## **Formal Methods**

Cinzia Bernardeschi

Department of Information Engineering University of Pisa, Italy

FMSS, 2020-2021

◆□ > ◆□ > ◆豆 > ◆豆 > □ = − の < ⊙

#### **Overview**

- Formal methods definition
- Specification languages
- Verification techniques
  - model checking
  - abstract interpretation
  - theorem proving
- Case studies:
  - Malware analysis
  - Data leakage
  - Cyber-physical systems attacks

▲□▶▲□▶▲□▶▲□▶ □ のQ@

Formal methods are mathematically-based techniques that can be used in the design and development of computer-based systems.

Formal methods

- allow the analysis of all possible executions of the system
- improve the current techniques based on simulation and testing (mathematical proof that the system behaves as expected)
- offer the possibility for detecting vulnerability in systems or for building more secure systems by design.

▲□▶▲□▶▲□▶▲□▶ □ のQ@

## Formal methods and Safety critial systems

Software used in safety critical systems is required a very high reliability

- For example, for flight-critical avionics systems in civil transport airplaines the requirement is probability of failure 10<sup>-9</sup> per hour of operation. 1 failure at most every 1.000.000.000 hours of operation.
- Provide evidence of reliability. We cannot run the system for 10<sup>-9</sup> hours. Long operational testing and doubts about the representativeness of the testing.

One of the best hope we have for technological support in building dependable software is the use of formal methods. (J. Knight, Fundamentals of Dependable Computing, 2012.)

A formal method consists of

 a language a mathematical notation or a computer language with a formal semantics

▲□▶▲□▶▲□▶▲□▶ □ のQ@

- a set of tools for proving properties without real executions
- a methodology for its application in industrial practice.

#### Formal methods definition

Formal methods are more ambitious than traditional approaches, they are more complex and the associated tools are more difficult to build.

A major problem comes from undecidability result of Computational complexity theory

- some verification problems are either impossible or very difficult to solve automatically
- For example, termination problem of programs (in the general case, there is no decision procedure that can determinate whether a program P may terminate or not)
- Another example: there is no decision procedure that can determinate whether a given instruction of program P will ever be executed

#### Formal methods definition

Strategies to deal with undecidability:

- expressiveness restrictions identify classes of systems for which verification is decidable at least in principle (example, finite systems)
- accuracy restrictions weaker formulation of the problem which is decidable and of interest (approximation instead of exact solutions)
- automatic restrictions semi-automatic verification, with human intervention at some point semi-decision procedures, which may terminate giving the correct result or not terminate.

## **Specifications**

Specifications are means to describe a system, its components, its global/local environment.

#### declarative specifications

describe what a system should do, but not how it should do it it is not possible to derive automatically from these constraints an implementation of the system

operational specifications possibly define what a system should do, and definitely define how it should do it it should be possible to derive automatically at least a skeleton of an implementation of the system

# Formal methods for modelling and verification

- formal methods for the analysis of programs
- formal methods for the analysis of systems

# Formal methods for the analysis of programs

- Semantics of the programming language the semantics of the language describes mathematically the meaning of the programs. Operational semantics as semantics of the programming language.
- Models of programs. (A model of a system is an abstraction that focus on some part of the system.) *Control flow graphs* as models of a program. The control flow graph represents the control structure of the program.

In the following:

- structured high level language (e.g., C language)
- Iow level languages (e.g, Java bytecode, assembly code)

# A simple high level language

We consider a simple sequential language with the following syntax, where *op* stands for the usual arithmetic and logic operations

▲□▶▲□▶▲□▶▲□▶ □ のQ@

Instruction set

exp::= const | var | exp op exp
com: = var := exp | if exp then com else com |
while exp do com | com; com | halt

A program *P* is a sequence of instructions  $\langle c \rangle$ , where  $c \in com$ .

# Standard Operational semantics

State of the program:

x 2 y 5 z 3 k 4 memory m

\$\langle c, m \rangle\$
 where c is a command and m is a memory

► ⟨*m*⟩

a single memory, in the final state.

The semantics of programs is given by means of a transition system and we call this semantics *execution semantics*.

- ► Var(c) denote the set of variables occurring in c.
- ▶  $\mathcal{V}^{\epsilon}$  is the domain of constant values, ranged over by k, k', ...,
- ▶ for each  $X \subseteq var$ , the domain  $\mathcal{M}_X^{\epsilon} = X \to \mathcal{V}^{\epsilon}$  of memories defined on X, ranged over by  $m, m', \ldots$

#### Transitions

The semantic rules define a relation

 $\longrightarrow^{\epsilon} \subseteq \mathcal{Q}^{\epsilon} \times \mathcal{Q}^{\epsilon}$ 

where  $\mathcal{Q}^{\epsilon}$  is a set of states.

A separate transition

$$\longrightarrow_{\mathit{expr}}^{\epsilon} \subseteq (\mathit{exp} imes \mathcal{M}^{\epsilon}) imes \mathcal{V}^{\epsilon}$$

is used to compute the value of the expressions.

With m[k/x] we denote the memory m' which agrees with m on all variables, except on x, for which m'(x) = k.

# **Operational semantics**

**Expr**<sub>const</sub> 
$$\overline{\langle k, m \rangle \longrightarrow_{expr}^{\epsilon} k}$$

Expr<sub>var</sub> 
$$\overline{\langle x, m \rangle \longrightarrow_{expr}^{\epsilon} m(x)}$$

$$\mathsf{Expr}_{op} \quad \frac{\langle e_1, m \rangle \longrightarrow_{expr}^{\epsilon} k_1 \ \langle e_2, m \rangle \longrightarrow_{expr}^{\epsilon} k_2 \ k_1 \ op \ k_2 = k_3}{\langle (e_1 \ op \ e_2), m \rangle \longrightarrow_{expr}^{\epsilon} k_3}$$

▲□▶▲□▶▲≡▶▲≡▶ ≡ のへで

## **Operational semantics**

Ass 
$$\frac{\langle e, m \rangle \longrightarrow_{expr}^{\epsilon} k}{\langle x := e, m \rangle \longrightarrow^{\epsilon} m[k/x]}$$

halt 
$$\overline{\langle \text{halt}, m \rangle \longrightarrow^{\epsilon} \langle m \rangle}$$

$$\mathbf{Seq}_{2} \quad \frac{\langle \boldsymbol{c}_{1}, \boldsymbol{m} \rangle \longrightarrow^{\epsilon} \langle \boldsymbol{c}_{2}, \boldsymbol{m}' \rangle}{\langle \boldsymbol{c}_{1}; \boldsymbol{c}_{3}, \boldsymbol{m} \rangle \longrightarrow^{\epsilon} \langle \boldsymbol{c}_{2}; \boldsymbol{c}_{3}, \boldsymbol{m}' \rangle}$$

▲□▶▲圖▶▲≣▶▲≣▶ ≣ のQ@

# **Operational semantics**

$$If_{true} \quad \frac{\langle e, m \rangle \longrightarrow_{expr}^{\epsilon} true}{\langle \text{if } e \text{ then } c_1 \text{ else } c_2, m \rangle \longrightarrow^{\epsilon} \langle c_1, m \rangle }$$

$$\begin{array}{c} \text{If}_{\textit{false}} & \frac{\langle e, m \rangle \longrightarrow_{expr}^{\epsilon} \textit{false}}{\langle \text{if } e \text{ then } c_1 \text{ else } c_2, m \rangle \longrightarrow^{\epsilon} \langle c_2, m \rangle} \end{array}$$

While true 
$$\langle e, m \rangle \longrightarrow_{expr}^{\epsilon} true$$
  
 $\langle while e do c, m \rangle \longrightarrow^{\epsilon} \langle c; while e do c, m \rangle$ 

While 
$$_{false} \xrightarrow{\langle e, m \rangle \longrightarrow_{expr}^{\epsilon} false} \langle while e do c, m \rangle \longrightarrow^{\epsilon} \langle m \rangle}$$

Given a program  $P = \langle c \rangle$  and an initial memory  $m \in \mathcal{M}^{\epsilon}_{Var(c)}$ , we denote by E(P, m)

the transition system defined by  $\rightarrow^{\epsilon}$  starting from the initial state  $\langle c, m \rangle$ .

Actually, since the program is deterministic, there exists at most one final state, i.e. a state  $\langle m \rangle$  for some memory *m*.

・ロト・日本 キョン・ヨン ヨー シック

#### **Transition system**

A transition system is a tuple  $TS = (S, Act, \rightarrow, I)$  such that

- S is a set of states (the state space)
- Act is a set of actions
- ►  $\rightarrow \subseteq S \times Act \times S$  is a transition relation

 $s \rightarrow_a s'$  if from state *s* the system moves to state *s'* by executing action *a* 

•  $I \subseteq S$  is the set of initial states

Often, transition systems are drawn as directed graphs with states represented by vertices and transitions represented by edges

## An example

E(P, m) x:=2; y:=5; х ..... > 2 х if (x>0) then У < y:=y+y; y:=y+y; < halt; у 5 else y:=0; s0 s3 halt; x:=2 1: x:=2; y:=y+y 2: y:=5; 2 10 y:=5; 2 3: if (x>0) х halt; > < if (x>0) then v < 4: then y:=y+y; y:=y+y; s4 else y:=0; 5: else y:=0; s1 halt; 6: halt; skip y:=5 2 10 x < х 2 if (x>0) then ⇒ s5 ٧ < y:=y+y; else y:=0; s2 if (x>0) then halt; y:=y+y; else y:=0;

◆□▶ ◆□▶ ◆三▶ ◆三▶ ◆□▶ ◆□▶

The control flow graph of a program  $P = \langle c \rangle$  is a directed graph (*V*; *E*), where *V* is a set of nodes and *E* : *VxV* is a set of edges connecting nodes.

Nodes correspond to instructions. Moreover, there is an initial node and a final node that represent the starting point and the final point of an execution. E contains the edge (i; j) if and only if the instruction at address j can be immediately executed after that at address i.

The control flow graph does not contain information on the semantics of the instructions.

◆□▶ ◆□▶ ◆ □ ▶ ◆ □ ▶ ◆ □ ● ● ● ●

# An example



C++ generation of the CFG with Visual Studio. Gcc developer-options: -fdump-tree-cfg -blocks -vops

# An example



:

#### Java language

```
3 public class Hello {
    public static void main(String[] args) {
4
5
6
    public void stampa(String s){
           System.out.println("Outside");
8
9
      try { System.out.println("Inside try");}
10
      catch (ArithmeticException e) {System.out.println("Print catch");}
11
12
     System.out.println("Print ....");
13
14
15
```

#### Java bytecode

```
1 Compiled from "Hello, java"
2 public class Hello {
    public Hello();
      Code:
         0: aload_0
         1: invokespecial #1 // Method java/lang/Object."<init >":()V
         4: return
    public static void main(java.lang.String[]);
9
      Code:
10
         0: return
    public void stampa(java.lang.String);
      Code:
14
         0: getstatic #2 // Field java/lang/System.out: Ljava/io/ PrintStream;
         3: 1dc
                      #3 // String Outside
16
         5: invokevirtual #4 // Method java/io/ PrintStream. println:(Ljava/lang/String:)V
17
         8: getstatic #2 // Field java/lang/System.out: Ljava/io/PrintStream;
18
        11: 1dc
                           #5 // String Inside try
19
        13: invokevirtual #4 // Method java/io/PrintStream.println:(Ljava/lang/String;)V
20
        16: goto
                           28
        19: astore_2
        20: getstatic
                           #2 // Field java/lang/System.out: Ljava/io/PrintStream;
23
                           #7 // String Print catch
        23: Idc
24
        25: invokevirtual #4 // Method java/io/PrintStream.println:(Ljava/lang/String;)V
25
        28: getstatic
                         #2 // Field java/lang/System.out: Ljava/io/PrintStream;
26
        31: 1dc
                           #8 // String Print ...
27
        33: invokevirtual #4 // Method java/io/PrintStream.println:(Ljava/lang/String;)V
28
        36: return
29
30
```

Java bytecode

A subset of the instruction set

Bytecode instructions operate on the operand stack. Variables in memory are named registers.

pop Pop top operand stack element.

- dup Duplicate top operand stack element.
- op Pop two operands off the operand stack, perform the operation op ∈ { add, mult, compare .. }, and push the result onto the stack.
- const *d* Push constant *d* onto the operand stack.

### Java bytecode

load xPush the value of the register xonto the operand stack.

store *x* Pop a value off the operand stack and store it into local register *x*.

if cond j Pop a value off the operand stack, and evaluate it against the condition cond = { eq, ge, null, ... }; branch to j if the value satisfies cond.

▲□▶▲□▶▲□▶▲□▶ □ のQ@

goto j Jump to j.

# Standard Operational semantics

The domain of the states of the standard semantics is (A denotes the set of addresses):

 $\mathcal{Q}^{\epsilon} = \mathcal{A}^{\epsilon} \times \mathcal{M}^{\epsilon} \times \mathcal{S}^{\epsilon}.$ 

A bytecode B consists of a sequence of instructions (assume 1 is the address of first instruction), and is executed starting from a memory and an empty operand stack.

A state is given by the value of three variables, PC, MEM and STACK, where

- PC is the program counter,
- MEM is the memory, and
- *STACK* is the operand stack (  $\lambda$  is the empty stack).

We denote by (i, m, s) the state labeled by PC = i, MEM = m, STACK = s.

## Standard Operational semantics rules

$$op \quad \frac{c[i] = op}{\langle i, m, k_1 \cdot k_2 \cdot s \rangle \longrightarrow^{\epsilon} \langle i+1, m, (k_1 \text{ op } k_2) \cdot s \rangle }$$

$$\mathsf{pop} \quad \frac{c[i] = \mathsf{pop}}{\langle i, m, k \cdot s \rangle \longrightarrow^{\epsilon} \langle i+1, m, s \rangle}$$

$$\mathsf{push} \quad \frac{c[i] = \mathsf{push} \mathsf{k}}{\langle i, m, s \rangle \longrightarrow^{\epsilon} \langle i + 1, m, k \cdot s \rangle}$$

load 
$$\frac{c[i] = \text{load } x}{\langle i, m, s \rangle \longrightarrow^{\epsilon} \langle i + 1, m, m(x) \cdot s \rangle}$$

▲□▶▲□▶▲≡▶▲≡▶ ≡ のへで

## Standard Operational semantics rules

$$\begin{array}{l} \text{store} \quad \frac{c[i] = \text{store } \mathbf{x}}{\langle i, m, k \cdot s \rangle \longrightarrow^{\epsilon} \langle i + 1, m[k/x], s \rangle} \\ \\ \text{if}_{\textit{false}} \quad \frac{c[i] = \text{if } \text{j}}{\langle i, m, 0 \cdot s \rangle \longrightarrow^{\epsilon} \langle i + 1, m, s \rangle} \\ \\ \\ \text{if}_{\textit{true}} \quad \frac{c[i] = \text{if } \text{j}}{\langle i, m, k \neq 0 \cdot s \rangle \longrightarrow^{\epsilon} \langle j, m, s \rangle} \end{array}$$

goto 
$$c[i] = \text{goto } j$$
  
 $\langle i, m, s \rangle \longrightarrow^{\epsilon} \langle j, m, s \rangle$ 

▲□▶▲□▶▲≡▶▲≡▶ ≡ のへで

# Transition system

E(B, m, λ) 1 ←> .... < 1 7 > < 4 .... ..... s3 s0 1: load y 2: if 5 2 ·← > 1 <6 J 7 .... 3: push 1 7 < 2 7 > s4 4: goto 6 .... 5: push 0 s1 6: store x 1 7: halt <7. > .... s5 2 7 < 3 , ŷ > .... s2 > < 7 ....

◆□ ▶ ◆□ ▶ ◆三 ▶ ◆□ ▶ ◆□ ▶

s6

# A bytecode and its CFG



- 2: if 5
- 3: push 1
- 4: goto 6
- 5: push 0
- 6: store x
- 7: halt



The CFG contains the edge (i, j) if and only if the instruction at address *j* can be immediately executed after that at address *i*. ([6] is the point at which the branches starting at the conditional instruction [2] join)

▲□▶▲圖▶▲≣▶▲≣▶ ■ めんの

# Java bytecode Instructions

рор	Pop top operand stack element.		
dup	Duplicate top operand stack element.		
α <b>ορ</b>	Pop two operands with type $lpha$ off the operand stack,		
	perform the operation $op \in \{ \text{ add, mult, compare } \}$ ,		
	and push the result onto the stack.		
$lpha { m const} \; {\pmb d}$	Push constant d with type $\alpha$ onto the operand stack.		
$lpha$ load ${m x}$	Push the value with type $lpha$ of the register $m{x}$		
	onto the operand stack.		
$lpha$ store ${\it X}$	Pop a value with type $\alpha$ off the operand stack and		
	store it into local register x.		
if <i>cond j</i>	Pop a value off the operand stack, and evaluate it against		
-	<pre>the condition cond = { eq, ge, null, };</pre>		
	branch to <i>j</i> if the value satisfies <i>cond</i> .		
goto <b>j</b>	Jump to j.		

#### Java bytecode Instructions

getfield <b>C</b> .f	Pop a reference to an object of class C
	off the operand stack; fetch the object's
	field f and put it onto the operand stack.

putfield *C.f* Pop a value *k* and a reference to an object of class *C* from the operand stack; set field *f* of the object to *k*.

invoke *C.mt* Pop value *k* and a reference *r* to an object of class *C* from the operand stack; invoke method *C.mt* of the referenced object with actual parameter *k*.

 $\alpha$  return Pop the  $\alpha$  value off the operand stack and return it

from the method.

The bytecode of a method is a sequence *B* of instructions.

When a method is invoked (invoke instruction), it executes with a new empty stack and with an initial memory where all registers are undefined except for the first one, register x0, that contains the reference to the object instance on which the method is called, and register x1, that contains the actual parameter.

When the method returns, control is transferred to the calling method: the caller's execution environment (operand stack and local registers) is restored and the returned value, if any, is pushed onto the operand stack.

< □ > < 同 > < 三 > < 三 > < 三 > < ○ < ○ </p>

# An example

The bytecode corresponds to a method *mt* of a class A. Suppose that register x1 (the parameter of *A.mt*) contains a reference to an object of another class B. Note that register x0 contains a reference to A. After the bytecode has been executed, the final value of field f1 of the object of class A is 0 or 1 depending on the value of field f2 of the object of class B.

0:	aload	<i>x</i> 0
1:	aload	<i>x</i> 1
2 :	getfield	B.f2
3:	ifge	6
4:	iconst	0
<b>5</b> :	goto	7
6:	iconst	1
7:	putfield	A.fl
8:	iconst	1
9:	return	