# Safe Neurosymbolic Learning with Differentiable Symbolic Execution

**Chenxi Yang**, Swarat Chaudhuri

The University of Texas at Austin

# Motivation

- Learning techniques are being used in safety-critical tasks.

**Practical uses of Deep Neural Networks in Healthcare**

deepakvraghavan  May 2, 2018 · 9 min read

*The Costly Pursuit of Self-Driving Cars Continues On. And On. And On.*

Many in Silicon Valley promised that self-driving cars would be a common sight by 2021. Now the industry is resetting expectations and settling in for years of more work.

- In the real world, neural networks are often invoked by human-written code.

Provably Safe Neurosymbolic Learning

# Safe Neurosymbolic Learning

- We focus on imitation learning.

# Safe Neurosymbolic Learning

- We focus on imitation learning.
  - Human written code invoking neural networks

```
1  thermostat(x):
2      N = 20
3      isOn = 0.0
4      i = 0
5      while i < N:
6          if isOn ≤ 0.5:
7              isOn := π_θ^cool(x)
8              x := COOLING(x)
9          else:
10             isOn, heat := π_θ^heat(x)
11             x := WARMING(x, heat)
12         i := i + 1
13         assert(!EXTREME_TEMPERATURE(x))
14
15     return
```

# Safe Neurosymbolic Learning

- We focus on imitation learning.
  - Human written code invoking neural networks

```
1  thermostat(x):
2      N = 20
3      isOn = 0.0
4      i = 0
5      while i < N:
6          if isOn ≤ 0.5:
7              isOn := π_θ^cool(x)
8              x := COOLING(x)
9          else:
10             isOn, heat := π_θ^heat(x)
11             x := WARMING(x, heat)
12         i := i + 1
13         assert(!EXTREME_TEMPERATURE(x))
14
15     return
```

2 different NNs

# Safe Neurosymbolic Learning

- We focus on imitation learning.
  - Human written code invoking neural networks

```
1   thermostat(x):
2       N = 20
3       isOn = 0.0
4       i = 0
5       while i < N:
6           if isOn ≤ 0.5:
7               isOn := π_θ^cool(x)
8               x := COOLING(x)
9           else:
10              isOn, heat := π_θ^heat(x)
11              x := WARMING(x, heat)
12          i := i + 1
13          assert(!EXTREME_TEMPERATURE(x))
14
15      return
```

Symbolic code

# Safe Neurosymbolic Learning

- We focus on imitation learning
  - Human written code invoking neural networks
  - A trajectory dataset for imitation

```
1   thermostat(x):
2       N = 20
3       isOn = 0.0
4       i = 0
5       while i < N:
6           if isOn ≤ 0.5:
7               isOn := π_θ^cool(x)
8               x := COOLING(x)
9           else:
10              isOn, heat := π_θ^heat(x)
11              x := WARMING(x, heat)
12          i := i + 1
13          assert(!EXTREME_TEMPERATURE(x))
14
15      return
```

20-length trajectories

# Safe Neurosymbolic Learning

- We focus on imitation learning
  - Human written code invoking neural networks
  - A trajectory dataset for imitation
    - Trajectories are sequences of input-output pair of neural networks

```
1  thermostat(x):
2     N = 20
3     isOn = 0.0
4     i = 0
5     while i < N:
6        if isOn ≤ 0.5:
7           isOn := π_θ^cool(x)
8           x := COOLING(x)
9        else:
10          isOn, heat := π_θ^heat(x)
11          x := WARMING(x, heat)
12       i := i + 1
13       assert(!EXTREME_TEMPERATURE(x))
14
15    return
```

Input-output pairs

# Safe Neurosymbolic Learning

- We focus on imitation learning
    - Human written code invoking neural networks
    - A trajectory dataset for imitation
        - Trajectories are sequences of input-output pair of neural networks
        - In this example, trajectories are in the form of

```
<([x], [isOn], 'cool'), ([x], [isOn,
heat], 'heat'), ([x], [isOn, heat],
'heat'), … >
```

```
1   thermostat(x):
2       N = 20
3       isOn = 0.0
4       i = 0
5       while i < N:
6           if isOn ≤ 0.5:
7               isOn := π_θ^cool(x)
8               x := COOLING(x)
9           else:
10              isOn, heat := π_θ^heat(x)
11              x := WARMING(x, heat)
12          i := i + 1
13          assert(!EXTREME_TEMPERATURE(x))
14
15      return
```

Input-output pairs

# Safe Neurosymbolic Learning

- We focus on imitation learning
  - Human written code invoking neural networks
  - A trajectory dataset for imitation
    - Trajectories are sequences of input-output pair of neural networks
    - In this example, trajectories are in the form of

Input of NN

```
<([x], [isOn], 'cool'), ([x], [isOn,
heat], 'heat'), ([x], [isOn, heat],
'heat'), … >
```

Input-output pairs

```
1  thermostat(x):
2      N = 20
3      isOn = 0.0
4      i = 0
5      while i < N:
6          if isOn ≤ 0.5:
7              isOn := π_θ^cool(x)
8              x := COOLING(x)
9          else:
10             isOn, heat := π_θ^heat(x)
11             x := WARMING(x, heat)
12         i := i + 1
13     assert(!EXTREME_TEMPERATURE(x))
14
15     return
```

# Safe Neurosymbolic Learning

- We focus on imitation learning
  - Human written code invoking neural networks
  - A trajectory dataset for imitation
    - Trajectories are sequences of input-output pair of neural networks
    - In this example, trajectories are in the form of
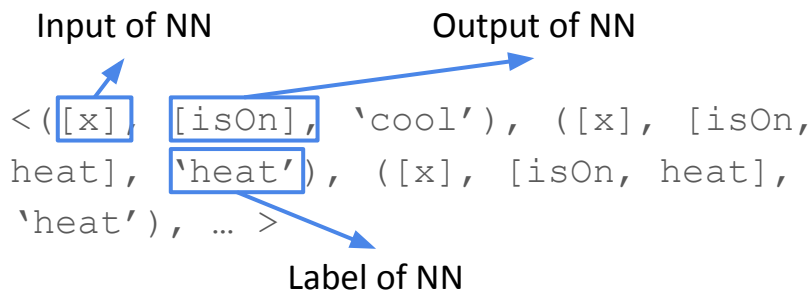
Input of NN          Output of NN

```
<([x], [isOn], 'cool'), ([x], [isOn,
heat], 'heat'), ([x], [isOn, heat],
'heat'), … >
```

Input-output pairs

```
1   thermostat(x):
2       N = 20
3       isOn = 0.0
4       i = 0
5       while i < N:
6           if isOn ≤ 0.5:
7               isOn := π_θ^cool(x)
8               x := COOLING(x)
9           else:
10              isOn, heat := π_θ^heat(x)
11              x := WARMING(x, heat)
12          i := i + 1
13          assert(!EXTREME_TEMPERATURE(x))
14
15      return
```

# Safe Neurosymbolic Learning

- We focus on imitation learning
  - Human written code invoking neural networks
  - A trajectory dataset for imitation
    - Trajectories are sequences of input-output pair of neural networks
    - In this example, trajectories are in the form of

Input of NN          Output of NN

```
<([x], [isOn], 'cool'), ([x], [isOn,
heat], 'heat'), ([x], [isOn, heat],
'heat'), … >
```

Label of NN

```
1   thermostat(x):
2       N = 20
3       isOn = 0.0
4       i = 0
5       while i < N:
6           if isOn ≤ 0.5:
7               isOn := π_θ^cool(x)
8               x := COOLING(x)
9           else:
10              isOn, heat := π_θ^heat(x)
11              x := WARMING(x, heat)
12          i := i + 1
13      assert(!EXTREME_TEMPERATURE(x))
14
15  return
```

Input-output pairs

# Safe Neurosymbolic Learning

- We focus on imitation learning.
  - Human written code invoking neural networks
  - A trajectory dataset for imitation
  - Constraints on the input

```
1  thermostat(x):
2     N = 20
3     isOn = 0.0
4     i = 0
5     while i < N:
6        if isOn ≤ 0.5:
7           isOn := π_θ^cool(x)
8           x := COOLING(x)
9        else:
10          isOn, heat := π_θ^heat(x)
11          x := WARMING(x, heat)
12       i := i + 1
13       assert(!EXTREME_TEMPERATURE(x))
14
15    return
```

$$x \in [55.0, 70.0]$$

# Safe Neurosymbolic Learning

- We focus on imitation learning
  - Human written code invoking neural networks
  - A trajectory dataset for imitation
  - Constraints on the input
  - Constraints on the trajectory of the system

```
1   thermostat(x):
2       N = 20
3       isOn = 0.0
4       i = 0
5       while i < N:
6           if isOn ≤ 0.5:
7               isOn := π_θ^cool(x)
8               x := COOLING(x)
9           else:
10              isOn, heat := π_θ^heat(x)
11              x := WARMING(x, heat)
12          i := i + 1
13          assert(!EXTREME_TEMPERATURE(x))
14
15      return
```

Safety constraints

# Formulation

Fitting the trajectory dataset
Data Signal: $Q(\theta)$

Satisfying the safety constraint
Safety Signal: $C(\theta)$

$$\min_{\theta} Q(\theta) \qquad \text{s.t. } C(\theta) \leq 0.$$

# Safe Learning Framework

Fitting the trajectory dataset
Data Signal: $Q(\theta)$

Satisfying the safety constraint
Safety Signal: $C(\theta)$

$$\min_{\theta} Q(\theta) \qquad \text{s.t.} \ \ C(\theta) \leq 0.$$

Using Lagrange multipliers

$$L(\theta, \lambda) = Q(\theta) + \lambda C(\theta)$$

High-level: the framework
repeatedly solves this
optimization problem

$$\min_{\theta} Q(\theta) + \lambda C(\theta)$$

Challenge: Estimate the
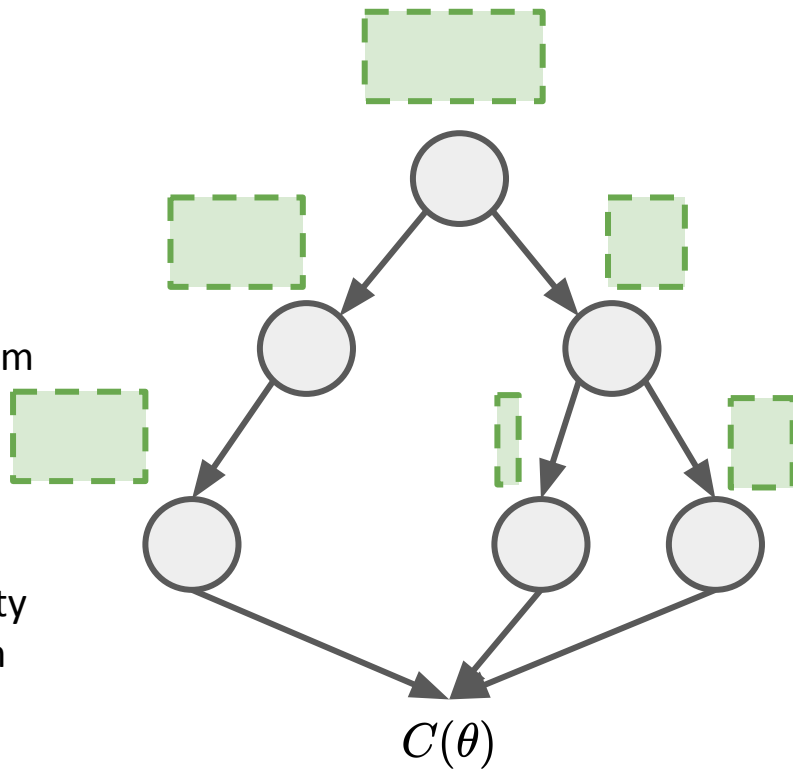gradient of safety signal
in a differentiable way.

# Verification Technique: Symbolic Execution



Symbolic input

Propagate the symbolic input through the program

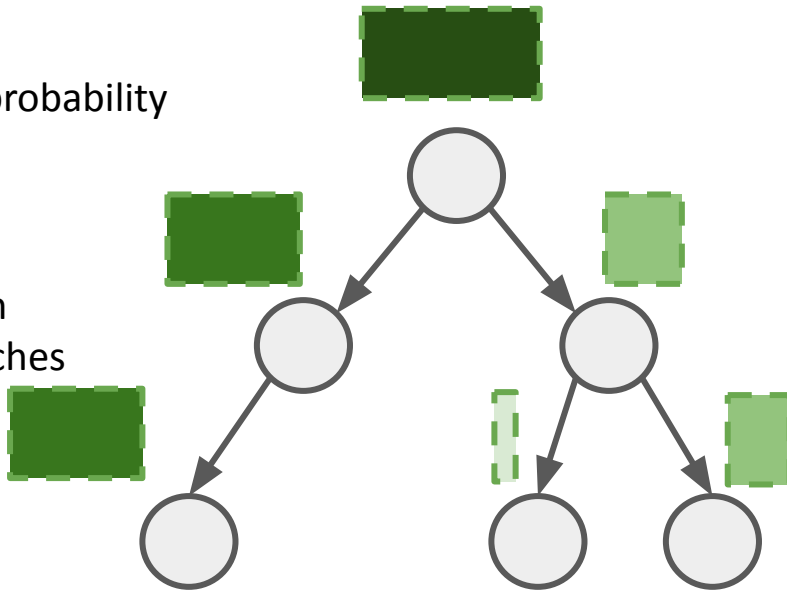Aggregate the safety measurement from all paths

*! Two Issues*

Path explosion: Branch X Loops

Non-differentiability: Discrete branches

$C(\theta)$

# DSE: Differentiable Symbolic Execution

Volume weighted probability

Sample paths when encountering branches

Use symbolic REINFORCE to estimate the gradient of safety loss.

$$p_\theta(t_j|\sigma_i) = \frac{\mathbf{Vol}(G_j \wedge V_i)}{\mathbf{Vol}(V_i)}$$

The volume of the symbolic state that satisfies G.

The entire volume of the symbolic state before the conditional.

$$p_\theta(\tau_\theta^{\#}) = p(\sigma_0) \prod_i p_\theta(t_i|\sigma_{i-1}).$$

$$C^{\#}(\theta) = \mathbf{E}_{\tau^{\#} \sim p_\theta(\tau^{\#})} Unsafe_\theta(\tau^{\#}).$$

# DSE: Differentiable Symbolic Execution

$$\nabla_\theta(C^\#(\theta))$$

Estimate the gradient

$$= \quad \nabla_\theta \mathbf{E}_{\tau^\# \sim p_\theta(\tau^\#)} Unsafe_\theta(\tau^\#)$$

$$= \quad \boxed{\mathbf{E}_{\tau^\# \sim p_\theta(\tau^\#)}[\nabla_\theta Unsafe_\theta(\tau^\#)]} + \boxed{\mathbf{E}_{\tau^\# \sim p_\theta(\tau^\#)}[Unsafe_\theta(\tau^\#)\nabla_\theta(\log p_\theta(\tau^\#))]}$$

Penalty is a function over NN's parameter
Exhibit the symbolic unsafety change of a path

Follow the classic REINFORCE estimator
Exhibit the probability change of a symbolic path

# Safe Learning Framework

Data Signal: $Q(\theta)$
Fitting the dataset

Safety Signal: $C(\theta)$
Satisfying the safety constraint

$$\min_{\theta} Q(\theta) \qquad \text{s.t.} \ \ C(\theta) \leq 0.$$

$$\min_{\theta} Q(\theta) + \lambda C(\theta)$$

Challenge: Estimate the gradient of safety signal in a differentiable way.

# Safe Learning Framework



Data Signal: $Q(\theta)$
Fitting the dataset

Safety Signal: $C(\theta)$
Satisfying the safety constraint

$$\min_{\theta} Q(\theta) \quad \text{s.t.} \ C(\theta) \leq 0.$$

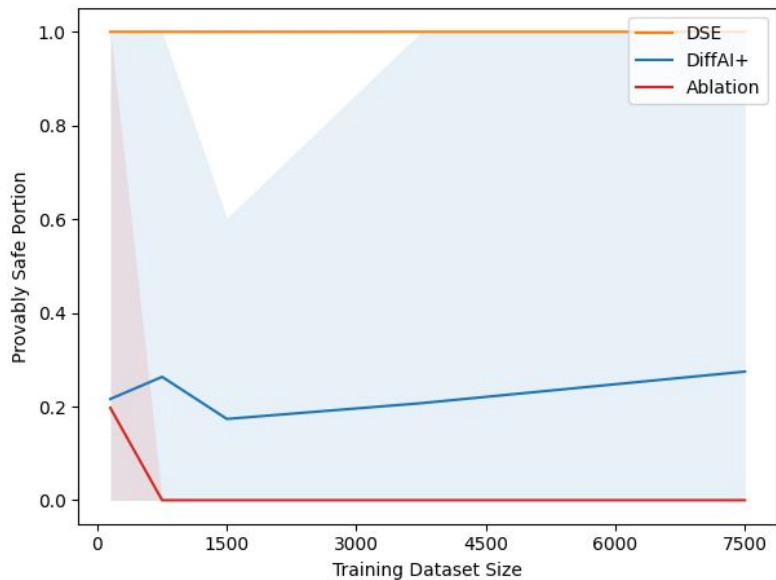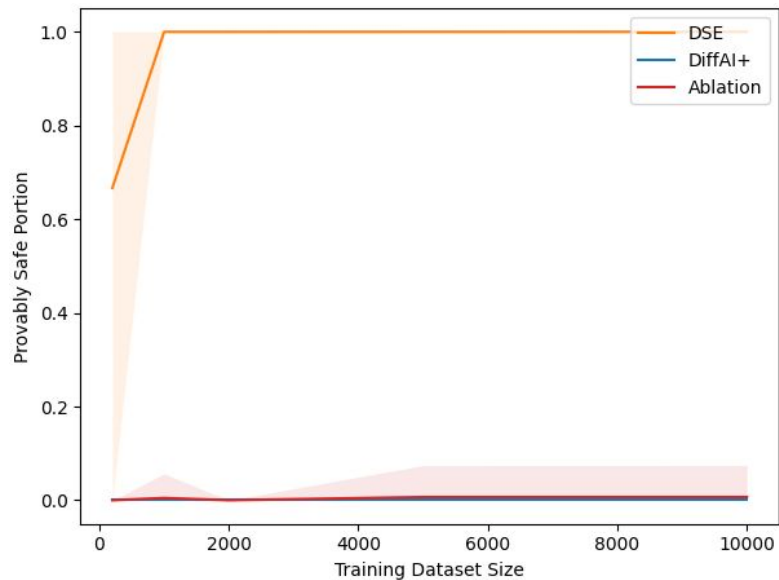$$\min_{\theta} Q(\theta) + \lambda \boxed{C^{\#}(\theta)}$$

**D**ifferentiable **S**ymbolic
Execution
(**DSE**)

Approximate Safety Loss

# Experiments

Four benchmarks

- Thermostat: 2 NN controllers and ~200 dynamic lines of code
- Racetrack: 2 competing vehicles controlled by NN and interactions with the map
- Aircraft collision: 1 aircraft NN controller and $4^{15}$ paths in maximum
- Cartpole: 1 cart pole controlled by NN

# Evalution



The larger y, the safer learnt progra. DSE performs much better than two baselines: Ablation and DiffAI+.

# Summary

- We provide DSE, a differentiable way to combine verification techniques over symbolic code with safe learning.
- In practice, DSE can learn provably safer programs than SOTA.

# Thank you!

Paper: http://arxiv.org/abs/2203.07671
Code: https://github.com/cxyang1997/DSE