification using Constraint Solving

ThanhVu (Vu) Nguyen











# **DNN EVERYWHERE**





### **DNN** Problems





Black person with hand-held thermometer = firearm. Asian person with hand-held thermometer = electronic device.

Computer vision is so utterly broken it should probably be started over from scratch.



Screenshot from 2020-03-31 11-23-45.png

Gun	88%
Photography	68%
Firearm	65%
Plant	59%









### **Robustness Properties**



L Diagnosis Result: Benign

**DL Diagnosis Result: Malignant** 

### Safety Properties



#### **DNN** Verification

Question: Given a network N and a property p, does N have p?

• p often has the form  $P \Rightarrow Q$  (precondition P, postcondition Q)

Answer: Yes / No

#### **DNN** Verification

Question: Given a network N and a property p, does N have p?

• p often has the form  $P \Rightarrow Q$  (precondition P, postcondition Q)

Answer: Yes / No



- Valid:  $x_1 \in [-1,1] \land x_2 \in [-2,2] \Rightarrow x_5 \leq 0$
- Invalid:  $x_1 \in [-1, 1] \land x_2 \in [-2, 2] \Rightarrow x_5 < 0$

#### Abstraction

- Overapproximate computation (e.g., ReLU) using abstract domains
  - interval (ReluVal), zonotopes (ERAN), polytopes ( $\alpha$ ,  $\beta$ -CROWN)



#### Abstraction

- Overapproximate computation (e.g., ReLU) using abstract domains
  - interval (ReluVal), zonotopes (ERAN), polytopes ( $\alpha, \beta$ -CROWN)



- Scale well, but *loose precision* (producing spurious cex's)
  - Newer work: iterative refine abstraction to filter spurious cex's

## Constraint Solving



# Constraint Solving



- Transform DNN verification into a constraint (satisfiability) problem
  - To prove  $N \Rightarrow p$  (where p is  $P \Rightarrow Q$ )
    - check if  $\neg(N \Rightarrow (P \Rightarrow Q))$ , i.e.,  $N \land P \land \neg Q$  is satisfiable
    - UNSAT: p is a property of N
    - SAT: p is not a property of N (also give counterexample inputs satisfiying P but not Q)

# Constraint Solving



- Transform DNN verification into a constraint (satisfiability) problem
  - To prove  $N \Rightarrow p$  (where p is  $P \Rightarrow Q$ )
    - check if  $\neg(N \Rightarrow (P \Rightarrow Q))$ , i.e.,  $N \land P \land \neg Q$  is satisfiable
    - UNSAT: p is a property of N
    - SAT: p is not a property of N (also give counterexample inputs satisfiying P but not Q)
- Solve the constraint(s)
  - SMT solvers (Planet, DLV) or customized simplex
  - MILP (Reluplex, Marabou)-based solvers
- Scalability is a HUGE problem

## Complexity and Scalability

Complexity: NP-Complete

- Scalability is the main problem
- State-of-the-art verification tools: networks with *138M* of parameters, 160K inputs
- Real-world networks: 3.5B parameters, 1.2M of inputs

## NeuralSAT: Our DNN Constraint Solver



Use NeuralSAT to prove  $N \Rightarrow (P \Rightarrow Q)$ 

- Call NeuralSAT( $N \land P \land \neg Q$ )
- Return UNSAT or SAT (and counterexample)

**Insight:** combines conflict clause learning in SAT solving and abstraction for scalability

### Example



To prove  $f: x_1 \in [-1,1] \land x_2 \in [-2,2] \Rightarrow x_5 \leq 0$ :

- NeuralSAT $(\neg f) =$ NeuralSAT $(N \land x_1 \in [-1, 1] \land x_2 \in [-2, 2] \land x_5 > 0)$
- NeuralSAT returns UNSAT, indicating f is valid





#### **Boolean Abstraction**

- Create 2 boolean variables v<sub>3</sub> and v<sub>4</sub> to represent *activation status* of x<sub>3</sub>, x<sub>4</sub>
  - $v_3 = T$  means  $x_3$  is active,  $-x_1 - 0.5x_2 - 1 > 0$
- Form two clauses  $\{v_3 \lor \overline{v_3} ; v_4 \lor \overline{v_4}\}$
- Find boolean values for v<sub>3</sub>, v<sub>4</sub> that satifies the clauses and their implications





#### Iteration 1

- Use abstraction to approximate upperbound  $x_5 \le 0.55$  (from  $x_1 \in [-1, 1], x_2 \in [-2, 2]$ )
- **Deduce**  $x_5 > 0$  *might be* feasible
- **Decide**  $v_3 = F$  (randomly)
  - new constraint  $-x_1 0.5x_2 1 < 0$





#### Iteration 2

- Approximate upperbound x<sub>5</sub> ≤ 0 (due to additional constraint from v<sub>3</sub> = F)
- **Deduce**  $x_5 > 0$  not feasible: CONFLICT
- Analyze conflict, backtrack and erase prev. decision  $v_3 = F$
- Learn new clause V3
  - *v*<sub>3</sub> will have to be *T* in next iteration





#### Iteration 3

• **Decide**  $v_3 = T$  (**BCP**, due to learned clause  $v_3$ )

• new constraint  $-x_1 - 0.5x_2 - 1 > 0$ 

- Approximate new upperbound for x<sub>5</sub> (using additional constraint from v<sub>3</sub> = T)
- **Deduce**  $x_5 > 0$  might be feasible
- **Decide**  $v_4 = T$  (randomly)
- :





#### After several iterations

- Learn clauses  $\{v_3, \overline{v_3} \lor v_4, \overline{v_3} \lor \overline{v_4}\}$
- **Deduce** not possible to satisfy the clauses

#### • Return UNSAT

- Cannot find inputs satisfying  $x_1 \in [-1, 1], x_2 \in [-2, 2]$  that cause N to return  $x_5 > 0$
- Hence, x<sub>5</sub> ≤ 0 holds (i.e., the original property is valid)

## NeuralSAT's Prototype and Preliminary Results

- Written in Python
- Accept standard DNN formats and specs
- Use DPLL/CDCL algorithms for clause learning and conflict analysis
- Use the polytope abstraction (can be replace with any other abstractions)

## ACAS XU Results

Much faster than the constraint solver Marabou

Prop	NeuralSAT	Marabou
$\phi_1$	1025.36	TO (3 hrs)
$\phi_2^*$	22.84	821.41
$\phi_3$	526.77	8309.09
$\phi_{ extsf{4}}$	330.83	133.97
$\phi_5$	83.51	1259.74
$\phi_{6}$	127.35	250.41
$\phi_7^*$	262.01	то
$\phi_8^*$	0.15	то
$\phi_{9}$	142.00	то
$\phi_{10}$	191.99	3134.35

#### ACAS XU Results Much faster than the constraint solver Marabou

-

Prop	NeuralSAT	Marabou
$\phi_1$	1025.36	TO (3 hrs)
$\phi_2^*$	22.84	821.41
$\phi_3$	526.77	8309.09
$\phi_4$	330.83	133.97
$\phi_5$	83.51	1259.74
$\phi_{6}$	127.35	250.41
$\phi_7^*$	262.01	ТО
$\phi_8^*$	0.15	ТО
$\phi_{9}$	142.00	ТО
$\phi_{10}$	191.99	3134.35

• Promising because NeuralSAT is a prototype with *no* optimizations

Still much slower than the abstraction tool nnenum

- nnenum applies a series of 7 optimizations
- comparable if **nnenum** runs using single thread

## Current Work / Future Directions

Current optimizations for NeuralSAT

- Parallize algorithms (e.g., Branch and Bound)
- Develop more precise (but still fast) abstraction
- Different search heuristics for boolean decisions

## Current Work / Future Directions

Current optimizations for NeuralSAT

- Parallize algorithms (e.g., Branch and Bound)
- Develop more precise (but still fast) abstraction
- Different search heuristics for boolean decisions

#### Future Directions

- Support richer specifications
- Mining specifications
- Apply formal reasoning (verification, specs. mining) to GNNs