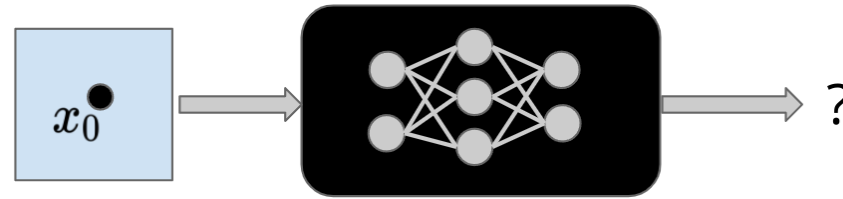


Introduction to Neural Network Verification

**Book Written by Prof. Aws Albarghouthi
UW-Madison**

**Presented by Yang Yang and Klaus Peng
2023/10/11**

Neural Network Verification

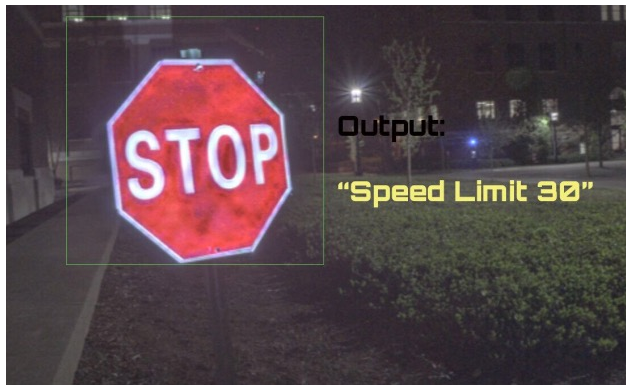


Neural network (f) works as a black box.

$$f : R^n \rightarrow R^m$$

Challenge

The network could behave **unexpectedly** and produce **wrong** results.



Goal

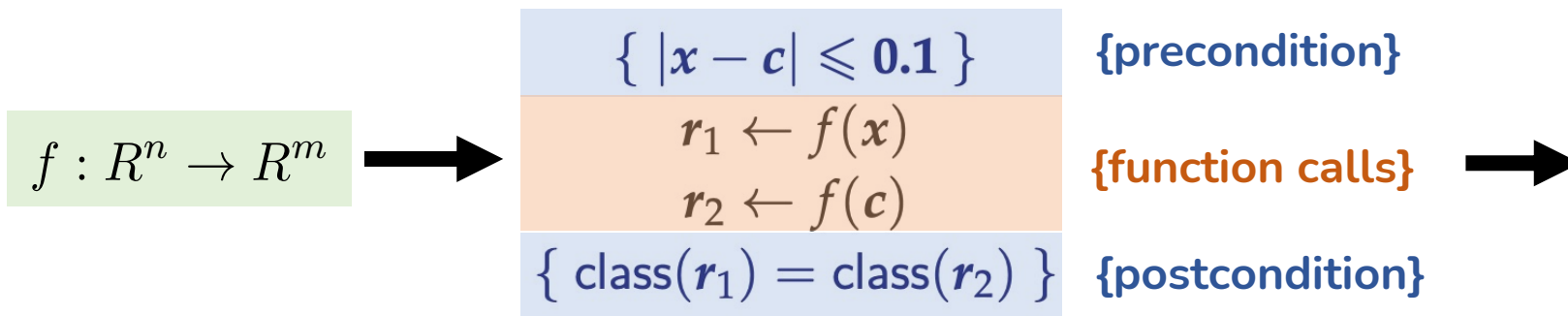
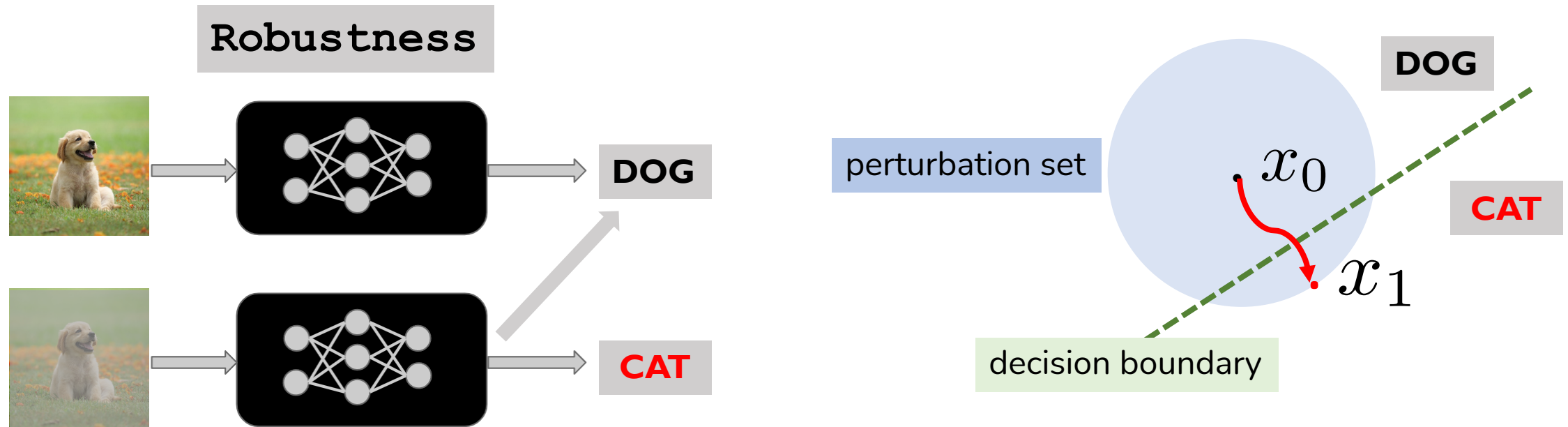
We want to make NN have some desired **properties** we can **formally trust**.



Neural Network Verification

Neural Network Properties

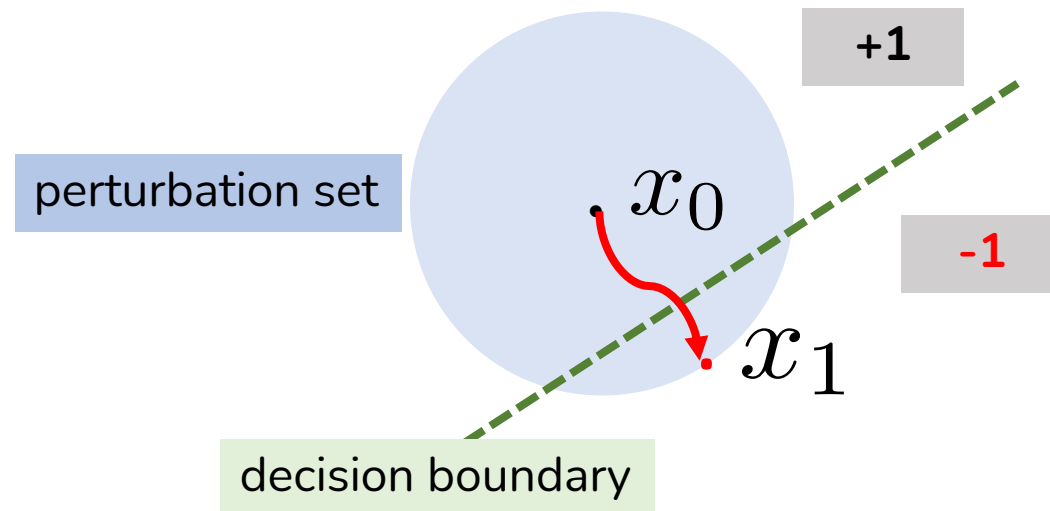
We want to make NN have some desired **properties** we can **formally** trust.



For any values of x_1, \dots, x_n *that* make the **precondition** **true**, let $r_1 = f(x_1), r_2 = g(x_2), \dots$ *Then* the **postcondition** is **true**.

$pre \wedge F \Rightarrow post$

Neural Network Verification



Assuming $f(x_0) > 0$, we solve the optimization problem to find the worst case:

$$f^* = \min \{ f(x) \}, x \in P$$

Where P is usually a perturbation set on x_0 , e.g.,

$$P = \{ x \mid \|x - x_0\|_p \leq \epsilon \}$$

$$f^* > 0$$

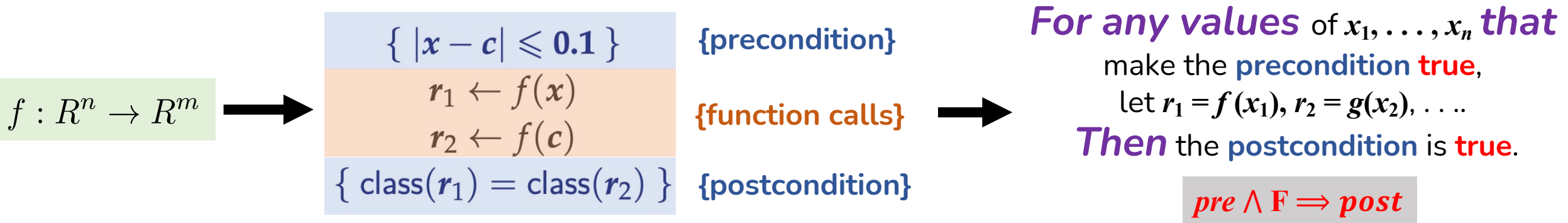
Challenge

In real-world case, non-linear, non-convex constraints and generally leads to **NP-complete**.

Review

Goal

We want to make NN
have some desired **properties** we can **formally** trust.



Challenge

In real-world case, non-linear, non-convex constraints and generally leads to **NP-complete**.

NEXT?

Precondition, postcondition, true, false, **NP-complete**, ...

$pre \wedge F \Rightarrow post$

Boolean Logic

SAT

Propositional Logic Recap

Proposition is a **sentence** that could be either **true** or **false**.

A: SE4ML is an interesting course. (X)

B: SE4ML is a graduate-level course. (✓)

Formula is a **string** that could be **generated** using proposition.

Implication: $A \Rightarrow B := \neg A \vee B$ *pre \wedge F \Rightarrow post*

$F := F \wedge F \mid F \vee F \mid \neg F$

$F \triangleq (p \wedge q) \vee \neg r$

Interpretation is a **set** that assigns true or false to **fv**'s elements.

$fv(F) = \{p, q, r\}$

$I = \{p \mapsto \text{true}, q \mapsto \text{true}, r \mapsto \text{false}\}$

Formula **F** is satisfiable (SAT) if there **exists** an interpretation **I** such that $\text{eval}(I(F)) = \text{true}$.

P: **Solve** the problem in polynomial time, bubble sort.

NP: **Validate** a solution in polynomial time, TSP.

NP-Complete: An NP problem that every NP problem is reducible to it, The Satisfiability Problem(SAT).

SAT: $\forall n_1, \dots, n_N$, whether $\text{eval}(F(n_1, \dots, n_N)) = \text{true}$

How you solve it?

$F := A \vee \neg A$

[STOC 1971, Stephen Cook] **The Complexity of Theorem-Proving Procedures**

Proved SAT to be the first NP-Complete using Turing Machine. \Rightarrow 1982 ACM Turing Award



Arithmetic Logic

Satisfiability Modulo Theories (SMT)

$$F \triangleq p \vee q, fv(F) = \{p, q\}$$

$$p \leftarrow a + b \leq 10$$

$$q \leftarrow c = d[1]$$

$$F' \triangleq a + b \leq 10 \vee c = d[1]$$

$$fv(F') = \{a, b, c, d\}$$

Boolean Variable

Boolean Expression

Linear Inequality

Linear Real Arithmetic (LRA)

$$\sum_{i=1}^n c_i x_i + b \leq 0 \text{ or } \sum_{i=1}^n c_i x_i + b < 0$$

Non-Linear Arithmetic

Arithmetic Logic

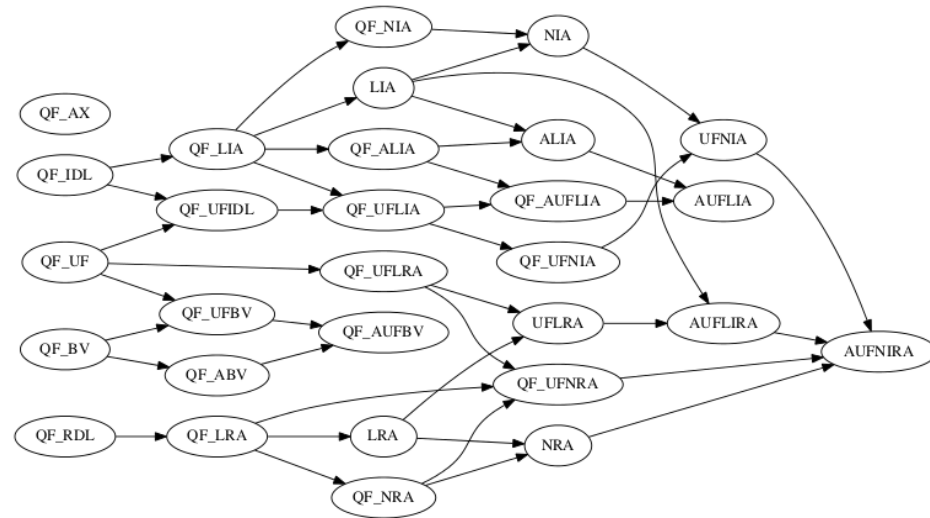


Figure 1: Overview of logics in SMT-LIB

Advantages of SMT Solvers

- Allows encoding using **various** theories/logics
- **Ability** to combine multiple theories

vs.

MILP Solver

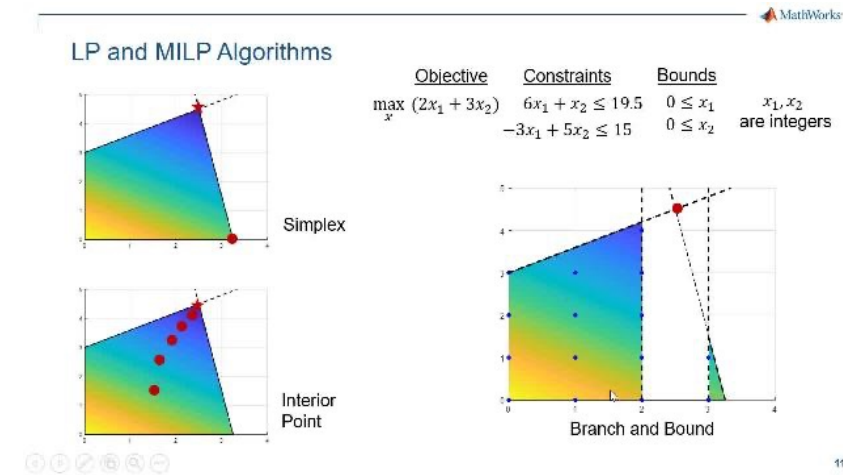


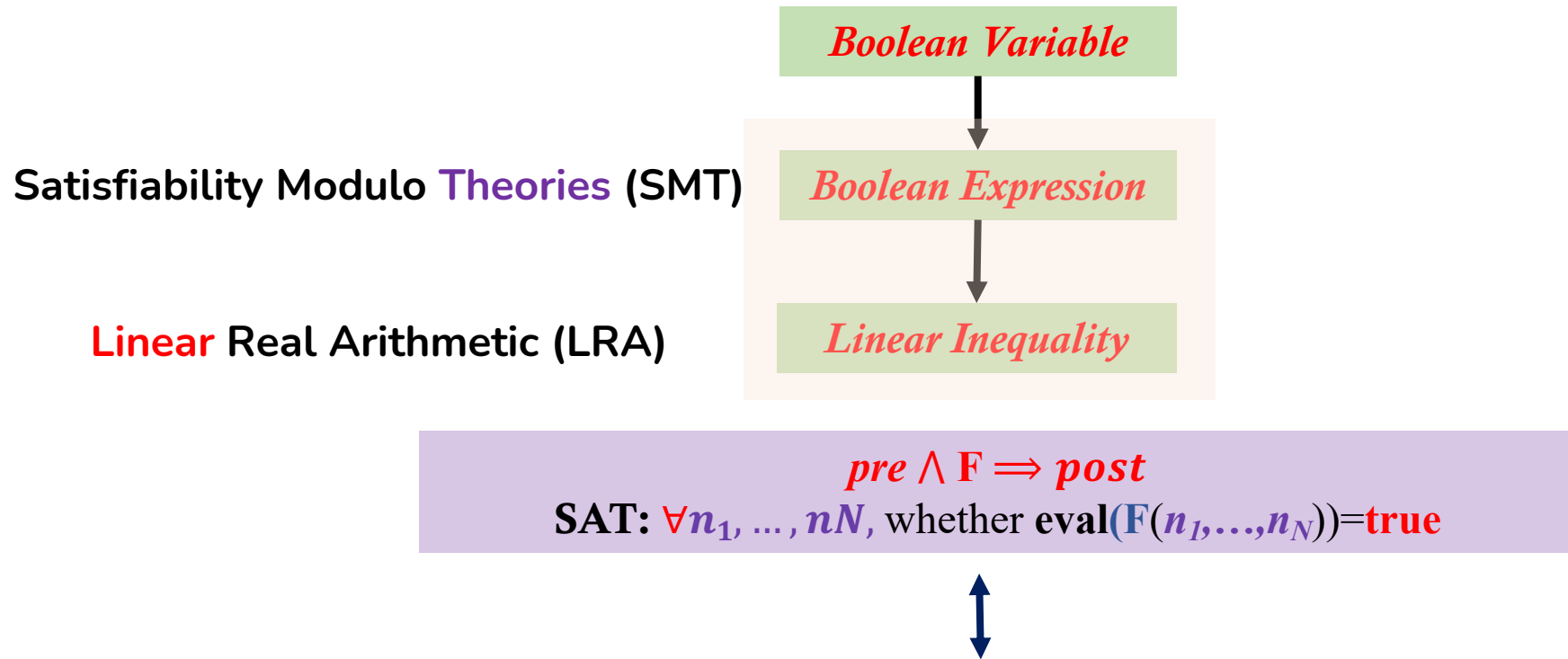
Figure 2: provided by Matlab

Advantages of MILP Solver

- Sometimes, **faster** compared to SMT solver

Review

We have tools **how** it could be used in NNV? Why we need such **logic**?



Use SMT to define the **specifications** of our neural network
Apply SMT solvers to **verify** if our model is robust or of other properties

NEXT?

Let's verify the real Neural Network !

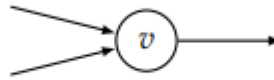
Encoding Nodes

Example of Nodes:



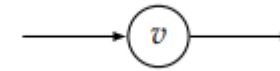
$$f_v(x) = 2x + 1$$

$$F_v \triangleq v^o = 2x + 1$$



$$f_v(x, y) = 3x - 2y$$

$$F_v \triangleq v^o = 3x - 2y$$



$$f_v(x) = \text{relu}(x) = \begin{cases} x & \text{if } x > 0 \\ 0 & \text{if } x \leq 0 \end{cases}$$

$$F_v \triangleq (x > 0 \Rightarrow v^o = x) \wedge (x \leq 0 \Rightarrow v^o = 0)$$

Formalized definition:

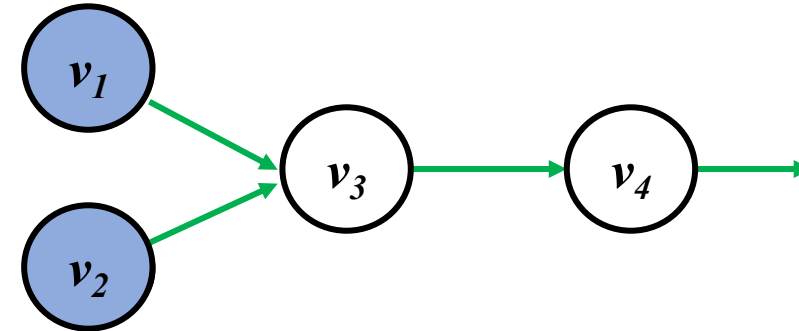
$$f_v(\mathbf{x}) = \begin{cases} \sum_{j=1}^{n_v} c_j^1 \cdot x_j + b^1 & \text{if } S_1 \\ \vdots \\ \sum_{j=1}^{n_v} c_j^l \cdot x_j + b^l & \text{if } S_l \end{cases} \quad \Rightarrow \quad F_v \triangleq \bigwedge_{i=1}^l [S_i \Rightarrow (v^o = \sum_{j=1}^{n_v} c_j^i \cdot v^{in,j} + b^i)]$$

v^o - output of node v
 $v^{in,i}$ - i th input to node v

Encoding Neural Network

1. Encoding Semantics of Nodes

$$F_V \triangleq \bigwedge_{v \in V \setminus V^{in}} F_v$$



$$V^{in} = \{v_1, v_2\}, V^o = \{v_4\}, f_{v_3} = 2x_1 - 3x_2, f_{v_4} = \text{relu}(x)$$

Formula for
each **non-input**
nodes:

$$F_{v_3} \triangleq v_3^o = 2v_3^{in,1} - 3v_3^{in,2}$$
$$F_{v_4} \triangleq (v_4^{in,1} > 0 \Rightarrow v_4^o = v_4^{in,1}) \wedge (v_4^{in,1} \leq 0 \Rightarrow v_4^o = 0)$$

Put them together:

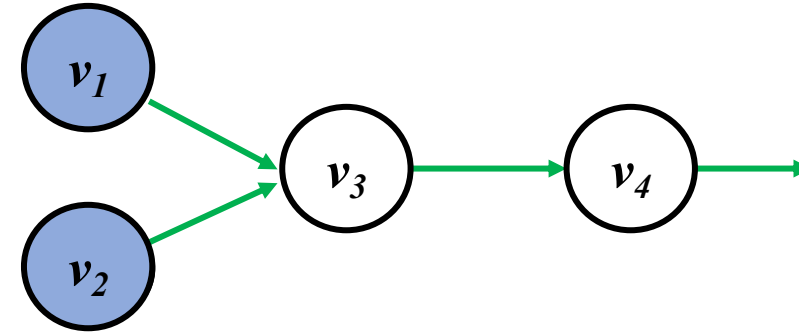
$$F_V \triangleq F_{v_3} \wedge F_{v_4}$$

Encoding Neural Network

2. Encoding Edges Between Nodes

$$F_{o \rightarrow v} \triangleq \bigwedge_{i=1}^n v^{in,i} = v_j^o$$

$$F_E \triangleq \bigwedge_{v \in V \setminus V^{in}} F_{o \rightarrow v}$$



$$V^{in} = \{v_1, v_2\}, V^o = \{v_4\}, f_{v_3} = 2x_1 - 3x_2, f_{v_4} = \text{relu}(x)$$

Formula for
each edge:

$$F_{o \rightarrow v_3} \triangleq (v_3^{in,1} = v_1^o) \wedge (v_3^{in,2} = v_2^o)$$

$$F_{o \rightarrow v_4} \triangleq v_4^{in,1} = v_3^o$$

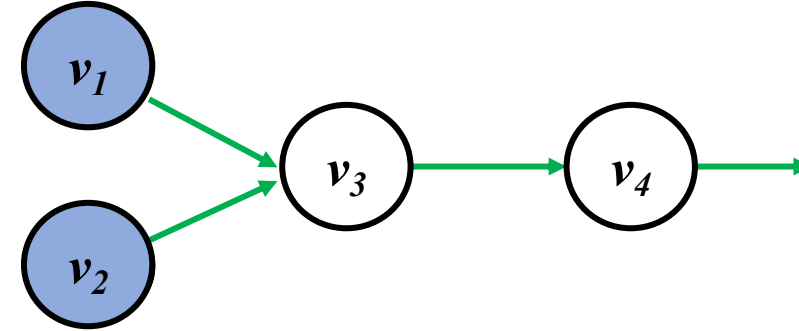
Put them together:

$$F_E \triangleq F_{o \rightarrow v_3} \wedge F_{o \rightarrow v_4}$$

Encoding Neural Network

3. All Together

$$F_G \triangleq F_V \wedge F_E$$



$$V^{in} = \{v_1, v_2\}, V^o = \{v_4\}, f_{v_3} = 2x_1 - 3x_2, f_{v_4} = \text{relu}(x)$$

Node & Edge:

$$F_{v_3} \triangleq v_3^o = 2v_3^{in,1} - 3v_3^{in,2}$$

$$F_{v_4} \triangleq (v_4^{in,1} > 0 \Rightarrow v_4^o = v_4^{in,1}) \wedge (v_4^{in,1} \leq 0 \Rightarrow v_4^o = 0)$$

$$F_V \triangleq F_{v_3} \wedge F_{v_4}$$

$$F_{o \rightarrow v_3} \triangleq (v_3^{in,1} = v_1^o) \wedge (v_3^{in,2} = v_2^o)$$

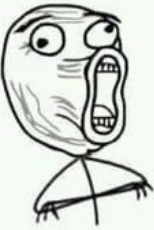
$$F_{o \rightarrow v_4} \triangleq v_4^{in,1} = v_3^o$$

$$F_E \triangleq F_{o \rightarrow v_3} \wedge F_{o \rightarrow v_4}$$

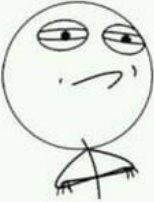
Finally:

$$F_G \triangleq F_{v_3} \wedge F_{v_4} \wedge F_{o \rightarrow v_3} \wedge F_{o \rightarrow v_4}$$


Encoding Correctness Properties



School: $2+2 = 4$



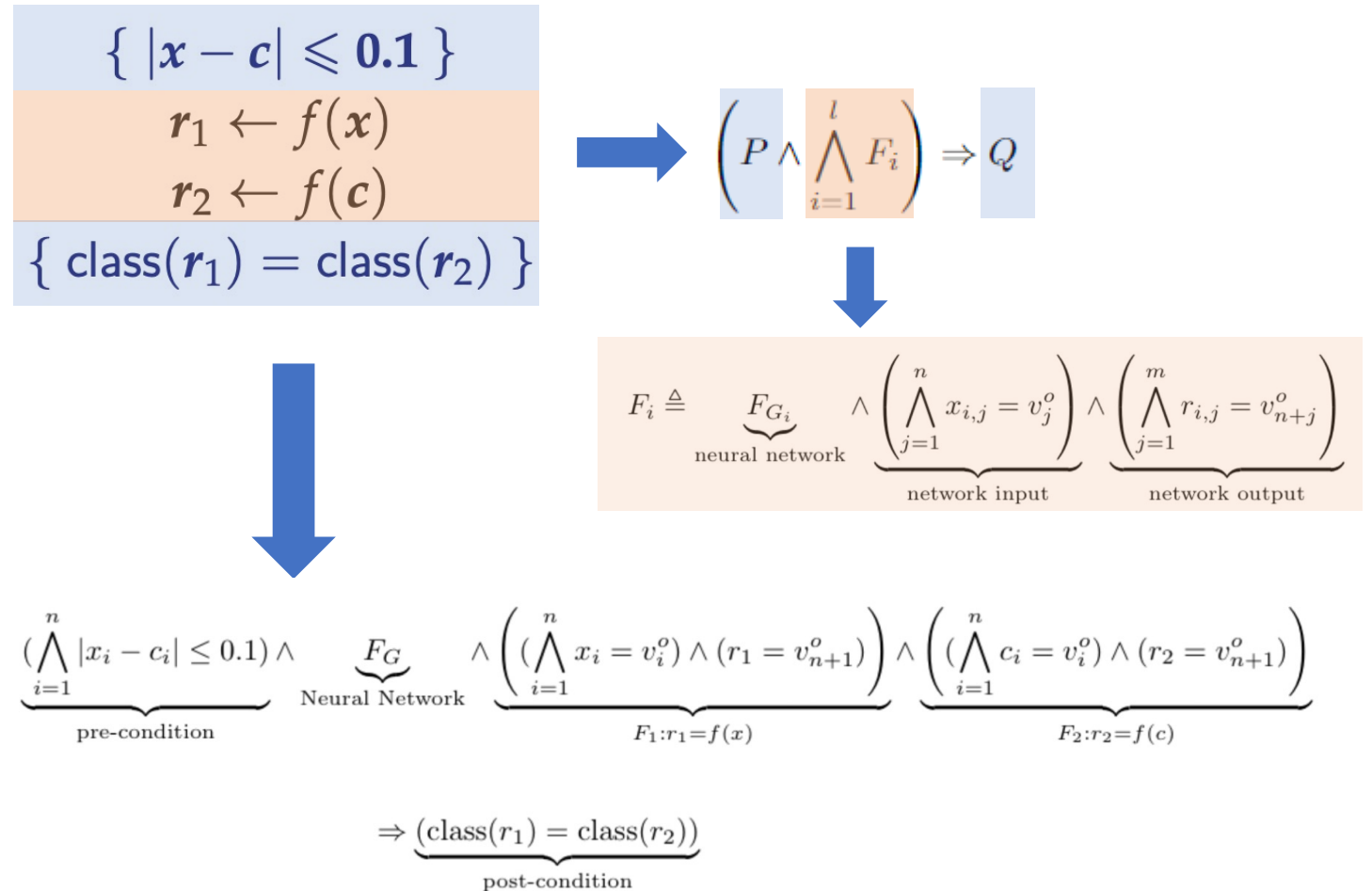
Homework: $2+3+4 = 9$



Exam: David has 4 apples, his train is 7 minutes early, calculate mass of the sun.

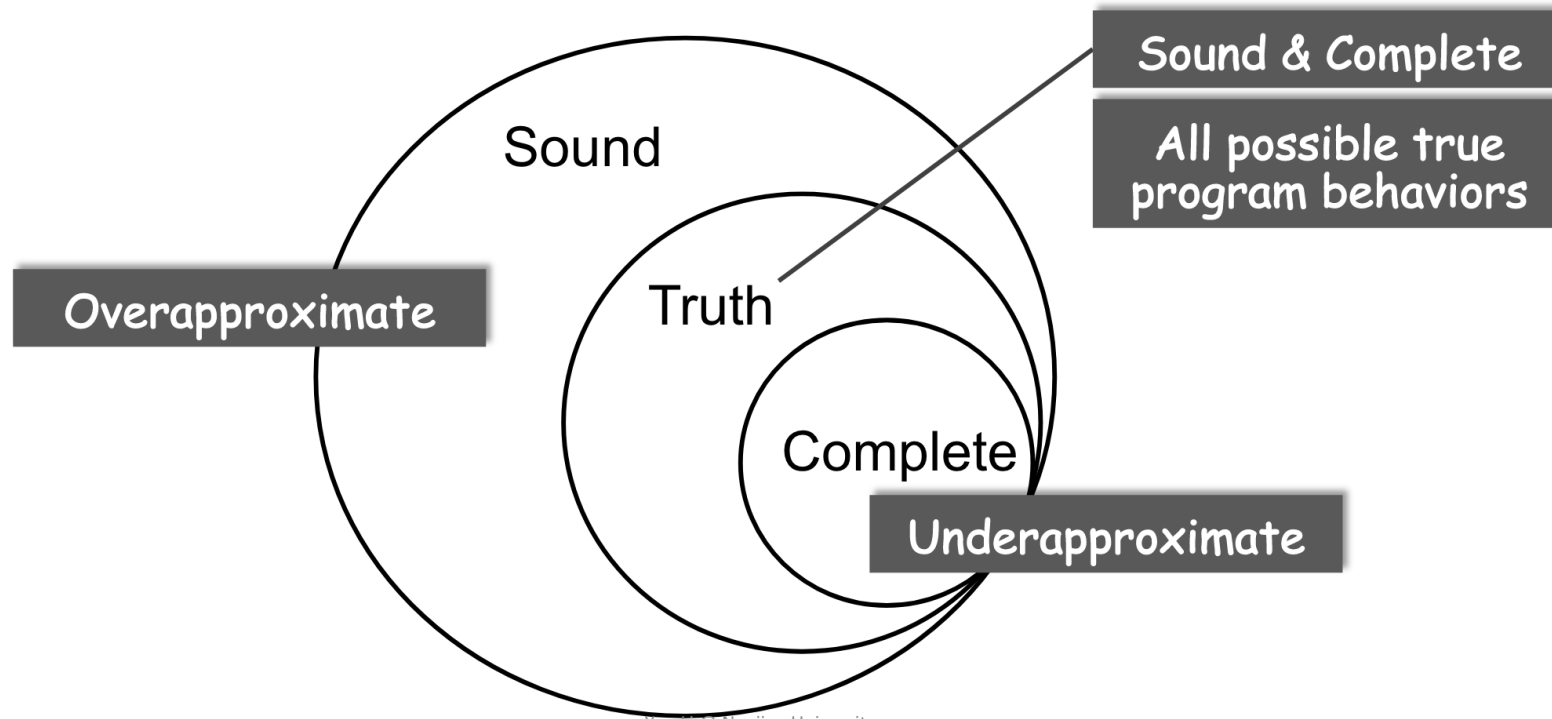
Recall previously:

$$pre \wedge F \Rightarrow post$$



Soundness & Completeness

If you are interested, check [Rice's Theorem](#), [PL](#) and [Static Analysis](#).



Conclusion

Why Neural Network Verification?

We want to make NN
have some desired **properties** we can **formally** trust.

How We Form it?

pre** \wedge **F** \Rightarrow **post

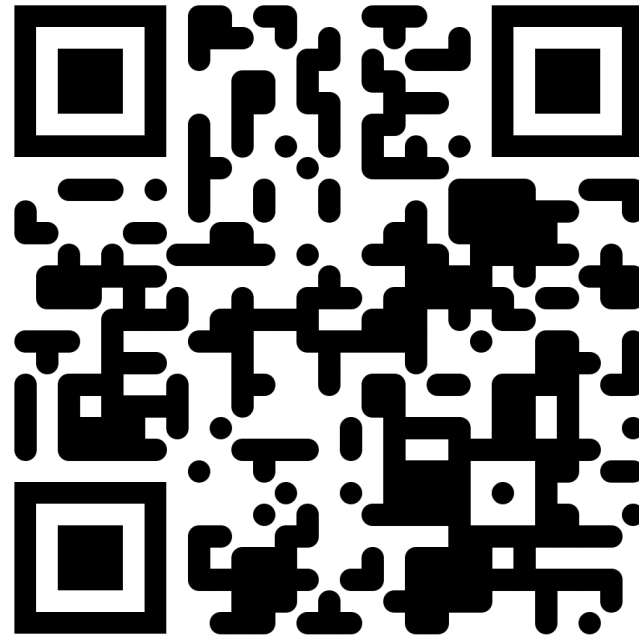
SAT: $\forall n_1, \dots, n_N$, whether $\text{eval}(\mathbf{F}(n_1, \dots, n_N)) = \text{true}$

How We Solve and Use it?

Use SMT to define the **specifications** of our neural network.
Apply SMT solvers to **verify** if our model is robust or of other properties.
Encode nodes, edges, properties and the network.

Questions?

Neural Network Verification IN ACTION



Instructions: <https://elio-yang.github.io/res/NNV>