

# BEHAVIOR TYPES OF SEVERAL JVM INSTRUCTIONS

JACOPO FREDDI, CHUN TIAN, MICHAEL CANELLA, FABIO BISELLI, GIULIA BACCOLINI

**ABSTRACT.** Based on the work in previous papers, we define a way to compute behavioural types for most Java bytecode instructions, in a way that they may be handled by automated verifiers to compute cost functions.

## 1. CONVENTIONS AND MATHEMATICAL CONSTRUCTS

**1.1. Datatypes.** The datatypes we consider are  $\mathbb{N}^+ = \mathbb{C} \cup \mathbb{E} \cup \{-\}$  for integer values (augmented with an undefined value to manage expressions that the verifier cannot handle exactly), and  $VM = \{\alpha, \beta, \gamma, \dots\} \cup \{\mathbf{this}\}$  as a set of names for virtual machines.  $\mathbb{C}$  stands for the integer constants and is equal to  $\mathbb{N}$ , while  $\mathbb{E}$  is the set of all expressions in which variables are involved, but can be handled by the verifier.

$\Gamma \in \underline{\Gamma}$  is a map  $VM \mapsto \sigma = \{\delta, \top, \perp, a, a \downarrow\}$  which holds information about the state of each virtual machine. A general datatype  $\mathbb{D} = \mathbb{N}^+ \cup VM$  is defined to compactly integrate all datatypes stored in the stack and the memory. We also consider the

**1.2. Machine abstractions.** We consider an abstraction of the bytecode program  $P$  with  $n$  instructions as a function defined  $\forall i \in [1; n]$ . No checks on the domain of  $P$  are made because the bytecode is only valid if all execution path end with a return statement, and in that case the analysis do not consider exceding lines of code.

The other abstractions we consider are:

- a memory  $F \in \mathbb{F}$  seen as a map  $N \mapsto \mathbb{D}$ ;
- a stack  $S \in \mathbb{S} = \emptyset \cup (\mathbb{D} \times \mathbb{S})$ ;
- a constant table **Const** where **Const**( $i$ ) denotes the constant  $\#i$  of the byte-code unit.

The JVM state is abstracted with a memory  $F$  initially filled with the parameters of the method or program, and a stack frame  $S$  initially empty. On the stack  $S$  are defined the operations with the trivial semantics, given the convention that  $a \cdot b = \{a, b\}$ :

$$top(S) : \mathbb{S} \mapsto \mathbb{D}. top(S) = x \text{ if } S = x \cdot S' \in \mathbb{S}$$

$$pop(S) : \mathbb{S} \mapsto \mathbb{S}. pop(S) = S' \text{ if } S = x \cdot S' \in \mathbb{S}$$

$$push(x, S) : \mathbb{D} \times \mathbb{S} \mapsto \mathbb{S}. push(x, S) = x \cdot S$$

being  $top$  and  $pop$  both intentionally partial functions, i.e. not being defined on empty stacks.

## 2. THE GREATEST LOWER BOUND FUNCTION

Every symbol used in the stack, memory or environment is involved in a partial ordering relationship with the other symbols of its type. The Greatest Lower Bound of two symbols, when defined, computes the greatest of the symbols lower or equal to them. The function is *not* commutative: when the arguments are different but hold the same informative value, it usually gives priority to the first one. However, the ordering will only affect the value assigned to formal parameters and is not semantic.

2.1.  $\sqcup$  **on datatypes.** In general we consider a “lowest type”  $\perp$  such that:

$$\forall x \in \mathbb{D} : x \sqcup \perp = \perp \sqcup x = \perp$$

This type marks the “undefined” position and is useful to define the greatest lower bound for the memory, since the domain for it must be equal in both operands. If a greatest lower bound must be computed between  $F$  and  $F'$  where  $\exists i.F(i) \in \mathbb{D} \wedge \nexists F'(i)$  we can say that  $\exists F'(i) = \perp \wedge F(i) \sqcup F'(i) = \perp$ .

On integers:

$$\sqcup : \mathbb{N}^+ \times \mathbb{N}^+ \mapsto \mathbb{N}^+ : a \sqcup b = \begin{cases} - & a = - \vee b = - \\ b & a \in \mathbb{C} \wedge b \in \mathbb{E} \\ a & \text{else} \end{cases}$$

On virtual machines:

$$\sqcup : VM \times VM \mapsto VM : \alpha \sqcup \beta = \alpha$$

On virtual machine states:

$$\sqcup : \sigma \times \sigma \mapsto \sigma : \alpha \sqcup \beta = \begin{cases} \alpha & \alpha = \beta \\ \delta & (\alpha = \top \wedge \beta = \perp) \vee (\alpha = \perp \wedge \beta = \top) \vee (\alpha = \delta \vee \beta = \delta) \\ a \downarrow & (\alpha = a \wedge \beta = \perp) \vee (\alpha = \perp \wedge \beta = a) \vee (\alpha = a \downarrow \vee \beta = a \downarrow) \end{cases}$$

2.2.  $\sqcup$  **on the memory.**

$$\sqcup : \mathbb{F} \times \mathbb{F} \mapsto \mathbb{F}. F_1 \sqcup F_2 = \overline{F}. \forall i \in D(\overline{F}) : \overline{F}(i) = F_1(i) \sqcup F_2(i)$$

2.3.  $\sqcup$  **on the stack.**

$$\sqcup : \mathbb{S} \times \mathbb{S} \mapsto \mathbb{S}. S_1 \sqcup S_2 = \begin{cases} \emptyset & S_1 = S_2 = \emptyset \\ push(top(S_1) \sqcup top(S_2), pop(S_1) \sqcup pop(S_2)) & \text{else} \end{cases}$$

Notice that, being *top* and *pop* both partial functions, the function defined on the stack is not total: if there are some  $S_1, S_2$  such that  $\nexists S_1 \sqcup S_2$ , the stack lengths are different and the program cannot be analysed due to stack overflow errors.

2.4.  $\sqcup$  **on the environment.** We define a function  $D : \underline{\Gamma} \mapsto \wp(VM)$ :

$$D(\Gamma) = d. \forall \alpha \in VM : \Gamma(\alpha) \in \sigma \implies \alpha \in d$$

that computes the domain of the environment.

$$\sqcup : \underline{\Gamma} \times \underline{\Gamma} \mapsto \underline{\Gamma}. \Gamma_1 \sqcup \Gamma_2 = \overline{\Gamma}. D(\overline{\Gamma}) = D(\Gamma_1) \wedge \forall \alpha \in D(\overline{\Gamma}) : \overline{\Gamma}(\alpha) = \Gamma_1(\alpha) \sqcup \Gamma_2(\alpha)$$

**2.5. Definition of informative ordering.** We define a partial ordering function in  $\mathbb{N}^+$ :

$$\begin{aligned} <_{\Sigma} \subset \mathbb{N}^+ \times \mathbb{N}^+ : a <_{\Sigma} b \iff \begin{cases} a = - \wedge b \neq - \\ a \in \mathbb{E} \wedge b \in \mathbb{C} \end{cases} \\ =_{\Sigma} \subset \mathbb{D} \times \mathbb{D} : \sigma \times \sigma \cup \mathbb{C} \times \mathbb{C} \cup \mathbb{E} \times \mathbb{E} \cup \{(-, -)\} \end{aligned}$$

Starting from this function, we define a partial ordering between stacks and between memories:

$$\begin{aligned} <_{\Sigma} \subset \mathbb{F} \times \mathbb{F} : F_1 <_{\Sigma} F_2 &\iff (F_1 \leq_{\Sigma} F_2 \wedge \exists i \in D(F_1) : F_1(i) <_{\Sigma} F_2(i)) \\ <_{\Sigma} \subset \mathbb{S} \times \mathbb{S} : S_1 <_{\Sigma} S_2 &\iff (top(S_1) <_{\Sigma} top(S_2) \vee ((top(S_1) =_{\Sigma} top(S_2) \wedge pop(S_1) <_{\Sigma} pop(S_2))) \end{aligned}$$

The relevance of making stacks and memories comparable with respect to the information value will become clear in next section.

### 3. BEHAVIOURAL TYPE FUNCTION

The behavioural type of a section of  $P$  is given by the function  $\mathbb{B}_x : \mathbb{F} \times \mathbb{S} \mapsto \mathbb{B}$  where  $\mathbb{B}_i(F, S)$  types the program from the instruction  $i$  to the end of the program considering a memory  $F$  and a stack  $S$ . This way jumps are easily typed and branches are created in a natural way.

We consider that,  $\forall i : \exists P[i] \implies \Gamma_i, F_i, S_i \vdash \mathbb{B}_i(F_i, S_i) = b_i(F_i, S_i), R_i \triangleright \Gamma'$ , being:

- $b_i$  the behavioural type from instruction  $i$  to the end of the program, parametric with respect to the stack and the memory (usually written recursively);
- $R_i$  the set of virtual machines released by the instruction (usually an empty set or a singleton);
- $\Gamma'$  the new environment which will type the next  $\mathbb{B}_x$ .

Save in particular cases, the main changes happen on the stack and the memory: therefore, if not otherwise stated, we will omit  $\Gamma'$  and  $R_i$  assuming that:

- $R_i = \emptyset$ ;
- $\Gamma' = \Gamma$ .

The analysis of the program is performed in iterations of sequential scans, in which each  $\mathbb{B}_i$  is computed. Ideally, when a scan discovers that  $\mathbb{B}_i$  is defined on parameters  $F_i, S_i$  and invoked recursively with some  $F_j, S_j$  such that  $F_j <_{\Sigma} F_i \vee S_j <_{\Sigma} S_i$ , the domain of  $\mathbb{B}_i$  is updated from  $F_i, S_i$  to  $F_i \sqcup F_j, S_i \sqcup S_j$ , all the types are recomputed to mirror the update and another scan is issued to check that all invocations parameters are at least as informative as the definition. This analysis ends because:

- the number of instruction of a program or method is finite;
- the chains defined by the  $<_{\Sigma}$  relationship are all finite;
- the stack and the memory are limited by definition.

The analysis will then create, at the end of the last iteration, a set of  $\mathbb{B}_i$  defined for all instructions, that will be analysed by the theorem prover and converted into cost functions.

### 3.1. Open problems:

- even if computation from created VM is considered not to affect the VM environment, one can still have multiple machines or threads running the same code. While a machine may not be aware of the VMs created by another one, different threads may access the same global variables, therefore creating race conditions on the access of these variables. In the current version the program is supposed to be non-concurrent, but in the future a way to detect synchronized areas should be considered.

4. TYPING RULES (MUST BE REVIEWED TO REFLECT THE USAGE OF  $\Sigma$ )

(T-Program)

$$\overline{\Gamma, F, \emptyset \vdash P : \mathbb{B}_1(F, \emptyset)}$$

(T-Return)

$$\frac{P[i] = \mathbf{return}}{\Gamma \vdash \mathbb{B}_i(F_i, S_i) = 0}$$

(T-If)

$$\frac{P[i] = \mathbf{if} \ [cond] \ L, \text{top}(S_i) \neq -}{\Gamma \vdash \mathbb{B}_i(F_i, S_i) = [cond](\mathbb{B}_L(F_i, \text{pop}(S_i))) + [\neg cond](\mathbb{B}_{i+1}(F_i, \text{pop}(S_i)))}$$

(T-If-Undefined)

$$\frac{P[i] = \mathbf{if} \ [cond] \ L, \text{top}(S_i) = -}{\Gamma \vdash \mathbb{B}_i(F_i, S_i) = \mathbb{B}_L(F_i, \text{pop}(S_i)) + \mathbb{B}_{i+1}(F_i, \text{pop}(S_i))}$$

(T-Goto)

$$\frac{P[i] = \mathbf{goto} \ L}{\Gamma \vdash \mathbb{B}_i(F_i, S_i) = \mathbb{B}_L(F_i, S_i)}$$

(T-New-VM)

$$\frac{P[i] = \mathbf{invokevirtual} \ \mathbf{createVM}, \beta \text{ fresh}}{\Gamma \vdash \mathbb{B}_i(F_i, S_i) = \nu\beta \ ; \ \mathbb{B}_{i+1}(F_i, \text{push}(\beta, \text{pop}(S_i))) \triangleright \Gamma[\beta \mapsto \top]}$$

(T-Release-VM)

$$\frac{P[i] = \mathbf{invokevirtual} \ \mathbf{releaseVM}, \beta = \text{top}(S_i), \Gamma(\beta) \neq \perp}{\Gamma \vdash \mathbb{B}_i(F_i, S_i) = \beta^\vee \ ; \ \mathbb{B}_{i+1}(F_i, \text{pop}(\text{pop}(S_i))) \triangleright \Gamma[\beta \mapsto \perp], \{\beta\}}$$

(T-Release-VM-Null)

$$\frac{P[i] = \mathbf{invokevirtual} \ \mathbf{releaseVM}, \beta = \text{top}(S_i), \Gamma(\beta) = \perp}{\Gamma \vdash \mathbb{B}_i(F_i, S_i) = \mathbb{B}_{i+1}(F_i, \text{pop}(\text{pop}(S_i)))}$$

(T-Load)

$$\frac{P[i] = \mathbf{load} \ n}{\Gamma \vdash \mathbb{B}_i(F_i, S_i) = \mathbb{B}_{i+1}(F_i, \text{push}(F(n), S_i))}$$

(T-Store)

$$\frac{P[i] = \mathbf{store} \ n}{\Gamma \vdash \mathbb{B}_i(F_i, S_i) = \mathbb{B}_{i+1}(F_i[n \mapsto \text{top}(S_i)], \text{pop}(S_i))}$$

(T-Integer-Increment-Defined)

$$\frac{P[i] = \mathbf{iinc} \ idx \ n, F_i(idx) \neq -}{\Gamma \vdash \mathbb{B}_i(F_i, S_i) = \mathbb{B}_{i+1}(F_i[idx \mapsto (F_i(idx) + n)], S_i)}$$

(T-Integer-Increment-Undefined)

$$\frac{P[i] = \mathbf{iinc} \ idx \ n, F_i(idx) = -}{\Gamma \vdash \mathbb{B}_i(F_i, S_i) = \mathbb{B}_{i+1}(F_i, S_i)}$$

(T-Iconst)

$$\frac{P[i] = \mathbf{iconst\_}n}{\Gamma \vdash \mathbb{B}_i(F_i, S_i) = \mathbb{B}_{i+1}(F_i, \text{push}(n, S_i))}$$

(T-LDC)

$$\frac{P[i] = \mathbf{l dc } x, \mathbf{Const}(x) \in \mathbb{C}}{\Gamma \vdash \mathbb{B}_i(F_i, S_i) = \mathbb{B}_{i+1}(F_i, \text{push}(n, S_i))}$$

(T-LDC-Undefined)

$$\frac{P[i] = \mathbf{l dc } x, \mathbf{Const}(x) = -}{\Gamma \vdash \mathbb{B}_i(F_i, S_i) = \mathbb{B}_{i+1}(F_i, \text{push}(-, S_i))}$$

More on the way...