

IA159 Formal Methods for Software Analysis

Symbolic Execution and Applications

Jan Strejček

Faculty of Informatics
Masaryk University

focus

- symbolic execution
- automated whitebox fuzz testing

sources

- J. C. King: *Symbolic Execution and Program Testing*, Communications of ACM, 1976.
- P. Godefroid, M. Y. Levin, and D. Molnar: *Automated whitebox fuzz testing*, NDSS 2008.

Motivation

```
1 int sum(int a, int b, int c) {  
2   int x = a + b;  
3   int y = b + c;  
4   int z = x + y - b;  
5   return z;  
6 }
```

testing checks that the program behaves correctly on selected inputs

- `sum(1, 1, 1)` returns 3
- `sum(1, 2, 3)` returns 6
- ...

Motivation

```
1 int sum(int a, int b, int c) {  
2   int x = a + b;  
3   int y = b + c;  
4   int z = x + y - b;  
5   return z;  
6 }
```

testing checks that the program behaves correctly on selected inputs

- `sum(1, 1, 1)` returns 3
- `sum(1, 2, 3)` returns 6
- ...

Motivation

```
1 int sum(int a, int b, int c) {  
2   int x = a + b;  
3   int y = b + c;  
4   int z = x + y - b;  
5   return z;  
6 }
```

we can execute the program with symbols $\alpha_1, \alpha_2, \alpha_3$ representing arbitrary input values

- `sum($\alpha_1, \alpha_2, \alpha_3$)` returns

Motivation

```
1 int sum(int a, int b, int c) {  
2   int x = a + b;  
3   int y = b + c;  
4   int z = x + y - b;  
5   return z;  
6 }
```

we can execute the program with symbols $\alpha_1, \alpha_2, \alpha_3$ representing arbitrary input values

■ `sum($\alpha_1, \alpha_2, \alpha_3$)` returns $\alpha_1 + \alpha_2 + \alpha_3$

(if `int` interpreted as \mathbb{Z})

→ symbolic execution

Symbolic execution semantics in general

each programming language has an execution semantics describing

- which data objects the program manipulates
- how statements manipulate data objects
- how control flows through the statements of a program

in **symbolic execution semantics**

- real data objects are replaced by symbolic ones, which are typically expressions over **symbols** $\alpha_1, \alpha_2, \dots$ representing arbitrary input values
- the semantics of statements is extended to accept symbolic input and produce symbolic output
- control flow is handled differently as some branching conditions can be evaluated to both *true* and *false* depending on the values of symbols

assumptions and notation

- consider a program that handles only integer (\mathbb{Z}) variables and it is built from assignments and branching statements
- a special assignment $x = *$ (or $x = \text{input}()$) corresponds to reading input
- *Vars* = the set of variables in the considered program
- *Sym* = $\{\alpha, \beta, \dots, \alpha_1, \alpha_2, \dots\}$ = countable set of symbols representing arbitrary input values
- *Exp(Sym)* = expressions over *Sym*, integers, and arithmetic operations
- for example, $2\alpha + \beta^3 - 7 \in \text{Exp}(Sym)$
- *Exp(Sym)* are symbolic data objects

assumptions and notation

- consider a program that handles only integer (\mathbb{Z}) variables and it is built from assignments and branching statements
- a special assignment $x = *$ (or $x = \text{input}()$) corresponds to reading input
- *Vars* = the set of variables in the considered program
- *Sym* = $\{\alpha, \beta, \dots, \alpha_1, \alpha_2, \dots\}$ = countable set of symbols representing arbitrary input values
- *Exp(Sym)* = expressions over *Sym*, integers, and arithmetic operations
- for example, $2\alpha + \beta^3 - 7 \in \text{Exp}(Sym)$
- *Exp(Sym)* are symbolic data objects

symbolic execution computes symbolic states consisting of

- 1 current program location
- 2 symbolic memory
- 3 path condition

Symbolic memory m

- $m: Vars \rightarrow Exp(Sym)$
- assigns expressions of $Exp(Sym)$ to program variables
- a symbol $\alpha \in Sym$ is called **fresh** if it was not used before in the considered computation
- **initial symbolic memory** m_0 assigns to each variable $x \in Vars$ a fresh symbol $m(x) = \alpha \in Sym$
- symbolic memory is modified by assignments
- for any program expression exp (Boolean or integer), by $m(exp)$ we denote the expression where each program variable $x \in Vars$ is replaced by $m(x)$

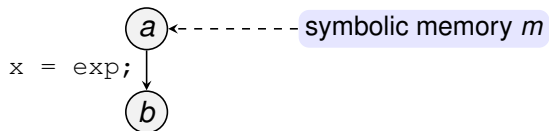
Symbolic memory m

- $m: Vars \rightarrow Exp(Sym)$
- assigns expressions of $Exp(Sym)$ to program variables
- a symbol $\alpha \in Sym$ is called **fresh** if it was not used before in the considered computation
- **initial symbolic memory** m_0 assigns to each variable $x \in Vars$ a fresh symbol $m(x) = \alpha \in Sym$
- symbolic memory is modified by assignments
- for any program expression exp (Boolean or integer), by $m(exp)$ we denote the expression where each program variable $x \in Vars$ is replaced by $m(x)$

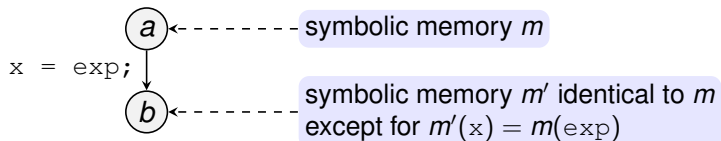
example

- let $m(x) = \alpha + 4$ and $m(y) = 2\alpha + \beta$
- $m(5x - y + 8) = 5(\alpha + 4) - (2\alpha + \beta) + 8 = 3\alpha - \beta + 28$
- $m(x \neq y) = (\alpha + 4 \neq 2\alpha + \beta) = (\alpha + \beta \neq 4)$

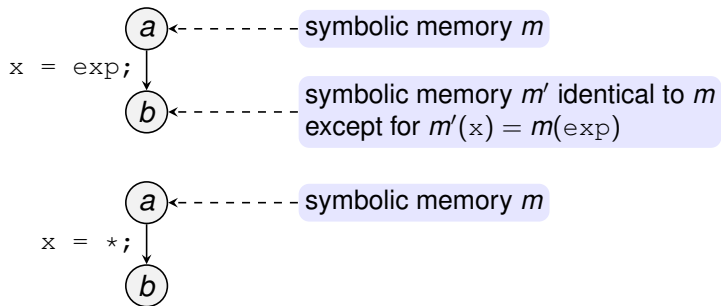
Symbolic execution of assignments



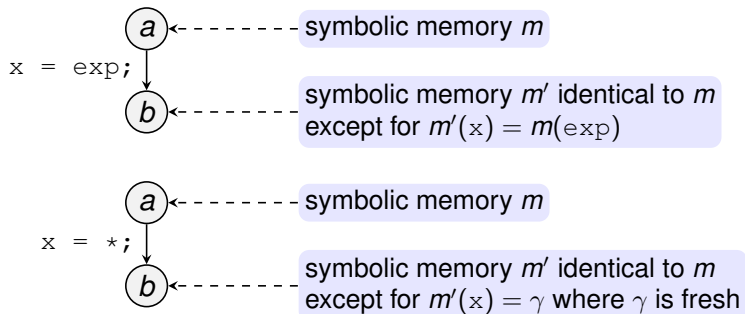
Symbolic execution of assignments



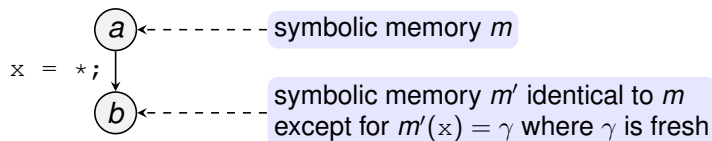
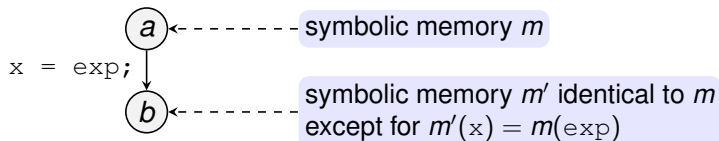
Symbolic execution of assignments



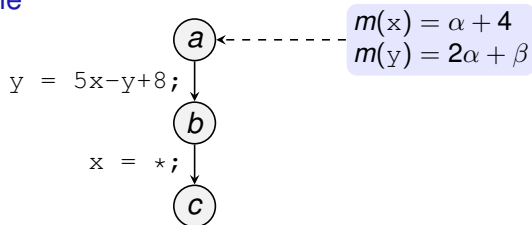
Symbolic execution of assignments



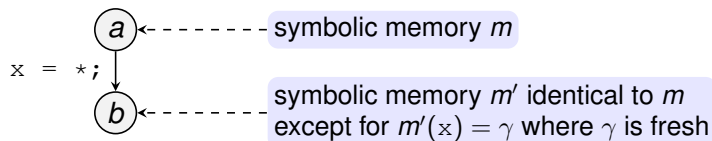
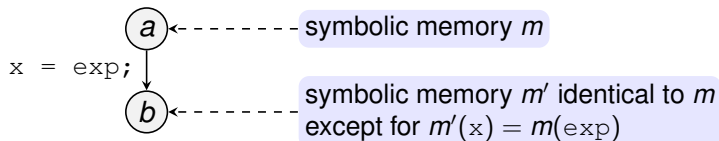
Symbolic execution of assignments



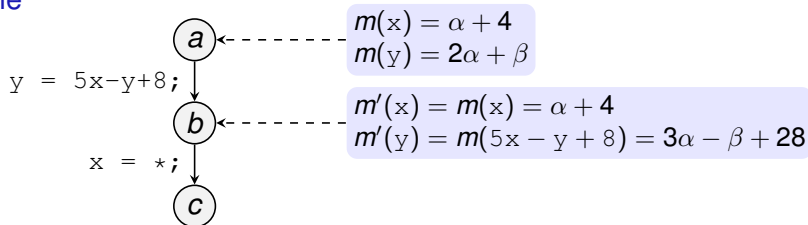
example



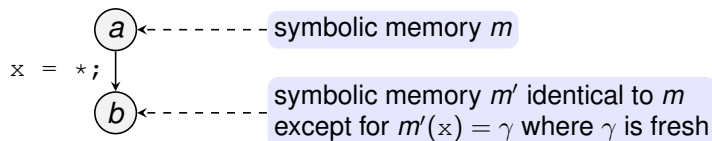
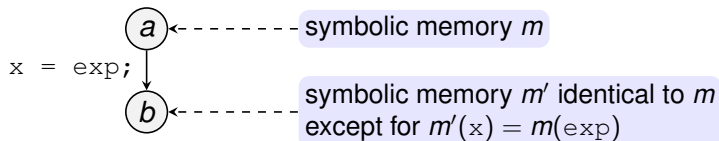
Symbolic execution of assignments



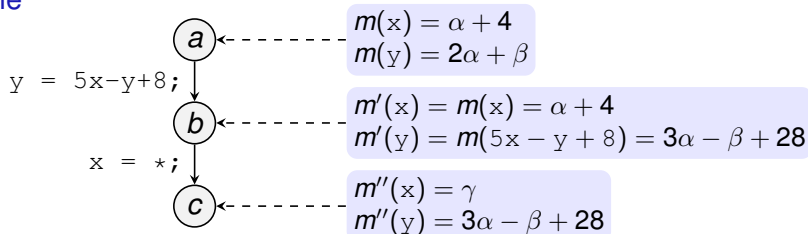
example



Symbolic execution of assignments



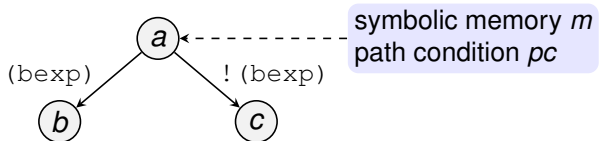
example



- a quantifier-free predicate formula over Sym corresponding to a program path
- pc is the necessary and sufficient condition on input values to navigate the program execution along the current path
- if pc is not satisfiable the corresponding path is unfeasible
- pc is initially set to *true*
- pc is modified by evaluation of branching statements

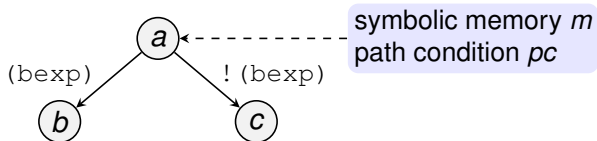
Symbolic execution of branching statements

```
if (bexp) {...} else {...}
```



Symbolic execution of branching statements

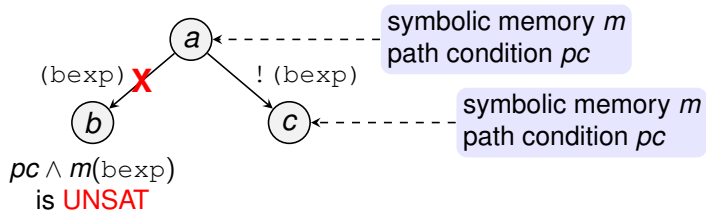
```
if (bexp) {...} else {...}
```



- 1 check feasibility of the *true* branch:
if $pc \wedge m(bexp)$ is not satisfiable, continue to the *false* branch

Symbolic execution of branching statements

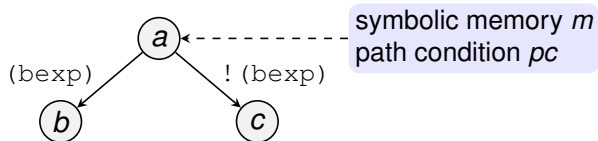
```
if (bexp) {...} else {...}
```



- 1 check feasibility of the *true* branch:
if $pc \wedge m(bexp)$ is not satisfiable, continue to the *false* branch

Symbolic execution of branching statements

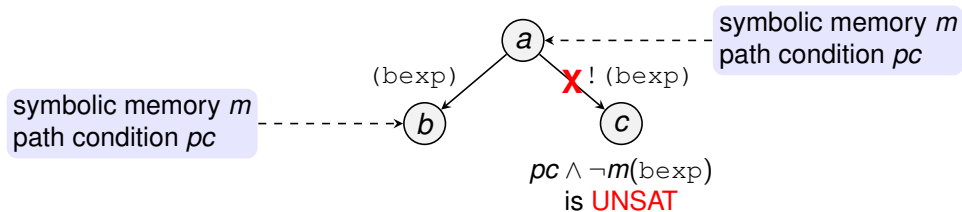
```
if (bexp) {...} else {...}
```



- 1 check feasibility of the *true* branch:
if $pc \wedge m(bexp)$ is not satisfiable, continue to the *false* branch
- 2 check feasibility of the *false* branch:
if $pc \wedge \neg m(bexp)$ is not satisfiable, continue to the *true* branch

Symbolic execution of branching statements

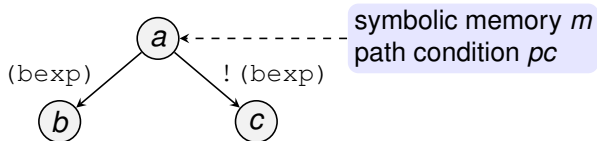
```
if (bexp) {...} else {...}
```



- 1 check feasibility of the *true* branch:
if $pc \wedge m(bexp)$ is not satisfiable, continue to the *false* branch
- 2 check feasibility of the *false* branch:
if $pc \wedge \neg m(bexp)$ is not satisfiable, continue to the *true* branch

Symbolic execution of branching statements

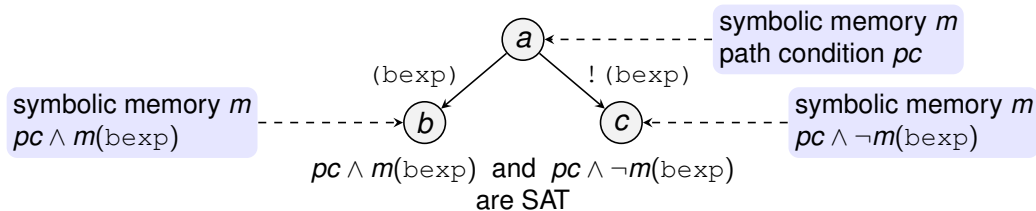
```
if (bexp) {...} else {...}
```



- 1 check feasibility of the *true* branch:
if $pc \wedge m(bexp)$ is not satisfiable, continue to the *false* branch
- 2 check feasibility of the *false* branch:
if $pc \wedge \neg m(bexp)$ is not satisfiable, continue to the *true* branch
- 3 if both are satisfiable, **fork** the symbolic execution
set pc in *true* branch to $pc \wedge m(bexp)$
set pc in *false* branch to $pc \wedge \neg m(bexp)$

Symbolic execution of branching statements

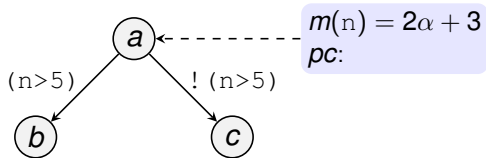
```
if (bexp) {...} else {...}
```



- 1 check feasibility of the *true* branch:
if $pc \wedge m(bexp)$ is not satisfiable, continue to the *false* branch
- 2 check feasibility of the *false* branch:
if $pc \wedge \neg m(bexp)$ is not satisfiable, continue to the *true* branch
- 3 if both are satisfiable, **fork** the symbolic execution
set pc in *true* branch to $pc \wedge m(bexp)$
set pc in *false* branch to $pc \wedge \neg m(bexp)$

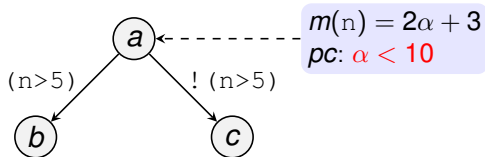
Symbolic execution of branching statements: example

```
if (n>5) {...} else {...}
```



Symbolic execution of branching statements: example

```
if (n>5) {...} else {...}
```

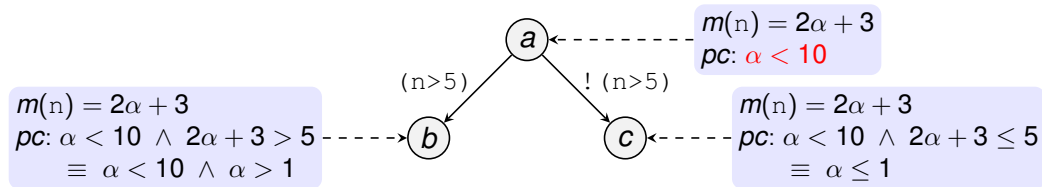


pc is $\alpha < 10$

- *true* branch is feasible as $\alpha < 10 \wedge m(n > 5) \equiv \alpha < 10 \wedge 2\alpha + 3 > 5$ is satisfiable (e.g. by $\alpha = 3$)
- *false* branch is feasible as $\alpha < 10 \wedge \neg m(n > 5) \equiv \alpha < 10 \wedge 2\alpha + 3 \leq 5$ is satisfiable (e.g. by $\alpha = 0$)
- **fork** execution and update pc on both branches

Symbolic execution of branching statements: example

```
if (n>5) {...} else {...}
```

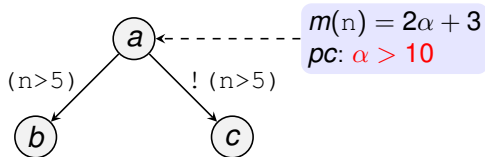


pc is $\alpha < 10$

- *true* branch is feasible as $\alpha < 10 \wedge m(n > 5) \equiv \alpha < 10 \wedge 2\alpha + 3 > 5$ is satisfiable (e.g. by $\alpha = 3$)
- *false* branch is feasible as $\alpha < 10 \wedge \neg m(n > 5) \equiv \alpha < 10 \wedge 2\alpha + 3 \leq 5$ is satisfiable (e.g. by $\alpha = 0$)
- **fork** execution and update pc on both branches

Symbolic execution of branching statements: example

```
if (n>5) {...} else {...}
```

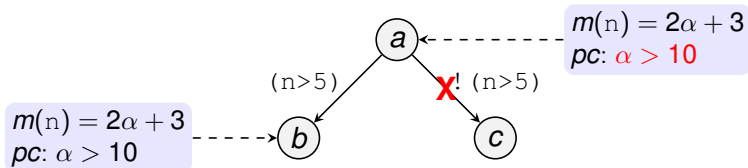


pc is $\alpha > 10$

- *false* branch is unfeasible as
 $\alpha > 10 \wedge \neg m(n > 5) \equiv \alpha > 10 \wedge 2\alpha + 3 \leq 5 \equiv \text{false}$
- continue to *true* branch with the same pc

Symbolic execution of branching statements: example

```
if (n>5) {...} else {...}
```

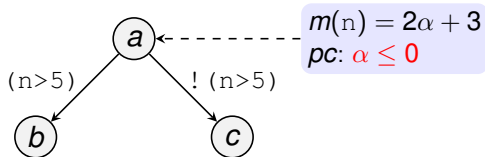


pc is $\alpha > 10$

- *false* branch is unfeasible as
$$\alpha > 10 \wedge \neg m(n > 5) \equiv \alpha > 10 \wedge 2\alpha + 3 \leq 5 \equiv \textit{false}$$
- continue to *true* branch with the same pc

Symbolic execution of branching statements: example

```
if (n>5) {...} else {...}
```

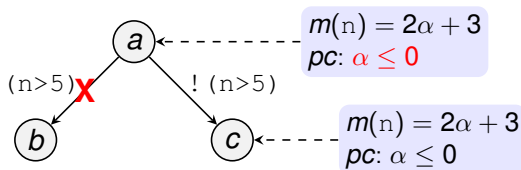


pc is $\alpha \leq 0$

- *true* branch is unfeasible as
 $\alpha \leq 0 \wedge m(n > 5) \equiv \alpha \leq 0 \wedge 2\alpha + 3 > 5 \equiv \text{false}$
- continue to *false* branch with the same pc

Symbolic execution of branching statements: example

```
if (n>5) {...} else {...}
```



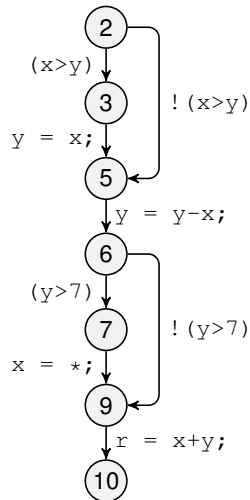
pc is $\alpha \leq 0$

- *true* branch is unfeasible as
 $\alpha \leq 0 \wedge m(n > 5) \equiv \alpha \leq 0 \wedge 2\alpha + 3 > 5 \equiv \text{false}$
- continue to *false* branch with the same pc

- **nodes** are states of symbolic execution, i.e., triples (l, m, pc) of program location l , symbolic memory m , and path condition pc
- **root** node $(l_0, m_0, true)$ consists of initial program location l_0 , initial symbolic memory m_0 assigning fresh symbols, and initial path condition $true$
- successors are computed by symbolic execution of the assignment or branching statement corresponding to the current location
- only locations with branching statement can have more successors

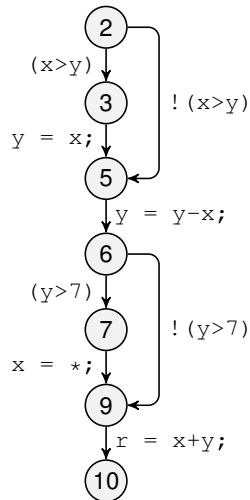
Symbolic execution tree: example 1

```
1 int foo(int x, int y) {  
2   if (x>y) {  
3     y = x;  
4   }  
5   y = y-x;  
6   if (y>7) {  
7     x = *;  
8   }  
9   return x+y;  
10 }
```

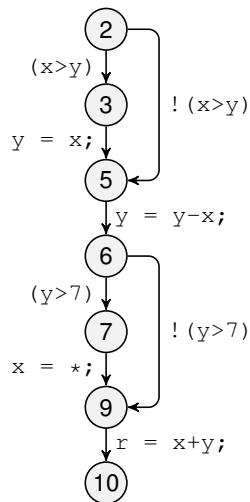
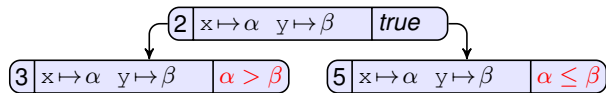


Symbolic execution tree: example 1

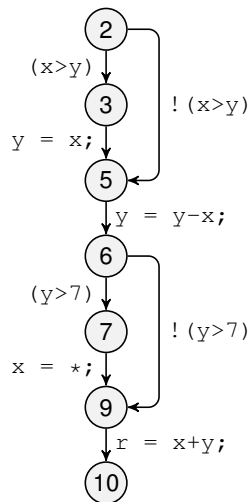
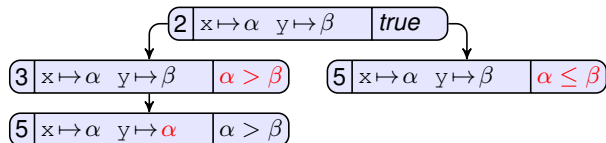
2	$x \mapsto \alpha \quad y \mapsto \beta$	<i>true</i>
---	--	-------------



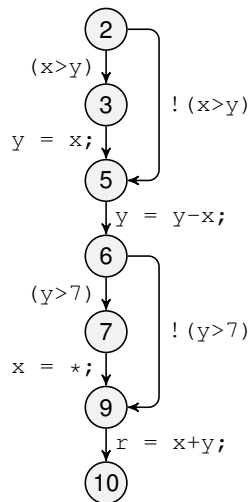
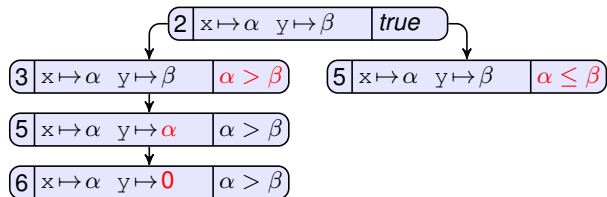
Symbolic execution tree: example 1



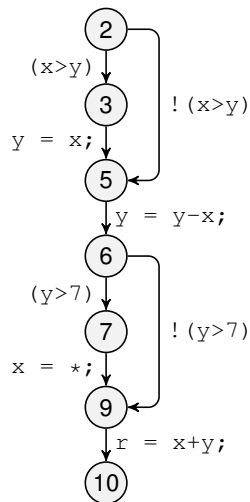
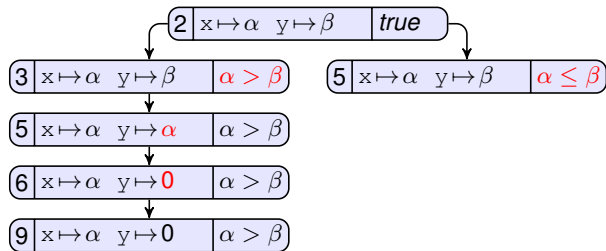
Symbolic execution tree: example 1



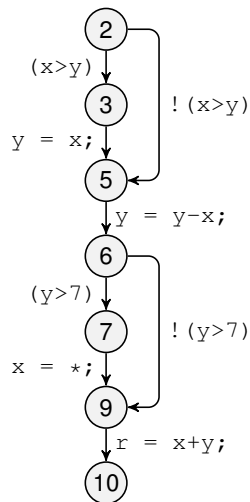
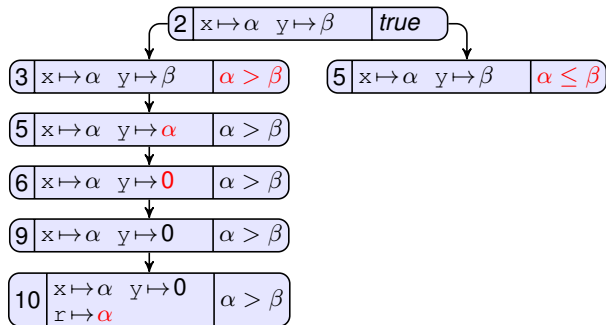
Symbolic execution tree: example 1



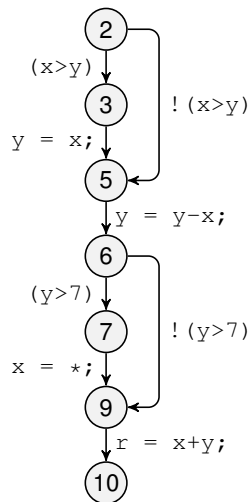
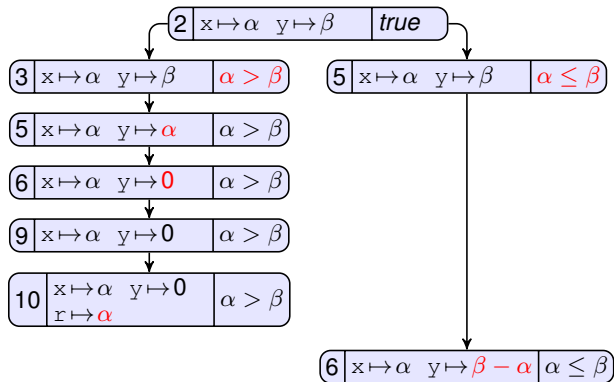
Symbolic execution tree: example 1



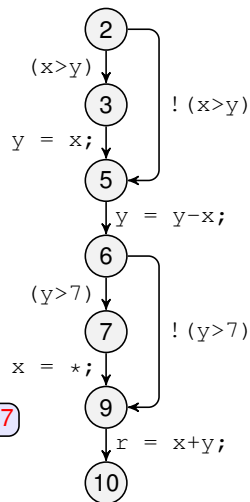
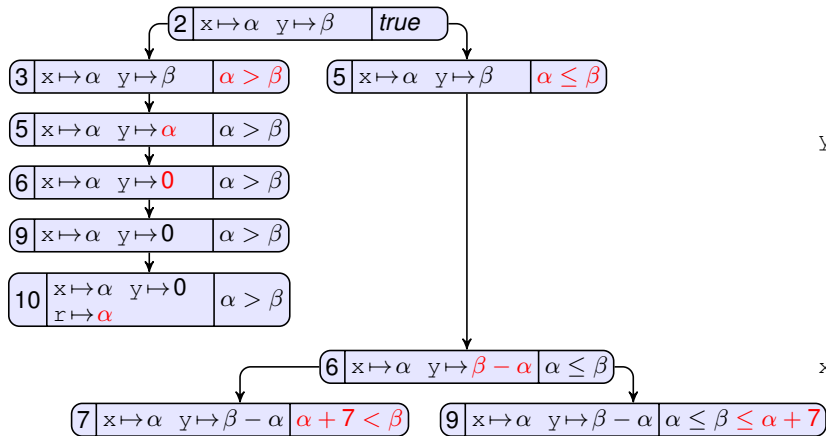
Symbolic execution tree: example 1



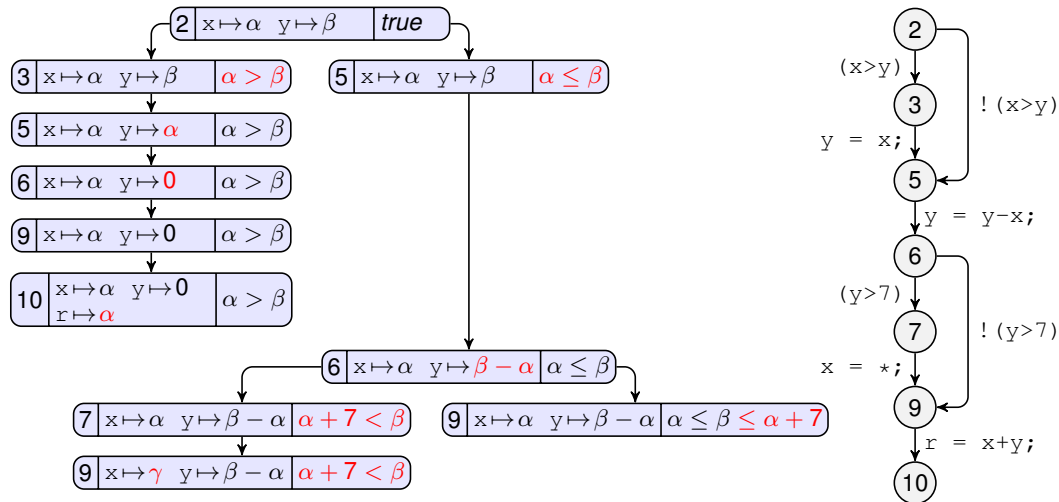
Symbolic execution tree: example 1



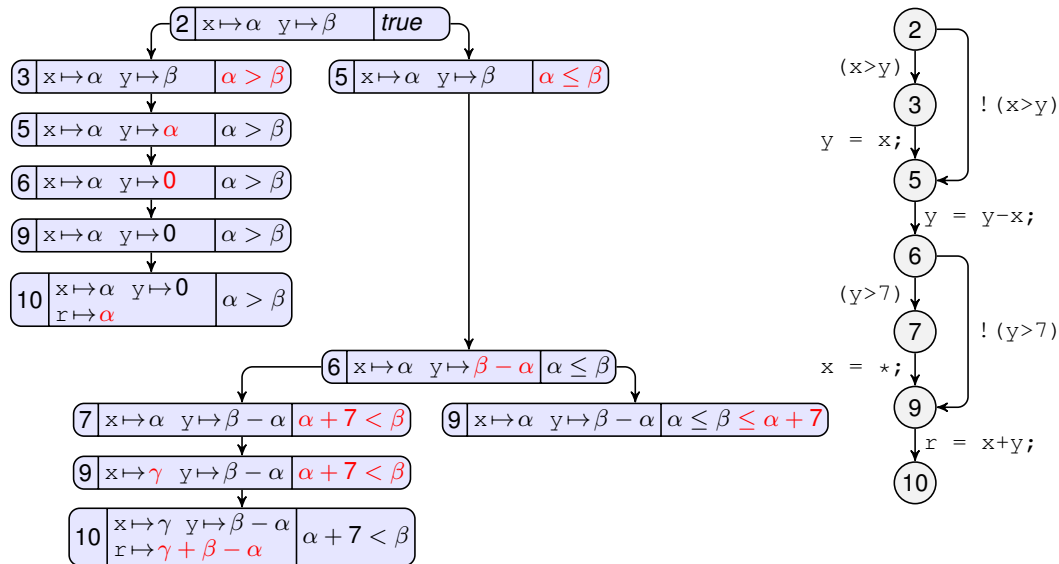
Symbolic execution tree: example 1



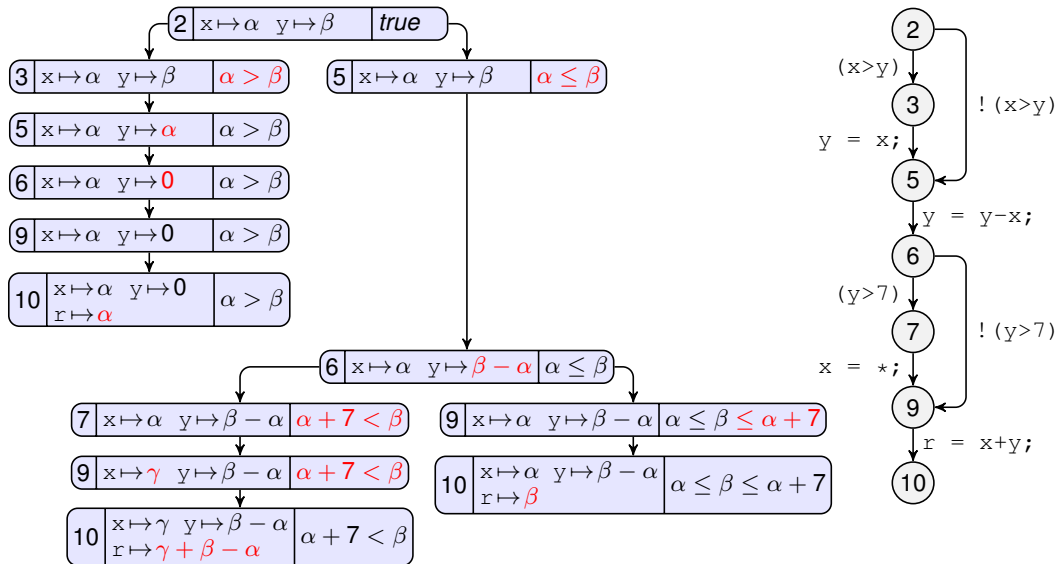
Symbolic execution tree: example 1



Symbolic execution tree: example 1

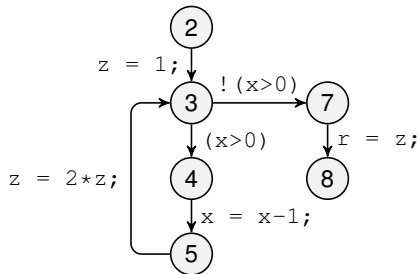


Symbolic execution tree: example 1



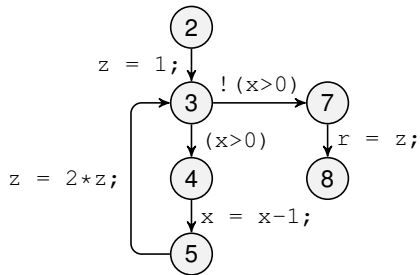
Symbolic execution tree: example 2

```
1 int power(int x) {  
2   int z = 1;  
3   while (x>0) {  
4     x = x-1;  
5     z = 2*z;  
6   }  
7   return z;  
8 }
```

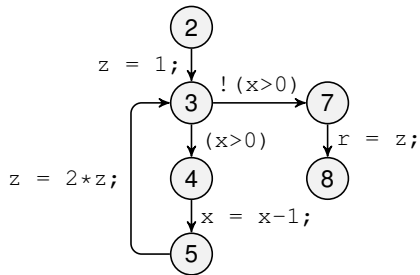
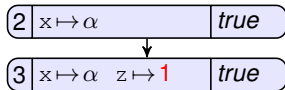


Symbolic execution tree: example 2

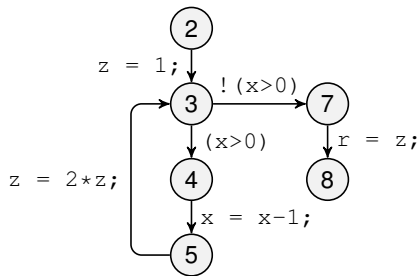
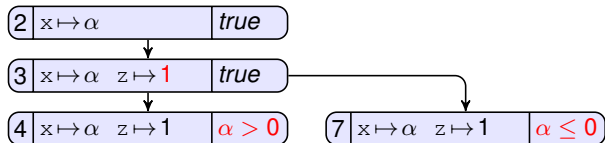
2	$x \mapsto \alpha$	<i>true</i>
---	--------------------	-------------



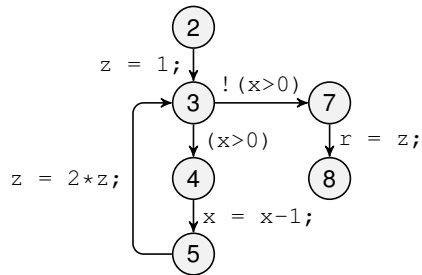
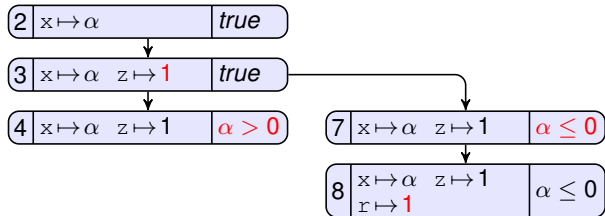
Symbolic execution tree: example 2



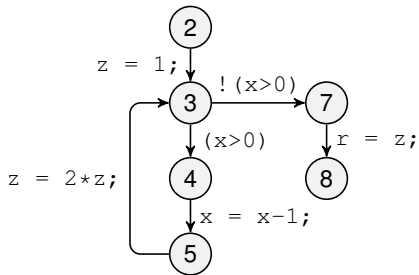
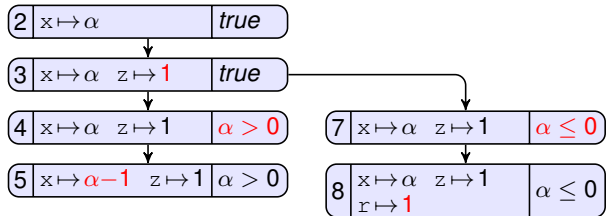
Symbolic execution tree: example 2



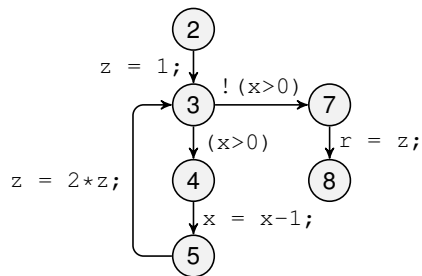
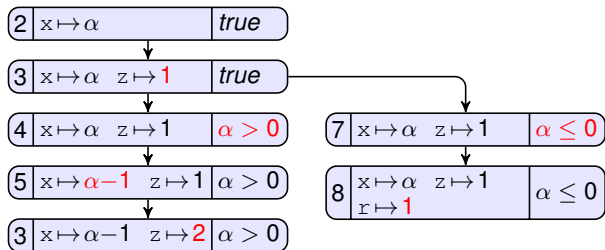
Symbolic execution tree: example 2



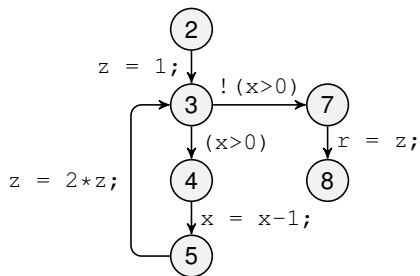
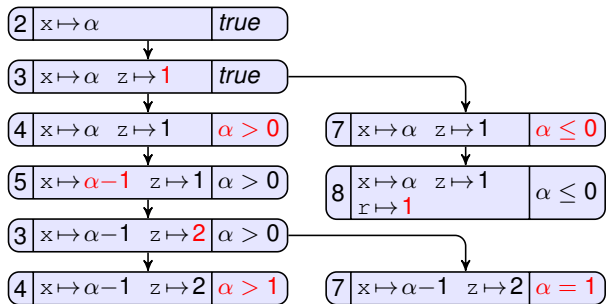
Symbolic execution tree: example 2



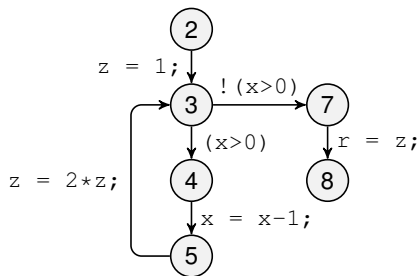
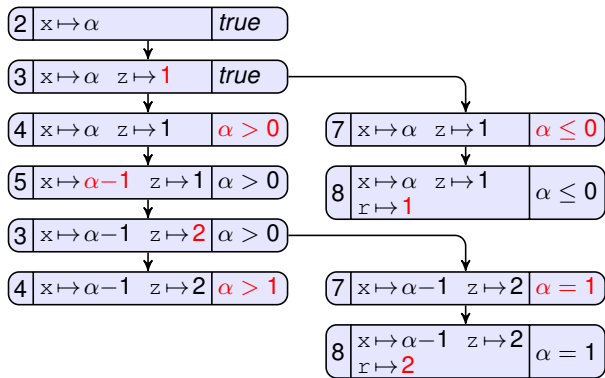
Symbolic execution tree: example 2



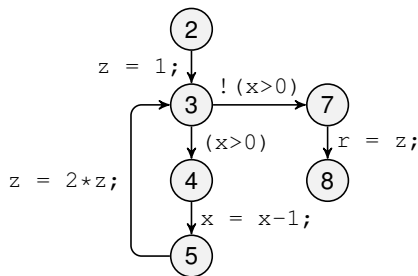
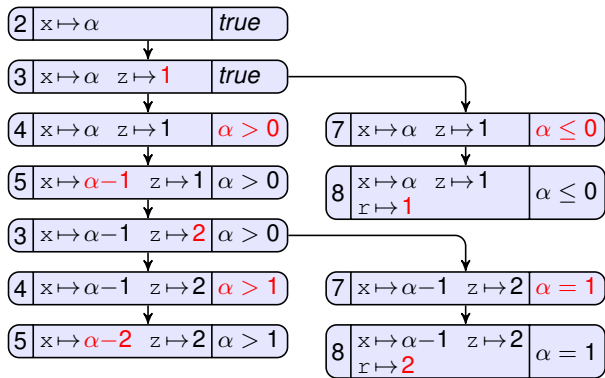
Symbolic execution tree: example 2



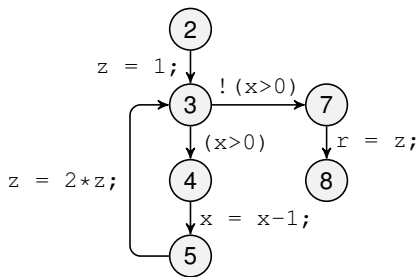
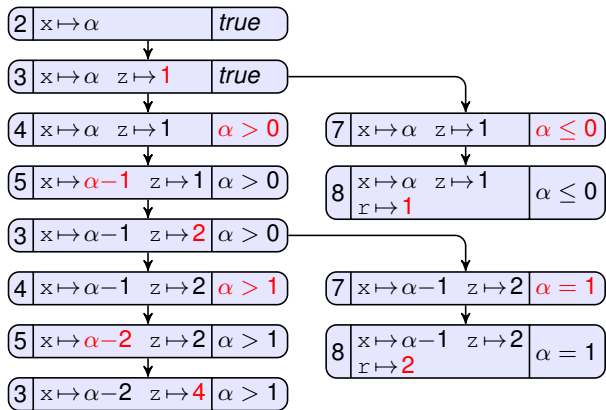
Symbolic execution tree: example 2



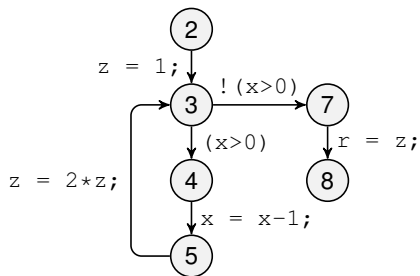
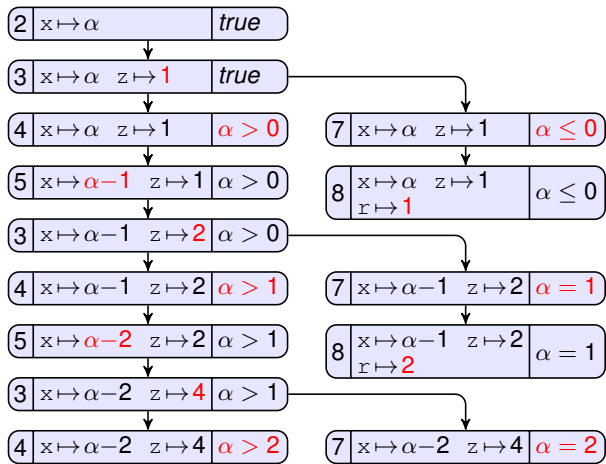
Symbolic execution tree: example 2



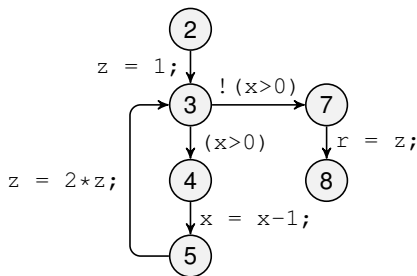
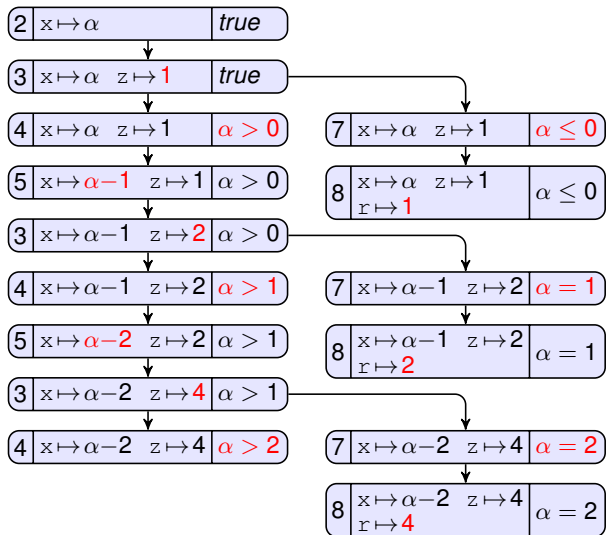
Symbolic execution tree: example 2



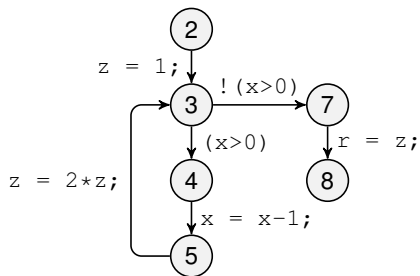
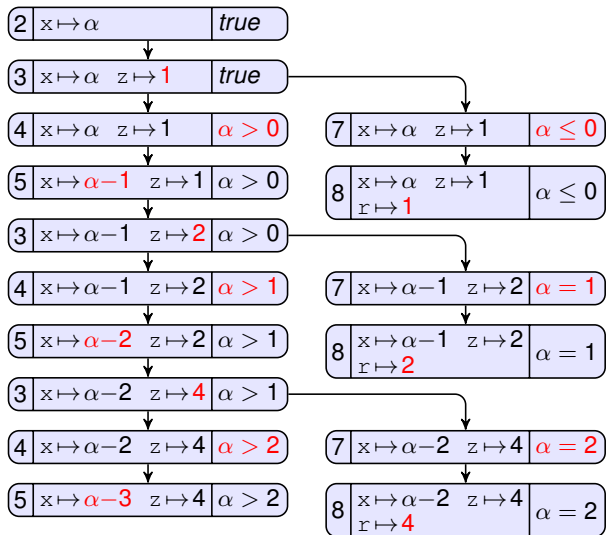
Symbolic execution tree: example 2



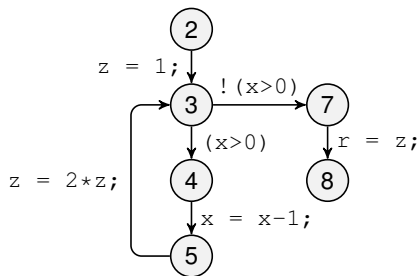
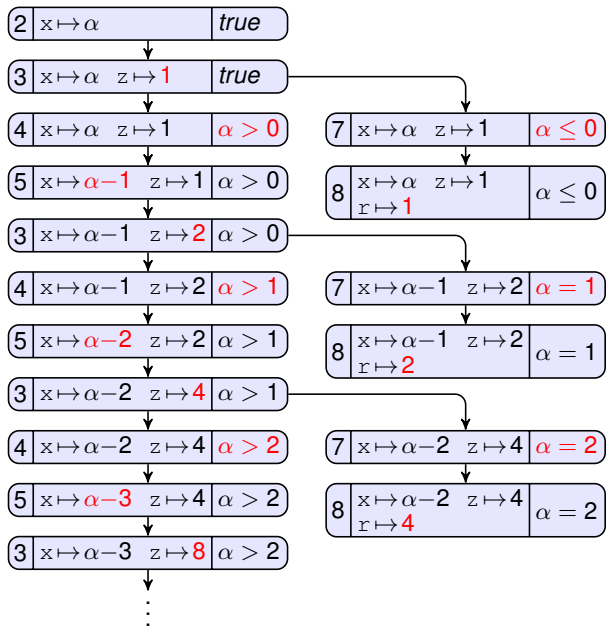
Symbolic execution tree: example 2



Symbolic execution tree: example 2



Symbolic execution tree: example 2



Properties of symbolic execution

- there is a bijection between paths in the symbolic execution tree (starting in its root) and feasible execution paths of the program
- the path condition gives the necessary and sufficient condition on input values to drive the execution along the corresponding path

Properties of symbolic execution

- there is a bijection between paths in the symbolic execution tree (starting in its root) and feasible execution paths of the program
- the path condition gives the necessary and sufficient condition on input values to drive the execution along the corresponding path
- in each symbolic execution, the path condition is satisfiable
 - initially, pc is set to *true*
 - pc is changed only when both branches of a branching statements are feasible
 - pc is extended with a conjunct corresponding to the corresponding branch and the new conjunct is satisfiable as the branch is feasible
- path conditions pc_1, pc_2 corresponding to two distinct leaves of the symbolic execution tree are mutually exclusive, i.e., $pc_1 \wedge pc_2 \equiv false$
- if the symbolic execution tree is finite, then the disjunction of all path conditions in its leaves is equivalent to *true*

Programs can be enriched with `assume(φ)` and `assert(φ)` statements. When symbolic execution passes through

- `assume(φ)`, it executes $pc \leftarrow pc \wedge \varphi$.
- `assert(φ)` and $pc \implies \varphi$ is not valid, it reports an **error**.

With these constructs, symbolic execution can be used with a modification of Floyd's proof method to prove program correctness.

This application is straightforward for any program whose symbolic execution tree is finite.

- deciding validity or satisfiability of formulas can be expensive or even impossible (e.g. for our simple language with unbounded data types)
- in practice, symbolic execution uses expressions and formulas over **bitvector theory** (operations and relations correspond to CPU instructions, e.g. arithmetic operations with overflows, bitwise operations, etc.), where validity and satisfiability are decidable (but expensive)

variable storage referencing problem

- when i is dependent on input, then $A[i]$ can point to various locations in memory
- unsatisfactory solution:
handle $A[i]$ as $ITE(i = 1, A[1], ITE(i = 2, A[2], \dots))$

variable storage referencing problem

- when i is dependent on input, then $A[i]$ can point to various locations in memory
- unsatisfactory solution:
handle $A[i]$ as $ITE(i = 1, A[1], ITE(i = 2, A[2], \dots))$

other memory related problems

- reading/writing via pointers
- comparison of addresses (inner program nondeterminism)
- allocation of memory blocks of symbolic size

variable storage referencing problem

- when i is dependent on input, then $A[i]$ can point to various locations in memory
- unsatisfactory solution:
handle $A[i]$ as $ITE(i = 1, A[1], ITE(i = 2, A[2], \dots))$

other memory related problems

- reading/writing via pointers
- comparison of addresses (inner program nondeterminism)
- allocation of memory blocks of symbolic size

solution: fully symbolic memory model

- performance issues

- **path explosion problem**
 - the number of branches in the symbolic execution tree can be extremely high or even infinite
 - typical for program cycles with the number of iterations depending on the input (symbolic execution forks again and again)
 - construction of full symbolic execution tree is often infeasible
- issues with complex arithmetic operations (e.g. in hashing, encryption or decryption), calls to the operating system and libraries
- practical solutions
 - concretization
 - **concolic execution**

- **concolic** = **concrete** + **symbolic**
- program is executed on a real input and on symbolic input simultaneously
- symbolic execution does not fork, it always follows the concrete execution and computes pc
- if a symbolic value is not available, we can switch to a concrete one

- typical applications
 - bug finding
 - test generation
 - analysis of abstract error traces
- often combined with other techniques
- used in many tools including Klee, PEX, SAGE, SLAM, Ultimate Automizer, Symbiotic

Automated whitebox fuzz testing

Automated whitebox fuzz testing

- an example of modern and sophisticated testing method
- implemented in **SAGE** (Scalable, Automated, Guided Execution)
- discovered 30+ new bugs in large-shipped (and thus intensively tested) file-reading Windows applications including image processors, media players, file decoders
- combines **fuzz testing** and **symbolic execution**

- symbolic execution is expensive compared to running tests
- thus we want to generate as many new inputs from one symbolic execution as possible
- input for the next symbolic execution is selected by some scoring function applied to all generated inputs
- in particular, the input that explored the most (so-far uncovered) pieces of code is chosen for the next symbolic execution

The main algorithm

```
1 procedure GenerateInputs(inputSeed)
2   inputSeed.bound  $\leftarrow$  0
3   workList  $\leftarrow$  {inputSeed}
4   Run&Check(program, inputSeed)
5   while workList  $\neq$   $\emptyset$  do
6     input  $\leftarrow$  PickFirstItem(workList)
7     childInputs  $\leftarrow$  ExpandExecution(input)
8     foreach newInput  $\in$  childInputs do
9       Run&Check(program, newInput)
10      Score(newInput)
11      workList  $\leftarrow$  workList  $\cup$  {newInput}
```

- Score(newInput) counts the newly covered blocks
- workList is ordered by the score of inputs
- PickFirstItem(workList) returns the input with the highest score

Application of symbolic execution

```
1 procedure ExpandExecution(input)
2   childInputs  $\leftarrow \emptyset$ 
3   PC  $\leftarrow$  SymbolicExecution(program, input)
4   for  $j \leftarrow$  input.bound to |PC| - 1 do
5     if  $\bigwedge_{i=0}^{j-1} \text{PC}[i] \wedge \neg \text{PC}[j]$  has solution  $M$  then
6       newInput  $\leftarrow$  Combine(input,  $M$ )
7       newInput.bound  $\leftarrow j$ 
8       childInputs  $\leftarrow$  childInputs  $\cup$  {newInput}
9   return childInputs
```

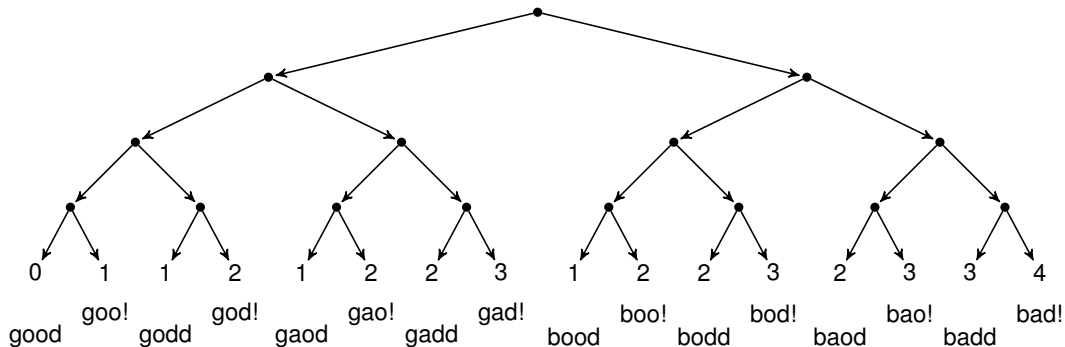
- `Combine(input, M)`
 - creates a new input from the original input and M
 - `Combine("abcde", input[3] \mapsto "F")` returns "abcFe"
- path conditions are represented as arrays PC of conjuncts

Example

```
1 void top(char input[4]) {
2   int cnt=0;
3   if (input[0] == 'b') cnt++;
4   if (input[1] == 'a') cnt++;
5   if (input[2] == 'd') cnt++;
6   if (input[3] == '!') cnt++;
7   if (cnt >= 3) abort(); // error
8 }
```

Example

```
1 void top(char input[4]) {  
2   int cnt=0;  
3   if (input[0] == 'b') cnt++;  
4   if (input[1] == 'a') cnt++;  
5   if (input[2] == 'd') cnt++;  
6   if (input[3] == '!') cnt++;  
7   if (cnt >= 3) abort(); // error  
8 }
```



- the algorithm can be parallelized: only workList and the overall block coverage need to be shared
- SAGE recovers easily from divergencies (situations when an execution deviates from the assumed execution path) induced e.g. by inner program nondeterminism
- SAGE runs 24/7 on large clusters, available for Microsoft developers