
A Type System for the JVM

Motivation: Security

Scenarios

- Web browser downloads and executes JVM applet.
In your own address space.
- SmartCard applet (one of many!) is executed in shared address space.

The problem

Should this bytecode be allowed:

..., Push (Addr 217), Push v, Putfield F D, ...

No! Must not allow access to arbitrary addresses.

Why: no hardware protection, software checks expensive.

⇒ **no Addresses in bytecode!**

How about

..., Push (Intg 217), Push v, Putfield F D, ...

No! Semantics underdefined, anything may happen, avoid!

The solution

Aim: Security without runtime checks (eg no type tags).

Method: Check statically against all kinds of attacks:
address spoofing, stack over/underflow, . . .

JVM needs proper type system
to detect malicious or erroneous bytecode

Overview

1. Define type system (declaratively)
2. Prove type safety:
When executing well-typed programs,
check-instr is always true when *exec-instr* is called.
3. Implement type system as a static analyzer
("bytecode verifier", later)

Jinja versus JVM

- Jinja: local variables + result
JVM: local variables + stack
- Jinja: type of local variables is declared and fixed
JVM: type of local variables is not declared and may change from instruction to instruction; same for stack.

The JVM type system describes for each program point the types in the local variables and on the stack.

Example

Code	Stack	Registers
Push (<i>Intg 1</i>)	[]	[<i>Class C, Integer</i>]
Load <i>1</i>		
IAdd		
Store <i>0</i>		

Motto

Types are abstract descriptions of values/states

Instructions (operating on values)
can be abstracted to operate on types

“Abstract interpretation”

Restriction

We treat only the exception-free fragment of the JVM
No system exceptions, no user exception

Type System

Stack type

- $ty_s = ty\ list$ $(ST :: ty_s)$

At each program point,
the Jinja type of each stack element
must be uniquely defined.

Register types

datatype *'a err = Err | OK 'a*

- *ty_l = ty err list* (*LT :: ty_l*)

For a given program point and a given register:

Err : the Jinja type of the register is not defined uniquely.
Not an error, but register cannot be read (is “unusable”).

OK T : the Jinja type of the register is definitely *T*.

Example: OK vs Err

Code	Stack	Registers
IfFalse 4	[<i>Boolean</i>]	[<i>OK (Class C)</i>]
Push (<i>Intg 1</i>)		
Store 0		
Goto 3		
Push (<i>Bool True</i>)		
Store 0		

Reasons for Err

- Uninitialized local variables
- Storage optimization:
two distinct local variables share the same register

Example: `if (b) { $V_1:T_1$; _ } else { $V_2:T_2$; _ }`

Instruction/state type, method type

- $ty_i = ty_s \times ty_l \quad (\tau :: ty_i)$
- Instructions can be unreachable:
 $ty_i' = ty_i \text{ option} \quad (\tau :: ty_i')$
None : instruction is unreachable
 $\lfloor \tau \rfloor$: state before/after instruction has type τ .
- Method type: $ty_i' \text{ list} \quad (\tau s :: ty_i' \text{ list})$

An example of the full truth

Code	Stack Registers
Goto 2	[([<i>Boolean</i>], [<i>Err</i>])]
Pop	
Push (<i>Bool True</i>)	
Store 0	

Well-typed Method

A fixed context

We consider a fixed method in a fixed JVM program P :

C_M	:: <i>cname</i>	the class the method is defined in
Ts	:: <i>ty list</i>	parameter types
T_r	:: <i>ty</i>	return type
mxs	:: <i>nat</i>	maximum stack size
mxl_0	:: <i>nat</i>	number of local variables (w/o parameters)
mxl	:: <i>nat</i>	number of registers ($mxl = 1 + Ts + mxl_0$)
is	:: <i>instr list</i>	instructions

Well-typed method (informal)

The method is well-typed iff there is a method type τs such that

- $|\tau s| = |is|$
- Each instruction is *applicable* to the corresponding instruction type in τs .
Eg: no `POP` when $ST = []$.
- The state after the execution of the instruction is “compatible” with all successor positions.
Eg: if $ST = [T]$ before `POP` then $ST = []$ afterwards.

Well-typed method (almost formal)

⋮	
$pc: l$	τ
⋮	
$pc': l'$	τ'
⋮	

For every pc we must have

- $app\ l\ \tau$ is true
- for every successor pc' of pc : $eff\ l\ \tau \sqsubseteq \tau'$

Applicability (wrt types!)

$app_i :: instr \Rightarrow ty_i \Rightarrow bool$

$app_i (\text{Load } n) (ST, LT) = (n < |LT| \wedge LT_{[n]} \neq Err \wedge |ST| < mxs)$

$app_i (\text{Store } n) (T \cdot ST, LT) = (n < |LT|)$

$app_i (\text{Push } v) (ST, LT) = (|ST| < mxs \wedge \text{typeof } v \neq None)$

$app_i \text{Pop} (T \cdot ST, LT) = True$

$app_i (\text{Getfield } F D) (T \cdot ST, LT) =$

$(\exists T_f. P \vdash D \text{ sees } F:T_f \text{ in } D \wedge P \vdash T \leq \text{Class } D)$

$app_i (\text{Putfield } F D) (T_1 \cdot T_2 \cdot ST, LT) =$

$(\exists T_f. P \vdash D \text{ sees } F:T_f \text{ in } D \wedge P \vdash T_2 \leq \text{Class } D \wedge P \vdash T_1 \leq T_f)$

$app_i (\text{New } C) (ST, LT) = (\text{is-class } P C \wedge |ST| < mxs)$

$app_i (\text{Checkcast } C) (T \cdot ST, LT) = (\text{is-class } P C \wedge \text{is-refT } T)$

Applicability (wrt types!)

$app_i \text{ IAdd } (Integer \cdot Integer \cdot ST, LT) = True$

$app_i \text{ CmpEq } (T_1 \cdot T_2 \cdot ST, LT) = (T_1 = T_2 \vee is-refT T_1 \wedge is-refT T_2)$

$app_i \text{ Goto } b (ST, LT) = True$

$app_i \text{ IfFalse } b (Boolean \cdot ST, LT) = True$

$app_i \text{ Return } (T \cdot ST, LT) = (P \vdash T \leq T_r)$

Applicability (wrt types!)

app_i (Invoke $M\ n$) (ST, LT) =

$n < |ST| \wedge$

$(ST_{[n]} \neq NT \longrightarrow$

$(\exists C\ D\ Ts\ T\ m.$

$ST_{[n]} = \text{Class } C \wedge P \vdash C \text{ sees } M: Ts \rightarrow T = m \text{ in } D \wedge$

$P \vdash \text{rev}(\text{take } n\ ST) [\leq] Ts))$

In all other cases:

$app_i\ i\ (ST, LT) = \text{False}$

Effect (wrt types!)

$eff_i :: instr \Rightarrow ty_i \Rightarrow ty_i$

$eff_i (\text{Load } n) (ST, LT) = (\text{let } OK\ T = LT_{[n]} \text{ in } (T \cdot ST, LT))$

$eff_i (\text{Store } n) (T \cdot ST, LT) = (ST, LT[n := OK\ T])$

$eff_i (\text{Push } v) (ST, LT) = (\text{let } [T] = \text{typeof } v \text{ in } (T \cdot ST, LT))$

$eff_i \text{Pop} (T \cdot ST, LT) = (ST, LT)$

$eff_i (\text{Getfield } FD) (T \cdot ST, LT) = (\text{snd } (\text{field } P\ D\ F) \cdot ST, LT)$

$eff_i (\text{Putfield } FD) (T_1 \cdot T_2 \cdot ST, LT) = (ST, LT)$

$eff_i (\text{New } C) (ST, LT) = (\text{Class } C \cdot ST, LT)$

$eff_i (\text{Checkcast } C) (T \cdot ST, LT) = (\text{Class } C \cdot ST, LT)$

Effect (wrt types!)

eff_i `IAdd` ($Integer \cdot Integer \cdot ST, LT$) = ($Integer \cdot ST, LT$)

eff_i `CmpEq` ($T_1 \cdot T_2 \cdot ST, LT$) = ($Boolean \cdot ST, LT$)

eff_i `Goto` b (ST, LT) = (ST, LT)

eff_i `IfFalse` b ($Boolean \cdot ST, LT$) = (ST, LT)

eff_i `Invoke` M n (ST, LT) =

let *Class* $C = ST_{[n]}$; (D, Ts, T_r, b) = *method* P C M *in*

($T_r \cdot drop$ ($n + 1$) ST, LT)

eff_i `Return` ($T \cdot ST, LT$) = don't care

Lifting applicability and effect

Remember: $ty_i' = ty_i$ option

$None :: ty_i'$ means “unreachable”

$app :: instr \Rightarrow ty_i' \Rightarrow bool$

$app\ instr\ None = True$

$app\ instr\ [\tau] = app_i\ instr\ \tau$

$eff :: instr \Rightarrow ty_i' \Rightarrow ty_i'$

$eff\ instr\ None = None$

$eff\ instr\ [\tau] = [eff_i\ instr\ \tau]$

Successor instructions

succs :: *instr* ⇒ *pc* ⇒ *pc*

succs (Goto *b*) *pc* = {*nat(int pc + b)*}

succs (IfFalse *b*) *pc* = {*pc + 1*, *nat(int pc + b)*}

succs Return *pc* = {}

In all other cases:

succs *i pc* = {*pc + 1*}

Orderings

\sqsubseteq $\tau \sqsubseteq \tau'$

τ describes a smaller set of JVM states than τ'

Examples:

- $([\mathbf{Class\ } C], []) \sqsubseteq ([\mathbf{Class\ } D], [])$ if $P \vdash C \preceq^* D$
- $([], [\mathbf{OK\ Integer}]) \sqsubseteq ([], [\mathbf{Err}])$

Construction of \sqsubseteq on $ty_i' = (ty\ list \times ty\ err\ list)option$:
by recursion over the type structure

Base case: \sqsubseteq on ty is subtype relation \leq

\sqsubseteq on pairs and lists

Given: $\sqsubseteq :: 'a \Rightarrow 'a \Rightarrow \text{bool}$

$\sqsubseteq :: 'b \Rightarrow 'b \Rightarrow \text{bool}$

On pairs: $\sqsubseteq :: 'a \times 'b \Rightarrow 'a \times 'b \Rightarrow \text{bool}$

$$(a,b) \sqsubseteq (a',b') \equiv (a \sqsubseteq b \wedge a' \sqsubseteq b')$$

On lists: $\sqsubseteq :: 'a \text{ list} \Rightarrow 'a \text{ list} \Rightarrow \text{bool}$

$$xs \sqsubseteq ys \equiv (|xs| = |ys| \wedge (\forall i < |xs|. xs_{[i]} \sqsubseteq ys_{[i]}))$$

\sqsubseteq on option

Given: $\sqsubseteq :: 'a \Rightarrow 'a \Rightarrow bool$

Extension: $\sqsubseteq :: 'a\ option \Rightarrow 'a\ option \Rightarrow bool$

$None \sqsubseteq x = True$

$[a] \sqsubseteq [b] = a \sqsubseteq b$

$\implies None$ is bottom element

⊆ on err

Given: $\underline{\subseteq} :: 'a \Rightarrow 'a \Rightarrow bool$

Extension: $\underline{\subseteq} :: 'a\ err \Rightarrow 'a\ err \Rightarrow bool$

$x \underline{\subseteq} Err = True$

$OK\ a \underline{\subseteq} OK\ b = a \underline{\subseteq} b$

$\implies Err$ is top element

Example

None \sqsubseteq $\lfloor ([\mathbf{Class\ C}], [OK\ Boolean, OK\ Integer]) \rfloor$
 \sqsubseteq $\lfloor ([\mathbf{Class\ C}], [OK\ Boolean, Err]) \rfloor$
 \sqsubseteq $\lfloor ([\mathbf{Class\ C}], [Err, Err]) \rfloor$
 \sqsubseteq $\lfloor ([\mathbf{Class\ Object}], [Err, Err]) \rfloor$

Remember the fixed method context

C_M	:: <i>cname</i>	the class the method is defined in
Ts	:: <i>ty list</i>	parameter types
T_r	:: <i>ty</i>	return type
mxs	:: <i>nat</i>	maximum stack size
mxl_0	:: <i>nat</i>	number of local variables (w/o parameters)
mxl	:: <i>nat</i>	number of registers ($mxl = 1 + Ts + mxl_0$)
is	:: <i>instr list</i>	instructions

Well-typed method

The method is **well-typed w.r.t. τs** iff

- $is \neq [] \wedge |\tau s| = |is|$
- $\forall ST LT. [(ST, LT)] \in set \tau s \longrightarrow |ST| < mxs \wedge |LT| = mxl$
- $[([], OK (Class C_M) \cdot map OK Ts @ replicate mxl_0 Err)] \sqsubseteq \tau S_{[0]}$
- $\forall p < |is|. \forall q \in succs is_{[p]} p. q < |is|$
- $\forall p < |is|. app is_{[p]} \tau S_{[p]}$
- $\forall p < |is|. \forall q \in succs is_{[p]} p. eff is_{[p]} \tau S_{[p]} \sqsubseteq \tau S_{[q]}$

Well-formed method

The method declaration $M : Ts \rightarrow T_r = (mxs, mxl_0, is, [])$ in class C_M is **well-formed w.r.t. τs** iff

- All parameter types are valid: $\forall T \in \text{set } Ts. \text{is-type } P T$
- The result type is valid: $\text{is-type } P T_r$
- The declaration is well-typed w.r.t. τs (as defined above).

Well-formed program

A JVM program P is well-formed w.r.t.

$\Phi :: \text{cname} \Rightarrow \text{mname} \Rightarrow \text{ty}_i' \text{ list}$

iff it is a well-formed Jinja program where each declaration of a method M in class C_M is well-formed w.r.t. $\Phi C_M M$ (as defined above, not as in Jinja!).

Type safety of the JVM

Conformance

JVM cannot get stuck but *check-instr* can be violated.

Let P be well-formed w.r.t. Φ .

1. Define conformance between JVM state and P, Φ .
2. Show that execution of P preserves conformance.
3. Show that if the state conforms to P, Φ , *check-instr* is true.

Conformance of stack and registers

- Stack:

$$P, h \vdash [v_1, \dots, v_m] [:\leq] [T_1, \dots, T_n] \equiv \\ m = n \wedge P, h \vdash v_1 : \leq T_1 \wedge \dots \wedge P, h \vdash v_n : \leq T_n$$

- Registers:

$$P, h \vdash [v_1, \dots, v_m] [:\leq_T] [U_1, \dots, U_n] \equiv \\ m = n \wedge P, h \vdash v_1 : \leq_T U_1 \wedge \dots \wedge P, h \vdash v_n : \leq_T U_n$$

where

$$P, h \vdash v : \leq_T \text{Err} \\ (P, h \vdash v : \leq_T \text{OK } T) = (P, h \vdash v : \leq T)$$

Conformance

- Frame (stk, loc, C, M, pc) conforms to P, Φ :
 - $P \vdash C$ sees $M : \dots \rightarrow \dots = (\dots, is, \dots)$ in C
 - $(\Phi \ C \ M)_{[pc]} = \lfloor (ST, LT) \rfloor$
 - $P, h \vdash stk \ [:\leq] \ ST \ \wedge \ P, h \vdash loc \ [:\leq_{\top}] \ LT$
 - if it is not the top frame: $is_{[pc]} = \text{Invoke } M_0 \ n_0$
and M_0 is the method in the frame above.
- JVM state:
 $P, \Phi \vdash (_, h, frs) \checkmark$ iff
 $P \vdash h \checkmark$ and every frame in frs conforms to P, Φ .

Type safety

Thm (Preservation of conformance) If P is well-formed w.r.t. Φ and $P, \Phi \vdash \sigma \checkmark$ and $\text{exec } P \sigma = [\sigma']$ then $P, \Phi \vdash \sigma' \checkmark$.
Proof by case distinction over the instructions.

Thm If P is well-formed w.r.t. Φ and $P, \Phi \vdash \sigma \checkmark$ and the computation of $\text{exec } P \sigma$ calls *exec-instr* ... then *check-instr* ... (with the same parameters) is true.
Proof by case distinction over the instructions.