

Operational Semantics

Lecture 2: (Java) Bytecode

Simon Castellan

SOS, Master Recherche Science Informatique, U. Rennes 1

2021-2022

Semantics of Java bytecode

An example of an operational semantics in the form of an abstract machine (the Java virtual machine).

References:

- ▶ Tim Lindholm , Frank Yellin, Gilad Bracha, Alex Buckley. Java Virtual Machine Specification, Java SE 8 Edition, 2015.
- ▶ Robert F. Stärk, Joachim Schmid, Egon Börger. Java and the Java Virtual Machine – Definition, Verification, Validation. 2001.
- ▶ S. Freund , J. Mitchell : A Type System for the Java Bytecode Language and Verifier, Journal of Automated Reasoning, Volume 30, Issue 3-4, Pages: 271 - 321. 2003.

Outline

- 1 Background on Java and JVM
- 2 Idealized bytecode syntax
- 3 Operational Semantics

Java

Java: a class-based, object-oriented programming language.

```
public class Bicycle{
    private int gear;
    private int id;
    private int speed;

    private static int numberOfBicycles = 0;

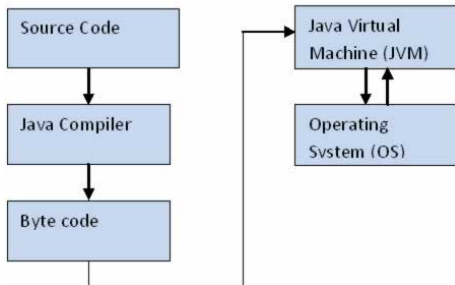
    public Bicycle(int startSpeed, int startGear){
        gear = startGear;
        speed = startSpeed;
        id = ++numberOfBicycles;
    }

    public void setGear(int newValue){
        gear = newValue;
    }
}

public class MountainBike extends Bicycle {
    // the MountainBike subclass has one field
    public int seatHeight;
    ...
}
```

JVM and Java bytecode

To enhance portability, Java defines a **bytecode** language that is interpreted by the **Java virtual machine**.



Java bytecode: factorial

The JVM is a stack machine with registers.

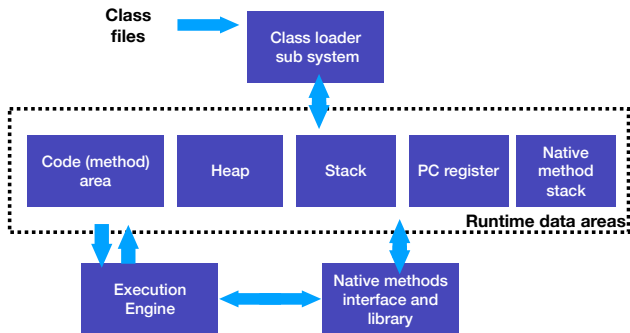
Source code

```
public static int factorial(int n) {
    int res = 1 ;
    for (int i = 1; i <= n; i++) { res = res * i; }
    return res;
}
```

Bytecode (as produced by `javap -v -c Factorial.class`)

```
int factorial(int); descriptor: (I)I flags: ACC_PUBLIC, ACC_STATIC, stack=2, locals=3, args_size=1
 0: iconst_1          // push the integer constant 1
 1: istore_1          // store it in register 1 (the res variable)
 2: iconst_1          // same for register 2 (the i variable)
 3: istore_2
 4: goto             14          // go to PC 14
 7: iload_1           // push value of register 1 (res)
 8: iload_2           // push value register 2 (i)
 9: imul             // multiply the two integers at top of stack
10: istore_1          // pop result and store it in register 1
11: iinc              2, 1        // increment register 2 (i) by 1
14: iload_2           // evaluate loop condition
15: iload_0           //
16: if_icmple        7          // iterate or exit loop
19: iload_1           // push value of register 1 (res)
20: ireturn          // return the top integer to the caller
```

The Java virtual machine state



The JVM has several components:

Code component includes bytecode of all methods of all classes loaded.

Program counter points to the next instruction to execute.

Heap (memory) of objects created while program is running.

Frame stack of methods under execution. Each frame contains an operand stack and local variables (registers). (*ignored in our modelisation*)

The Java virtual machine specification

6.5 Instructions

THE JAVA VIRTUAL MACHINE INSTRUCTION SET

*iadd**iadd***Operation** Add `int`**Format**

<i>iadd</i>

Forms *iadd* = 96 (0x60)**Operand** ..., *value1*, *value2* →**Stack** ..., *result***Description** Both *value1* and *value2* must be of type `int`. The values are popped from the operand stack. The `int` *result* is *value1* + *value2*. The *result* is pushed onto the operand stack.

The result is the 32 low-order bits of the true mathematical result in a sufficiently wide two's-complement format, represented as a value of type `int`. If overflow occurs, then the sign of the result may not be the same as the sign of the mathematical sum of the two values.

Despite the fact that overflow may occur, execution of an *iadd* instruction never throws a run-time exception.

Outline

- 1 Background on Java and JVM
- 2 Idealized bytecode syntax**
- 3 Operational Semantics

Idealized bytecode syntax

Instruction ::=

nop	<i>Does nothing</i>
Stack manipulation	
const <i>c</i>	<i>Push constant</i>
pop	<i>Pop and discard top value</i>
dup	<i>Duplicate top entry</i>
dup2	<i>Duplicate top two entries</i>
swap	<i>Swap top two entries</i>
numop <i>op</i>	<i>Arithmetic operation</i>
Local registers manipulation	
load <i>x</i>	<i>Push value of <i>x</i></i>
store <i>x</i>	<i>Pop top entry into <i>x</i></i>

Some bytecodes are typed (*iconst*, *aload*) to indicate whether they operate on integers, references, arrays, etc.

Idealized bytecode syntax (continued)

Instruction ::= ...

Manipulation of program counter

| `ifl` *pc*

Conditional jump

| `goto` *pc*

Inconditional jump

Object manipulation

| `new` *cl*

Creation of object

| `putfield` *f*

Object field modification

| `getfield` *f*

Object field access

Real Java bytecode: more “if”s: `if_eq`, `if_ge`, `if_nonnull`, `if_acmpeq`,...

Fields *f* are of the form `C : name : t` (*C*: class declaring the field; and *t* is its type)

Program structure

A program P is:

- ▶ a **set of instructions**: $\text{instructionAt}_P(m, pc)$ fetches the opcode at program counter pc in P
- ▶ a set of **classes** describing the list of fields: $\text{Fields}(c)$ describes the list of fields of class c .

A function in the code of method m in program P .

Outline

- 1 Background on Java and JVM
- 2 Idealized bytecode syntax
- 3 Operational Semantics**

Operational semantics

- ▶ Define semantic domains
- ▶ Define the SOS transition system
- ▶ Transition system parametrised by a specific program P .
- ▶ **Difficulty:** Low-level control flow (jumps).

Operational semantics

Semantic domains

Value = num n $n \in \mathbb{Z}$
 ref r $r \in \text{Reference}$
 null

Stack = Value* (a list of values, can be empty)

LocalVar = Var \hookrightarrow Value

CallFrame = ProgCount \times LocalVar \times Stack

Object = (FieldName \hookrightarrow Value)

Heap = Reference \hookrightarrow Object

State = Heap \times CallFrame

Operational semantics

Example rule (nop)

$$\frac{\text{instructionAt}_p(pc) = \text{nop}}{\langle\langle h, \langle pc, l, s \rangle \rangle\rangle \rightarrow \langle\langle h, \langle pc+1, l, s \rangle \rangle\rangle}$$

Operational semantics – Basic stack operations.

We use $c, n \dots$ for numeric values, loc for references, and v for arbitrary values.

$$\frac{\text{instructionAt}_p(pc) = \text{nop}}{\langle\langle h, \langle pc, l, s \rangle \rangle\rangle \rightarrow \langle\langle h, \langle pc+1, l, s \rangle \rangle\rangle}$$

$$\frac{\text{instructionAt}_p(pc) = \text{const } c}{\langle\langle h, \langle pc, l, s \rangle \rangle\rangle \rightarrow \langle\langle h, \langle pc+1, l, c :: s \rangle \rangle\rangle}$$

$$\frac{\text{instructionAt}_p(pc) = \text{pop}}{\langle\langle h, \langle pc, l, v :: s \rangle \rangle\rangle \rightarrow \langle\langle h, \langle pc+1, l, s \rangle \rangle\rangle}$$

$$\frac{\text{instructionAt}_p(pc) = \text{dup}}{\langle\langle h, \langle pc, l, v :: s \rangle \rangle\rangle \rightarrow \langle\langle h, \langle pc+1, l, v :: v :: s \rangle \rangle\rangle}$$

$$\frac{\text{instructionAt}_p(pc) = \text{dup2}}{\langle\langle h, \langle pc, l, v_1 :: v_2 :: s \rangle \rangle\rangle \rightarrow \langle\langle h, \langle pc+1, l, v_1 :: v_2 :: v_1 :: v_2 :: s \rangle \rangle\rangle}$$

$$\frac{\text{instructionAt}_p(pc) = \text{swap}}{\langle\langle h, \langle pc, l, v_1 :: v_2 :: s \rangle \rangle\rangle \rightarrow \langle\langle h, \langle pc+1, l, v_2 :: v_1 :: s \rangle \rangle\rangle}$$

Operational semantics – Arithmetic and local variables.

$$\frac{\text{instructionAt}_p(pc) = \text{numop } op}{\ll h, \langle pc, l, n_1 :: n_2 :: s \rangle \gg \rightarrow \ll h, \langle pc+1, l, \llbracket op \rrbracket(n_1, n_2) :: s \rangle \gg}$$

$$\frac{\text{instructionAt}_p(pc) = \text{load } x}{\ll h, \langle pc, l, s \rangle \gg \rightarrow \ll h, \langle pc+1, l, l(x) :: s \rangle \gg}$$

$$\frac{\text{instructionAt}_p(pc) = \text{store } x}{\ll h, \langle pc, l, v :: s \rangle \gg \rightarrow \ll h, \langle pc+1, l[x \mapsto v], s \rangle \gg}$$

Operational semantics – Control transfer.

$$\frac{\text{instructionAt}_p(pc) = \text{ifl e } pc' \quad n \leq 0}{\langle\langle h, \langle pc, l, n :: s \rangle \rangle\rangle \rightarrow \langle\langle h, \langle pc', l, s \rangle \rangle\rangle}$$

$$\frac{\text{instructionAt}_p(pc) = \text{ifl e } pc' \quad n > 0}{\langle\langle h, \langle pc, l, n :: s \rangle \rangle\rangle \rightarrow \langle\langle h, \langle pc+1, l, s \rangle \rangle\rangle}$$

$$\frac{\text{instructionAt}_p(pc) = \text{goto } pc'}{\langle\langle h, \langle pc, l, s \rangle \rangle\rangle \rightarrow \langle\langle h, \langle pc', l, s \rangle \rangle\rangle}$$

Operational semantics – Objects.

NB: we present here a simplified object initialization procedure.

$$\frac{\text{instructionAt}_p(pc) = \text{new } cl \quad \text{Fields}(cl) = \{f_1, \dots, f_n\} \\ \text{loc} \notin \text{dom}(h) \quad o := [f_1 \mapsto \diamond, \dots, f_n \mapsto \diamond] \quad h' = h[\text{loc} \mapsto o]}{\ll h, \langle \text{pc}, l, s \rangle \gg \rightarrow \ll h', \langle \text{pc}+1, l, \text{loc} :: s \rangle \gg}$$

$$\frac{\text{instructionAt}_p(pc) = \text{putfield } f \quad h(\text{loc}) = o \quad o' = o[f \mapsto v]}{\ll h, \langle \text{pc}, l, v :: \text{loc} :: s \rangle \gg \rightarrow \ll h[\text{loc} \mapsto o'], \langle \text{pc}+1, l, s \rangle \gg}$$

$$\frac{\text{instructionAt}_p(pc) = \text{getfield } f \quad h(\text{loc}) = o}{\ll h, \langle \text{pc}, l, \text{loc} :: s \rangle \gg \rightarrow \ll h, \langle \text{pc}+1, l, o(f) :: s \rangle \gg}$$

The initial value \diamond is 0 for integers and null for object references.

Operational semantics

- ▶ How would the initial state look like?
- ▶ Implement a simple bytecode interpreter (in OCaml, Scala. . .)
- ▶ Is this semantics deterministic, and why is it so?
- ▶ How much this semantics is blocking?
- ▶ How to ensure that some of these blocking configurations will never be reached?

Other language features

Aspects of the language not dealt with here include:

- ▶ Methods calls
- ▶ Visibility modifiers (`public`, `private`, `protected`, ...)
- ▶ Multi-threading (`monitorenter`, `monitorexit`)
- ▶ Exceptions
- ▶ The constant pool
- ▶ Arrays
- ▶ Long integers