

Neuro-Symbolic Execution: Augmenting Symbolic Execution with Neural Constraints

Shiqi Shen Shweta Shinde Soundarya Ramesh

Abhik Roychoudhury Prateek Saxena

National University of Singapore

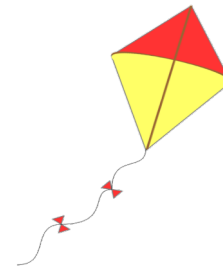
Symbolic Execution for Bug Finding



Pex



Angr



Kite

SAGE

jCUTE



Manticore

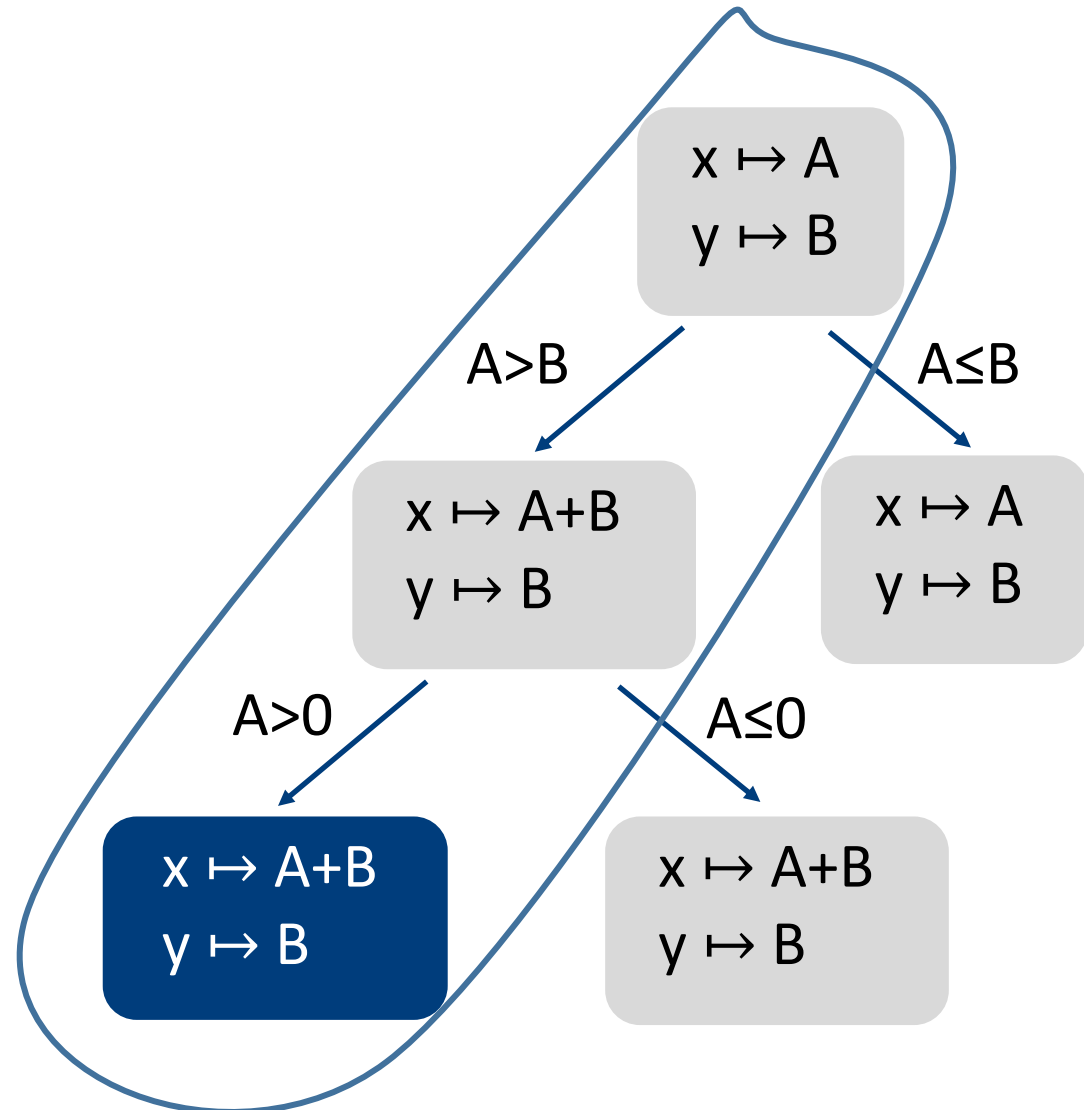
TRILION
Dynamic Binary Analysis

S²E

Recap: Symbolic Execution (SE)

```
1 def f (x, y):  
2   if (x>y)  
3     x = x+y  
4     if (x-y > 0)  
5       assert false  
6   return (x, y)
```

Dynamic Symbolic Execution (DSE):
A widely used variation of SE

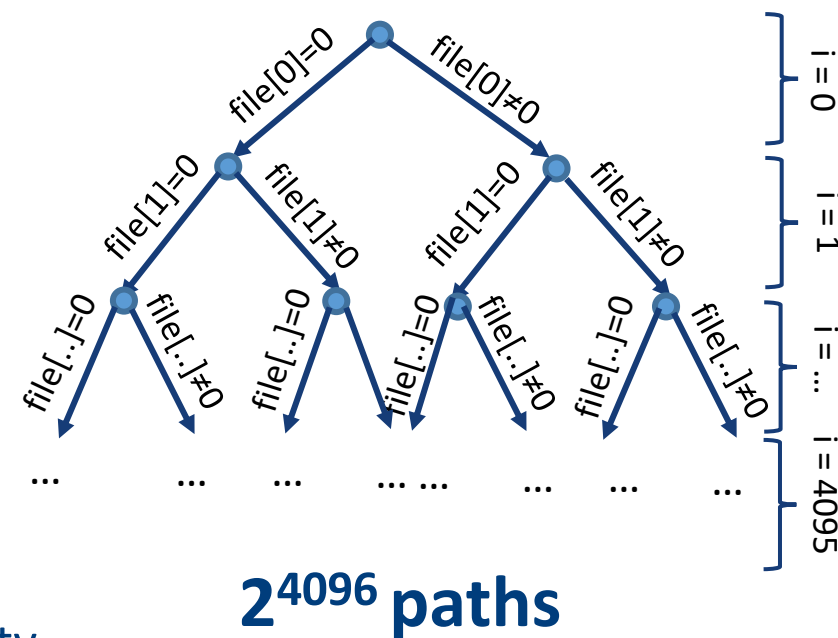


A and B are symbolic variables

Fundamental Limitations of Classic DSE

```
1 int main (...) {
2   if (strlen(filename)>1 && filename[0]=='-')
3     exit(1)
4   copy_data(...);
5   ...
6 }
7 void copy_data(..., int *file,...) {
8   static double data[4096], value;
9   read_double_value(file, ... );
10  value = fabs (data [0]);
11  for(i=0; i<4096; i++)
12    if(file[i] == 0.0) count++;
13  data[1] /= (value+count-3);
14  ...
15 }
```

- #1 Limitations of SMT Solvers
 - #2 Unmodeled Semantics
 - #3 Path Explosion
- Candidate Vulnerability Point (CVP): Divide-by-zero



Contributions

- **Neuro-symbolic execution**
 - A new approach to tackle the limitations of DSE
 - Reasons about exact (symbolic) & approximate constraints (neural nets)
- **A Tool – NeuEx**
 - Enhances the widely used DSE engine (KLEE)
- **Evaluation**
 - Finds 94% more bugs than KLEE in 12 hours

Problem

Inputs:

1. Source code
2. Symbolic Variables
(e.g., filename & file)
3. Candidate Vulnerability

Points (CVPs)

- Divide by zero
- Buffer overflow

Outputs:

Validated Exploits

```
1 int main (...) {
2     if (strlen(filename) > 1 && filename[0] == '-')
3         exit(1)
4     copy_data(...);
5     ...
6 }
7 void copy_data(..., int *file, ...) {
8     static double data[4096], value;
9     read_double_value(file, ... );
10    value = fabs (data [0]);
11    for(i=0; i<4096; i++)
12        if(file[i] == 0.0) count++;
13    data[1] /= (value+count-3); CVP: Divide-by-zero
14    ...
15 }
```

Key Insights

Values of Symbolic
Variables



Values of Vulnerable
Variables in CVP

Learn an approximation
with small number of I/O
examples

```
1 int main (...) {
2   if (strlen filename >1 && filename[0]=='-')
3     exit(1)
4   copy_data(...);
5   ...
6 }
7 void copy_data(..., int *file ...) {
8   [redacted]
9   [redacted]
10  [redacted]
11  [redacted]
12  [redacted]
13  data[1] /= (value+count-3); CVP: Divide-by-zero
14  ...
15 }
```

Key Insights

Approximation:

$$\begin{aligned} & \text{count} == \\ & \sum_{i \in [0, 4095]} \text{sign}(\text{file}[i] == 0) \\ & \wedge a == \text{file}[0] + 256 \text{file}[1] \\ & \wedge s == \text{sign}(\text{file}[1] < 127) \\ & \wedge \text{max} == \\ & (2 \times s - 1) \times a - 256^2 \times (s - 1) \end{aligned}$$

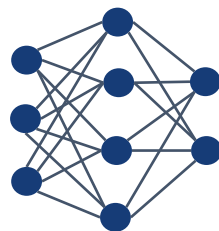
Machine learning can learn such an approximation

```
1 int main (...) {
2   if (strlen filename > 1 && filename[0] == '-')
3     exit(1)
4   copy_data(...);
5   ...
6 }
7 void copy_data(..., int *file, ...) {
8   static double data[4096], value;
9   read_double_value(file, ... );
10  value = fabs (data [0]);
11  for(i=0; i<4096; i++)
12    if(file[i] == 0.0) count++;
13  data[1] /= (value+count-3); CVP: Divide-by-zero
14  ...
15 }
```

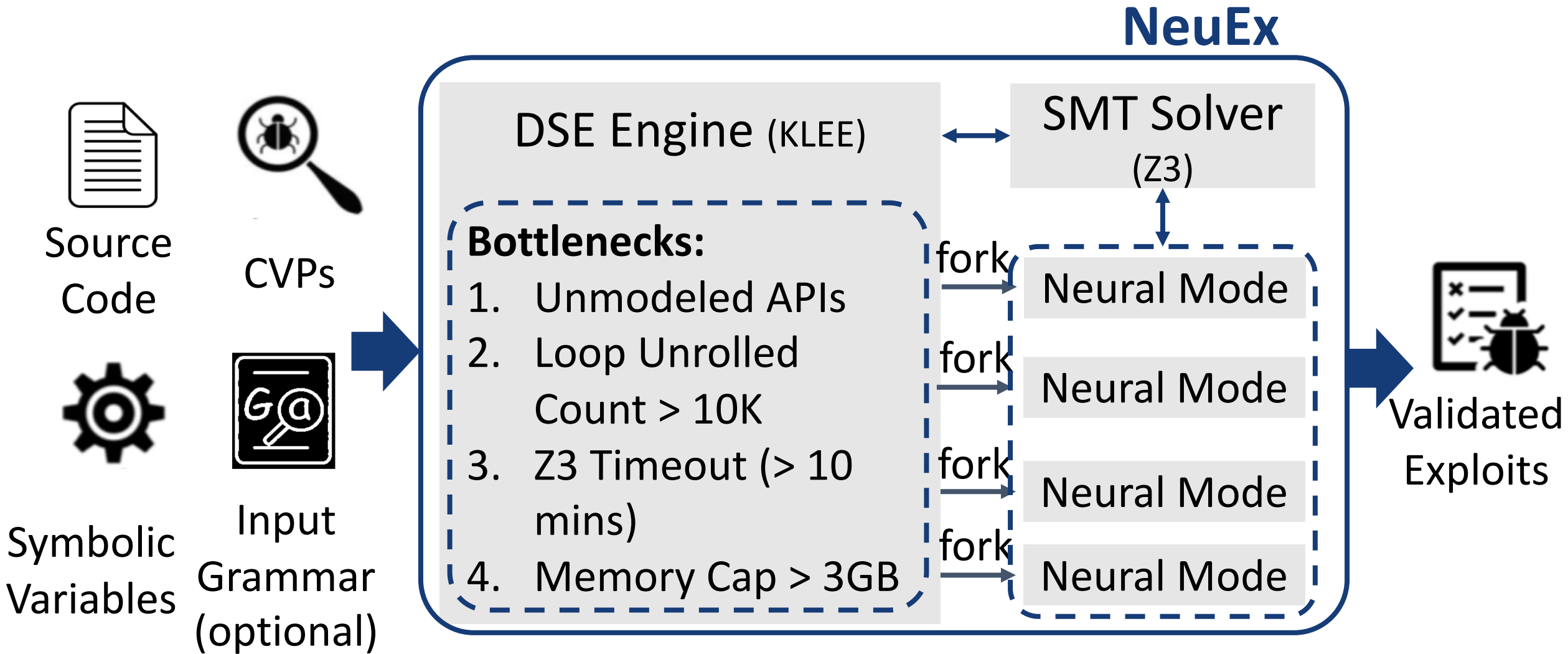

Approach

1. Neural nets can represent a large category of functions (universal approximation theorem).

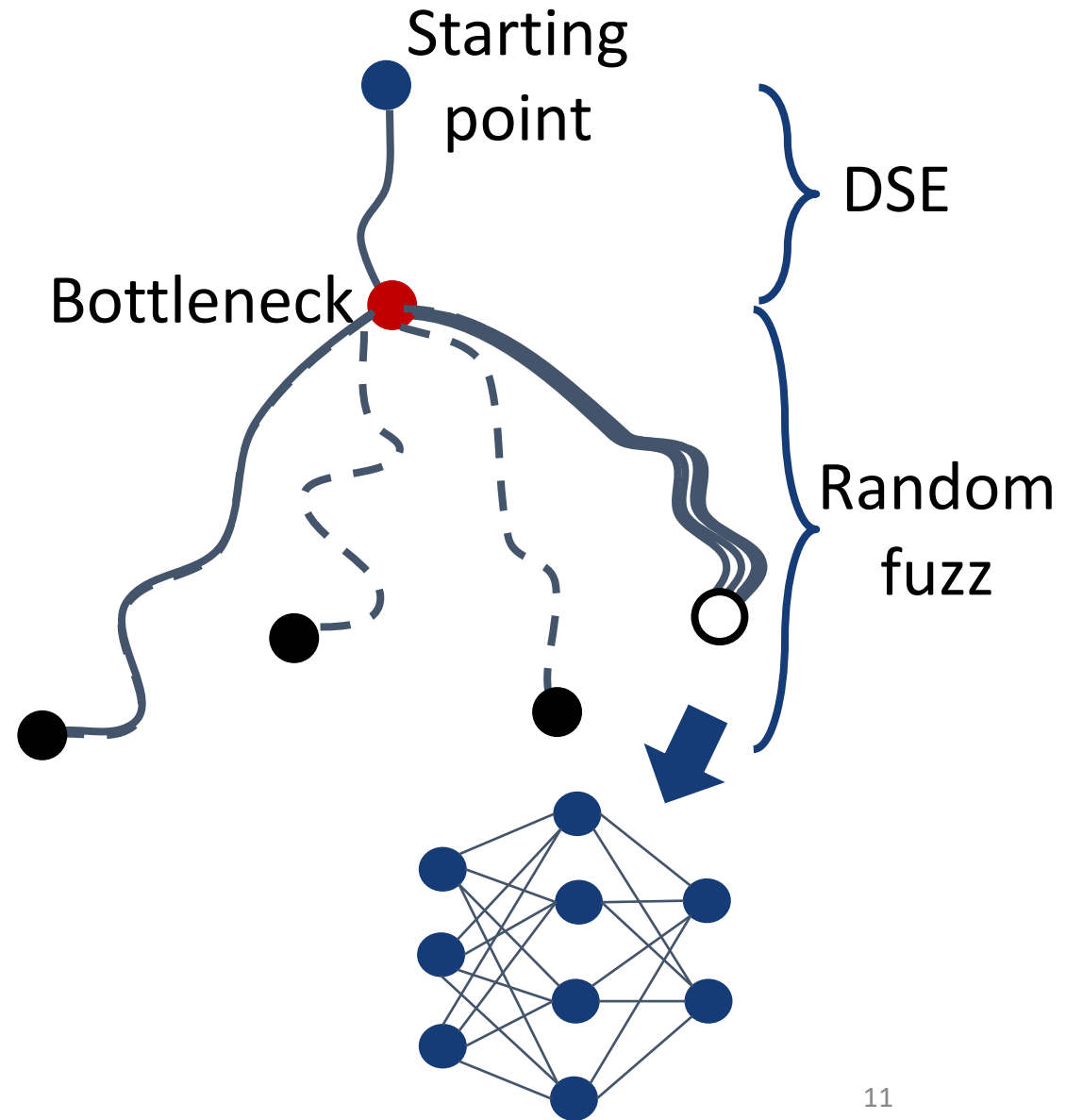
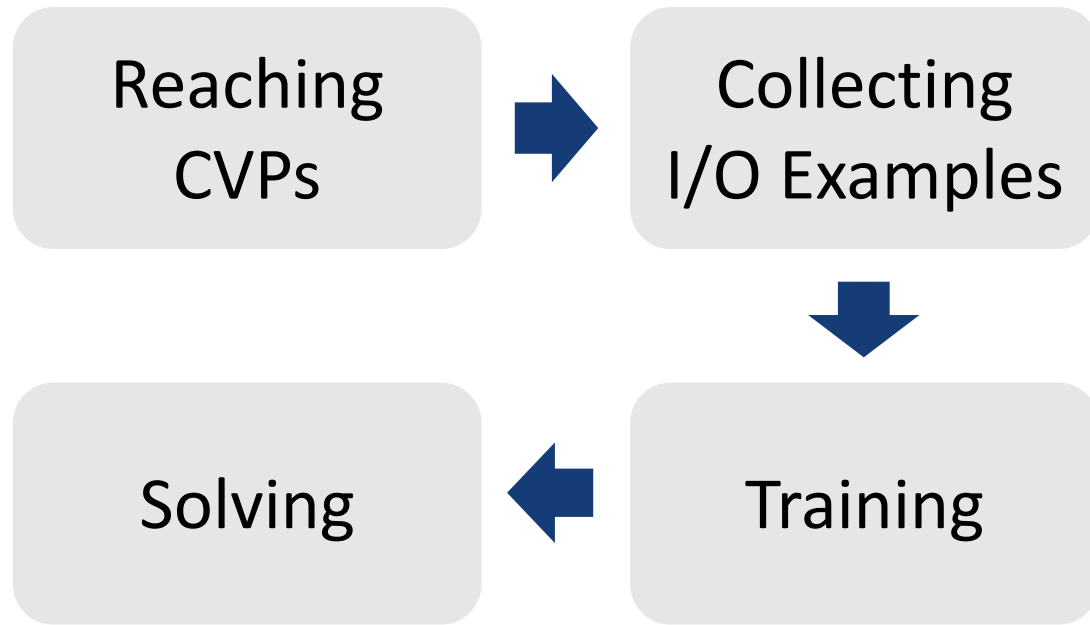
2. Multiple applications show that neural nets are learnable for many practical functions.

```
1 int main (...) {
2   if (strlen filename >1 && filename[0]=='-')
3     exit(1)
4   copy_data(...);
5   ...
6 }
7 void copy_data(..., int *file ..) {
8   Approximate Constraint (as a neural net):
9
10  file →  → count & value
11
12  data[1] /= (value+count-3); CVP: Divide-by-zero
13
14  ...
15 }
```

NeuEx Overview



Neural Mode



Generated Constraints

1. Reachability constraints:

$strlen(filename) \leq 1$
 $\forall filename \neq '-'$

\wedge

$N: infile \rightarrow (value, count)$

2. Vulnerability condition:

$value + count - 3 == 0$

```
1 int main (...) {
2   if (strlen(filename)>1 && filename[0]=='-')
3     exit(1)
4   copy_data(...);
5   ...
6 }
7 void copy_data(..., int *file,...) {
8   static double data[4096], value;
9   read_double_value(file, ... );
10  value = fabs (data [0]);
11  for(i=0; i<4096; i++)
12    if(file[i] == 0.0) count++;
13  data[1] /= (value+count-3); CVP: Divide-by-zero
14  ...
15 }
```

Constraint Solving

1. Reachability constraints:

$$\begin{aligned} & \text{strlen}(\text{filename}) \leq 1 \\ & \forall \text{filename} \neq \text{'-'}' \end{aligned}$$

\wedge

2. Vulnerability condition:

$$\text{value} + \text{count} - 3 == 0$$

Purely symbolic constraints:

➔ SMT solver

No variable shared with neural constraints

Mixed constraints:

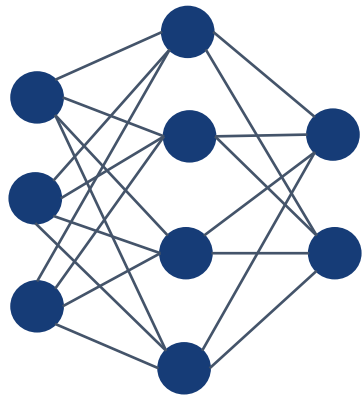
Including both neural constraints and symbolic constraints with shared variables



How to Solve Mixed Constraints?

Design Choice 1:

Neural net \rightarrow CNF



$$(X_1 \wedge Y_1) \vee (X_2 \wedge Y_2) \vee \dots \vee (X_n \wedge Y_n)$$

The number of clauses increases drastically with the neural net complexity

Design Choice 2:

Symbolic constraints \rightarrow Optimization objective of the neural net

Encoding Symbolic Constraints as an Optimization Objective

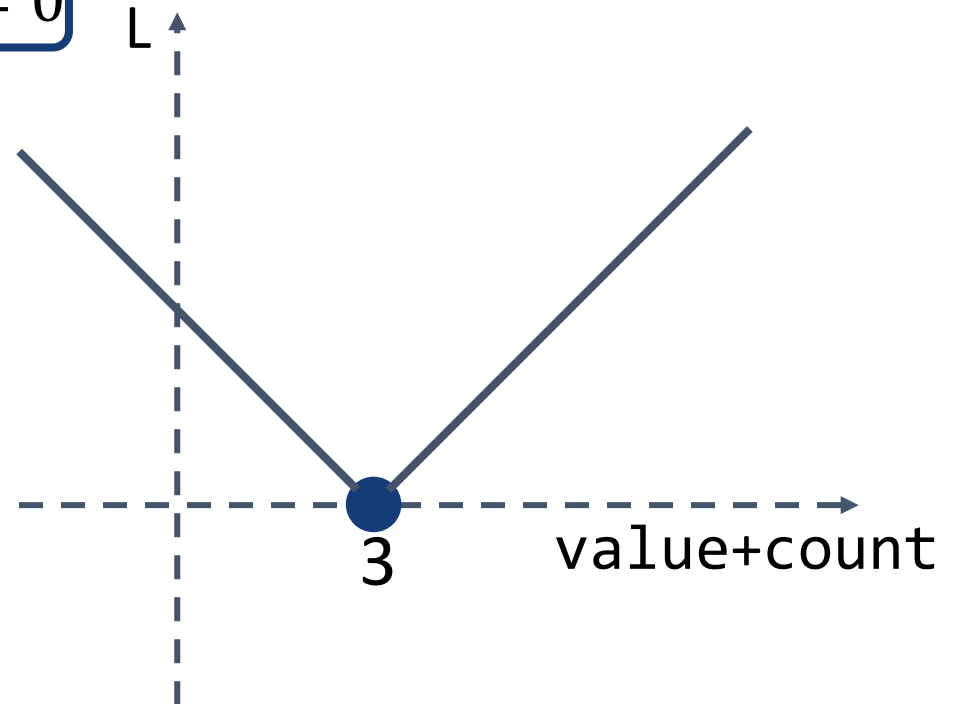
$N: infile \rightarrow (value, count) \wedge \boxed{value + count - 3 == 0}$
Symbolic constraint

Criterion for crafting the loss function:
The minimum point of the loss function satisfies the symbolic constraints.



One possible encoding:

$$L = \text{abs}(value + count - 3)$$



Encoding Symbolic Constraints as an Optimization Objective

NeuEx supports many symbolic constraints. Checkout the paper for the complete grammar and the corresponding loss functions.

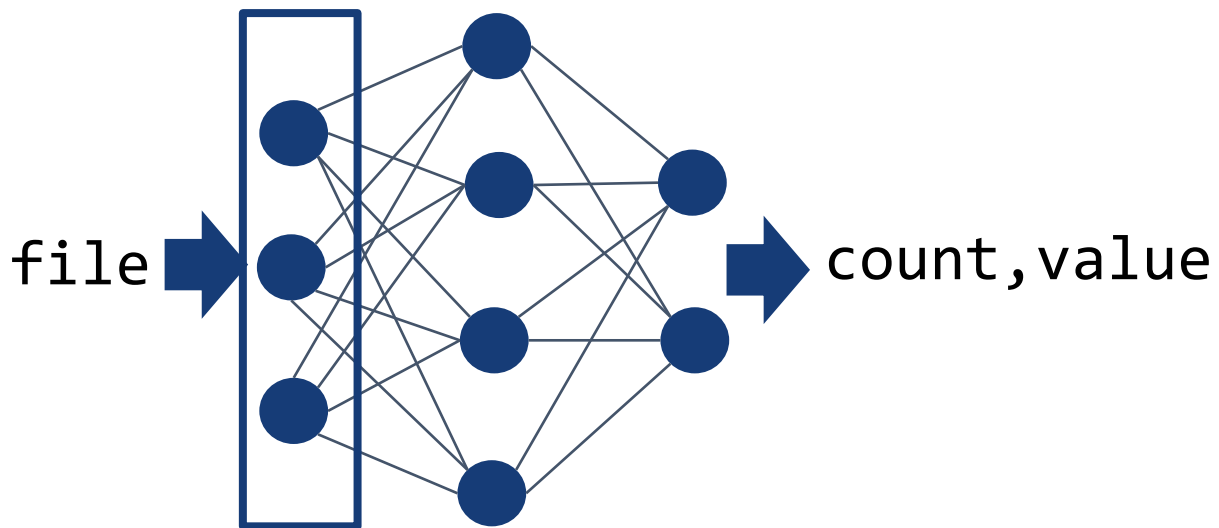
Neuro-Symbolic Execution: Augmenting Symbolic Execution with Neural Constraints

Shen Shiqi Shweta Shinde Soundarya Ramesh Abhik Roychoudhury Prateek Saxena
*Computer Science Department, School of Computing
National University of Singapore*
{shiqi04, shweta24, sramesh, abhik, prateeks}@comp.nus.edu.sg

Symbolic Constraint	Loss Function (L)
$S_1 ::= a < b$	$L = \max(a - b + \alpha, 0)$
$S_1 ::= a > b$	$L = \max(b - a + \alpha, 0)$
$S_1 ::= a \leq b$	$L = \max(a - b, 0)$
$S_1 ::= a \geq b$	$L = \max(b - a, 0)$
$S_1 ::= a = b$	$L = \text{abs}(a - b)$
$S_1 ::= a \neq b$	$L = \max(-1, -\text{abs}(a - b + \beta))$
$S_1 \wedge S_2$	$L = L_{S_1} + L_{S_2}$
$S_1 \vee S_2$	$L = \min(L_{S_1}, L_{S_2})$

Solving Mixed Constraints via Gradient Descent

Gradient: $\nabla_{file} L$



count	value	loss
257	20	274
1	20	18
0	20	17
...
0	3	0

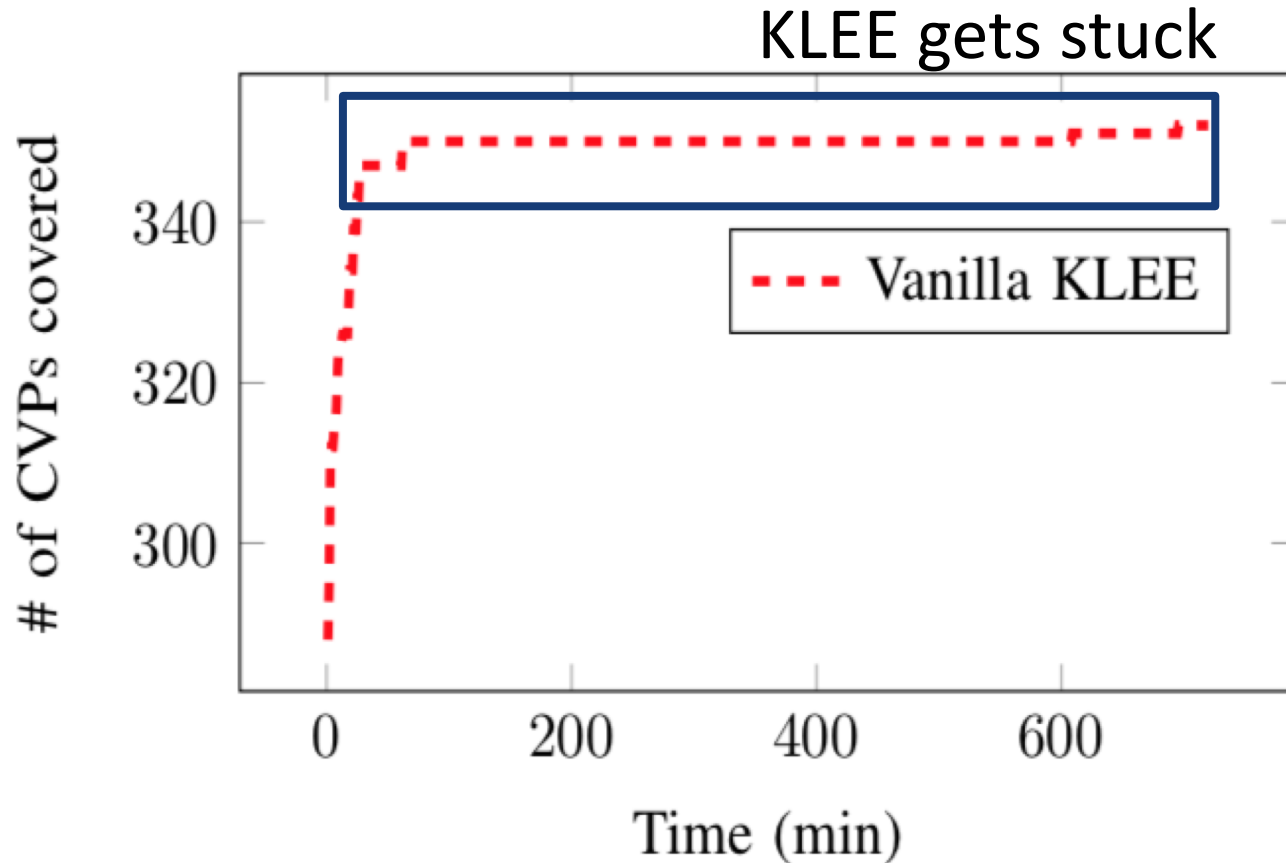
Concretely validate the exploit

file: 000...

Evaluation

- **Recall:** Neural mode is only triggered when DSE encounters bottlenecks
 - **Benchmarks:** 7 Programs known to be difficult for classic DSE
 - 4 Real programs
 - cURL: Data transferring
 - SQLite: Database
 - libTIFF: Image processing
 - libsndfile: Audio processing
 - LESE benchmarks
 - BIND, Sendmail, and WuFTP
- Include:
1. Complex loops
 2. Floating-point variables
 3. Unmodeled APIs

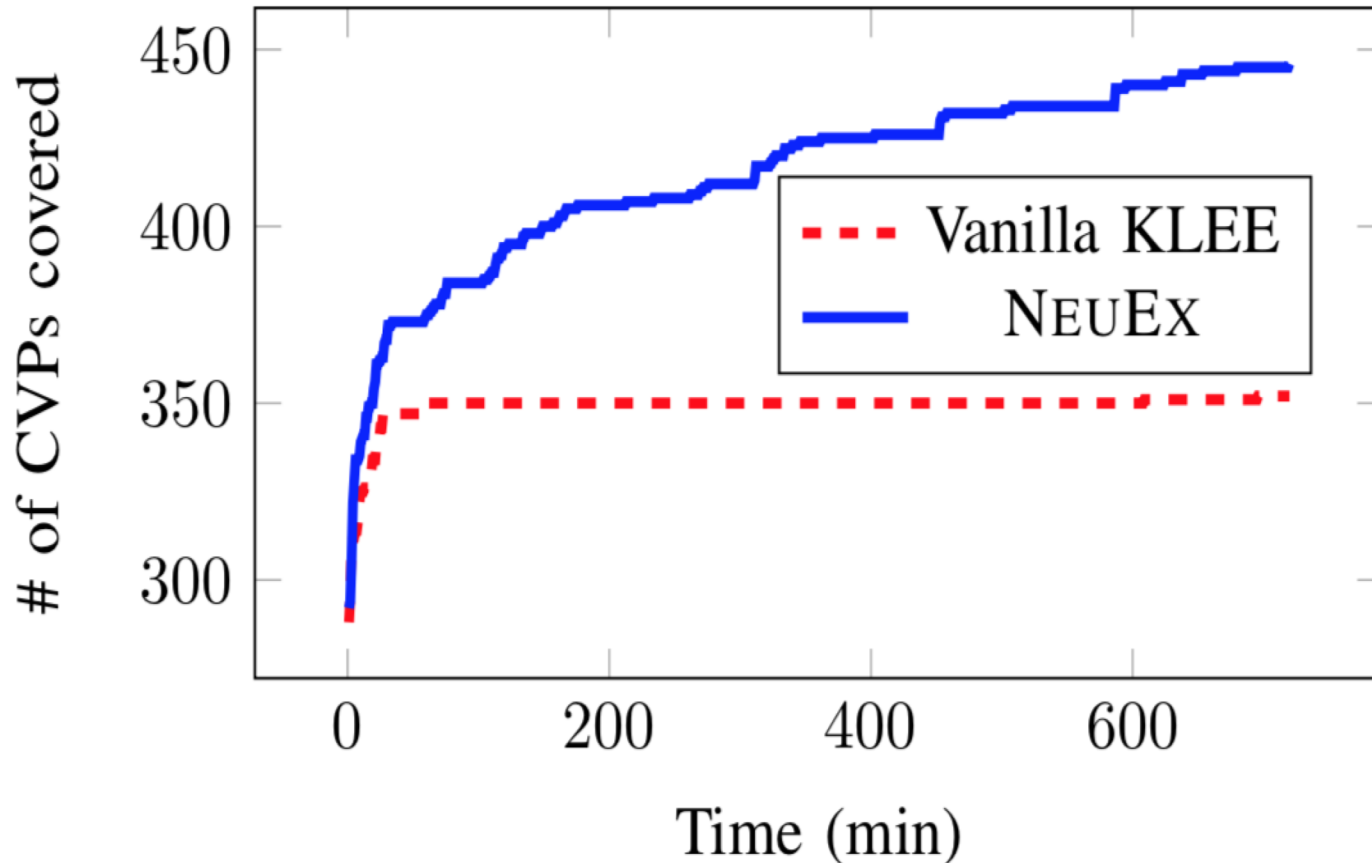
CVP Coverage & Bottlenecks for DSE



of bottlenecks: 61

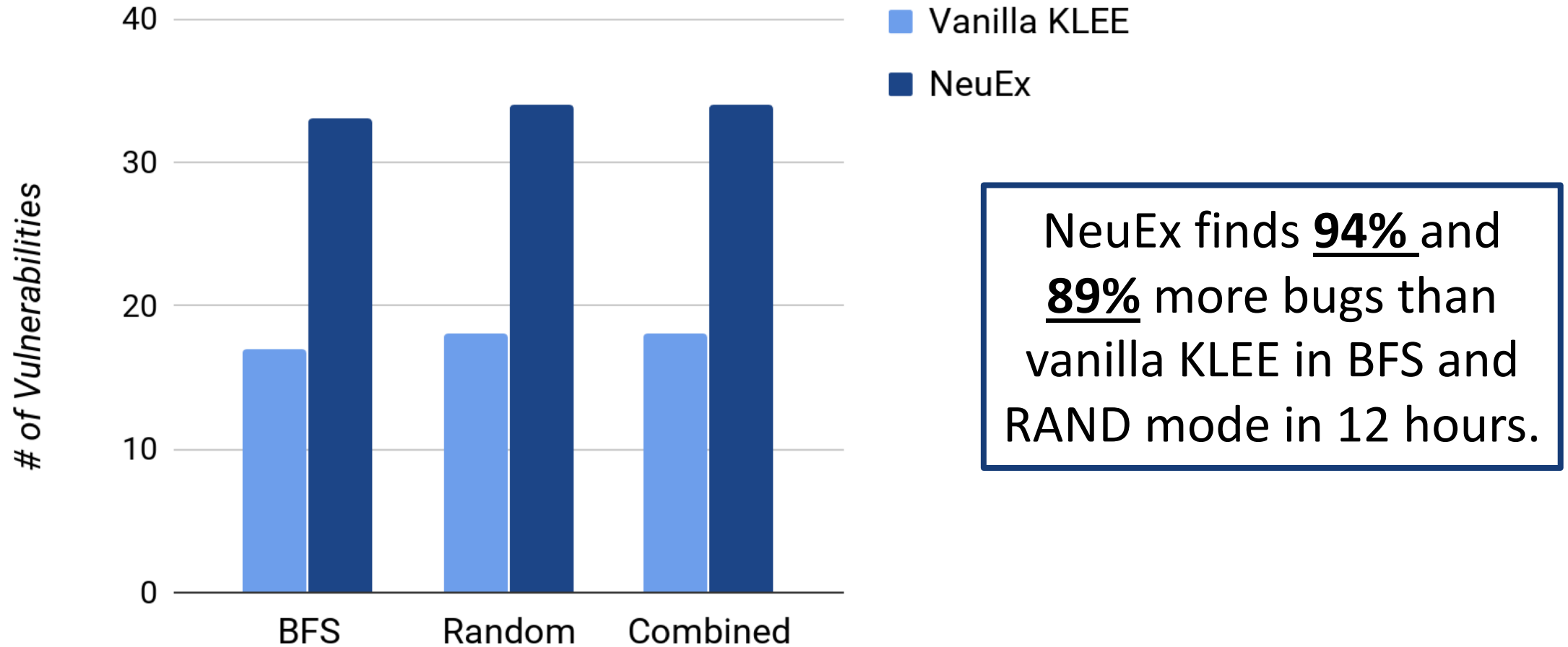
- Unmodeled APIs (6)
- Complicated loops (53)
- Z3 timeout (1)
- Memory exhaustion (1)

CVP Coverage of NeuEx vs KLEE



The number of CVPs reached or covered by NeuEx is **25%** higher than vanilla KLEE.

of Bugs Found by NeuEx vs KLEE



Comparison to DSE Extensions

LESE

- Structured Approach for Loop Reasoning
 - **Two orders** of magnitude slower on average

Veritesting

- Combination of DSE + static analysis
 - **Fails to find bugs** in 12 hours
 - Poor instruction coverage

Key Takeaways

A new approach: Neuro-Symbolic Execution

- Resolves fundamental bottlenecks of DSE
- First to learn unstructured representation not amenable to SMT solver
- Unique use of optimization techniques and SMT for solving constraints

Evaluation

- Finds **94%** more bugs than KLEE within 12 hours

Backup: New Bugs

NeuEx: 12

