

SymSan: Time and Space Efficient Concolic Execution via Dynamic Data-flow Analysis

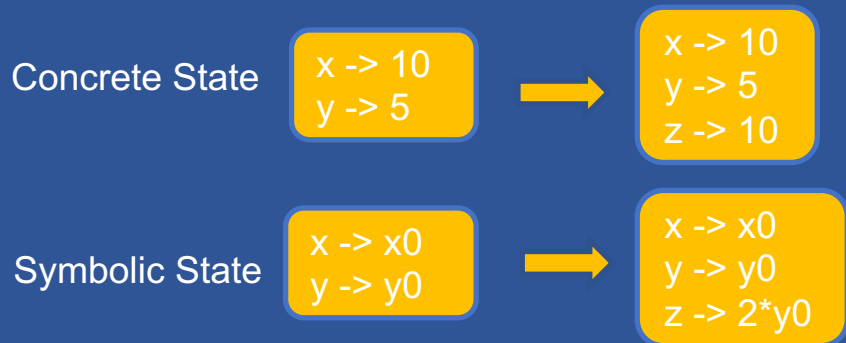
Ju Chen, Wookhyun Han, Mingjun Yin, Haochen Zeng, Chengyu Song, Byoungyoung Lee, Heng Yin, Insik Shin

31st USENIX SECURITY SYMPOSIUM
AUGUST 11, 2022



Concolic Execution: CONC(rete) + Symb(OLIC) execution

```
void foo(int x, int y) {  
  → z = 2 * y;  
  if (z == x) 2*y0 != x0 x: 10 y: 6 Testcase #1  
    if (x > y + 10) 2*y0 == x0 && x0 > y0 + 5 x: 20 y: 10 Testcase #2  
    assert(0);  
}
```



Overhead

$$z = x \oplus y$$

Parsing: Interpretation, slow (Angr, KLEE)
Instrumentation, faster! (QSYM, SymQEMU, SymCC, SymSan)

Locating: Reading $\text{Sym}(x)$ and $\text{Sym}(b)$ from symbolic state

Creating/
updating: Creating $\text{Sym}(z)$, updating symbolic state

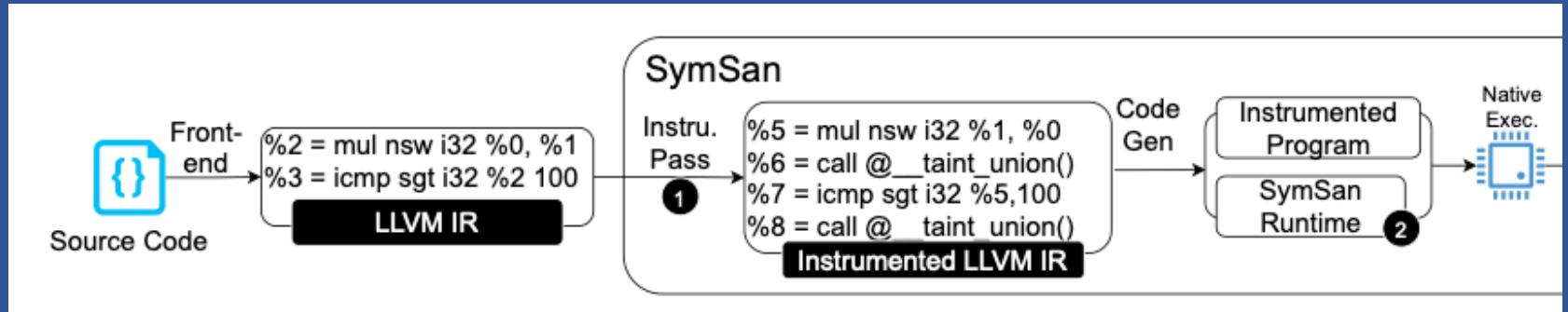
Locating/creating/updating have non-trivial
overhead

Insight

- Concolic execution is a special form of dynamic data-flow analysis, so...
- it can be simply implemented on top of LLVM DFSan (highly-optimized)

Use SymSan

CC=symsan CXX=symsan make



Concolic Execution is forward data-flow analysis

	Generic Data-flow analysis	Concolic Execution
Label Interpretation	Variables' properties	Variables' symbolic expressions
Label Introduction	How labels are introduced	Program inputs
Label Propagation	How labels are updated after executing an instruction	Compute symbolic expressions
Label Sinks	Where and how properties are used	Conditional branches. Update constraints

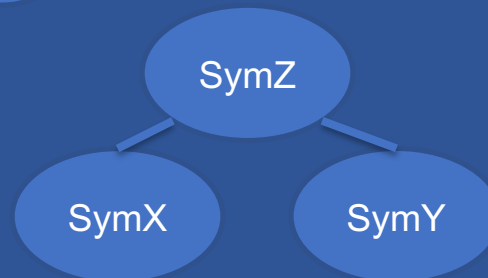
Var → Label (symbolic expression)

A label is (the index of) a symbolic expression

```
struct dfsan_label_info {  
    dfsan_label l1; // symbolic sub-expression  
    dfsan_label l2; // symbolic sub-expression  
    u64 op1; // concrete operand  
    u64 op2; // concrete operand  
    u16 op; // opcode, using LLVM IR operations  
    u16 size; // size of the result  
};
```

...	Lbl101	Lbl102	Lbl103	...
...	SymX	SymY	SymZ	...

Union Table



Running example

$$z = x + y$$

Shadow Mem

x	101
y	102

Before Exec.

Union Table (AST Table)

...	Lbl101	Lbl102	Lbl103	...
...	SymX	SymY	UNINIT	...

x	101
y	102
z	103

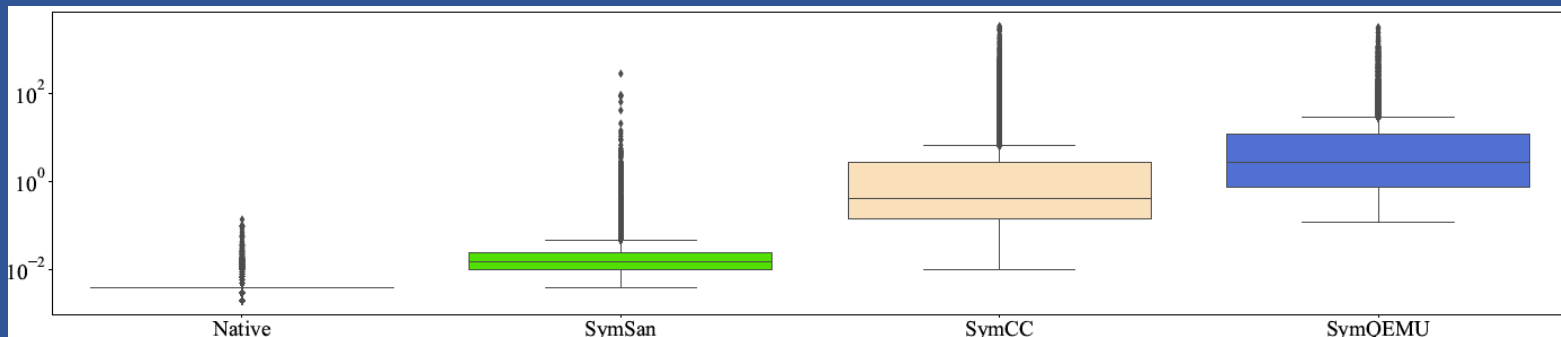
After Exec.

...	Lbl101	Lbl102	Lbl103	...
...	SymX	SymY	L1:101 L2:102 OP: ADD	...

Why is SymSan faster?

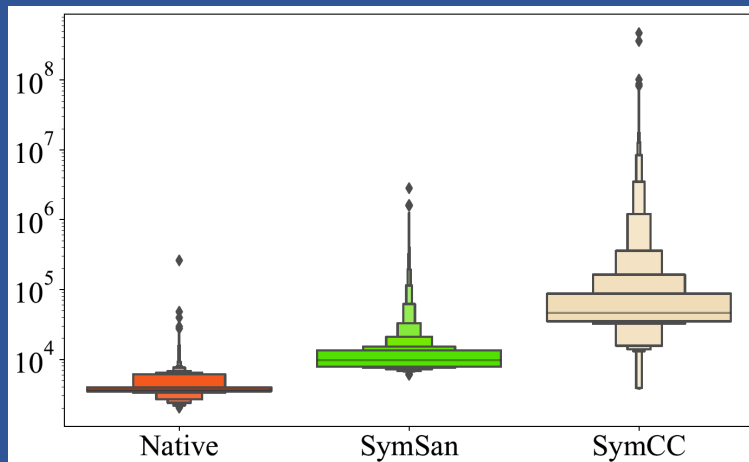
	SymCC	SymSan
Symbolic expression	std::shared_ptr	A 32-bit integer
Shadow memory access	std::map, $O(\log n)$	Direct mapping, $O(1)$
Symbolic expression allocation	new/malloc()	atomic_fetch_inc()
Arguments/Return value passing	std::array, multiple function calls	TLS, single MOV instruction

Efficiency (collecting constraints, no solving)



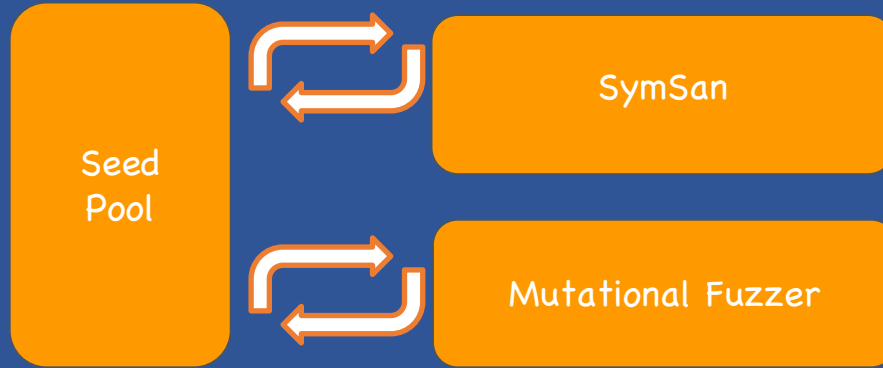
- Two orders of magnitude faster in constraints collecting
- Solves 2x more constraints given the same time budget

Memory Efficiency (no solving)



- Consumes one order of magnitude less memory
- Good for async-solving

Hybrid-fuzzing



- Fuzzbench (SymSan ranked #1 in average score)

Takeaways

- An efficient concolic executor, built on top of DFSan
- Doing source-based concolic execution? Try SymSan
 - near-optimal performance!
- World's fastest 😊 if pairing with JIGSAW (Oakland'22)

Thank you for listening!



<https://github.com/R-Fuzz/symsan>

