

# Defining the Java Virtual Machine as Platform for Provably Correct Java Compilation

Egon Börger<sup>1</sup> and Wolfram Schulte<sup>2</sup>

<sup>1</sup> Università di Pisa, Dipartimento di Informatica, I-56125 Pisa, Italy  
`boerger@di.unipi.it`

<sup>2</sup> Universität Ulm, Fakultät für Informatik, D-89069 Ulm, Germany  
`wolfram@informatik.uni-ulm.de`

**Abstract.** We provide concise abstract code for running the Java Virtual Machine (JVM) to execute compiled Java programs, and define a general compilation scheme of Java programs to JVM code. These definitions, together with the definition of an abstract interpreter of Java programs given in our previous work [3], allow us to prove that any compiler that satisfies the conditions stated in this paper compiles Java code correctly. In addition we have validated our JVM and compiler specification through experimentation.

The modularity of our definitions for Java, the JVM and the compilation scheme exhibit orthogonal language, machine and compiler components, which fit together and provide the basis for a stepwise and provably correct *design-for-reuse*. As a by-product we provide a challenging realistic case study for mechanical verification of a compiler correctness proof.

## 1 Introduction

Every justification showing that a proposed compiler behaves well is relative to a definition of the semantics of source and target language. In our previous work [3] we have developed a platform independent, rigorous yet easily manageable definition for an interpreter of Java programs, which captures the intuitive understanding Java programmers have of the semantics of their code. In this paper we provide a mathematical (read: rigorous and platform independent) yet practical model of an interpreter for the Java Virtual Machine, which formalizes the concepts presented in the JVM specification [6], as far as they are needed for the compilation of Java programs. We also extract from the JVM specification the definition of a scheme for the compilation of Java to JVM code and prove its correctness.

**Main Theorem.** *Every compiler that satisfies the conditions listed in this paper compiles Java programs correctly into JVM code.*

We split the JVM and the compilation function into an incremental sequence of four machines and functions—whose structure corresponds to the conservative extension relation among the modular components we exhibited for Java [3]—and define the JVM at two levels of abstraction: a ground model with an abstract

class file and abstract instructions, and a refined model where the abstract instructions are implemented by concrete JVM instructions. The structure of our Java machine is carried over *mutatis mutandis* to the basic structure of the abstract interpreter we are defining here for the JVM as target machine for Java compilation.

In sections 2 to 5 we define the sequence of successively extended JVM machines  $JVM_I$ ,  $JVM_C$ ,  $JVM_O$  and  $JVM_E$  for the compilation of programs from the imperative core  $Java_I$  of Java and its extensions  $Java_C$  (by classes, eg. procedures),  $Java_O$  (by object-oriented features, eg. class instances) and  $Java_E$  (by exceptions). We discuss here only the single threaded JVM, although our approach could easily include also multiple threads (see our multi-agent Java model with threads in [3]). We skip those language constructs which can be reduced by standard program transformation techniques to the core constructs dealt with explicitly in our Java models. We still do not consider Java packages, compilation units, visibility of names, strings, arrays, input/output, loading, linking and garbage collection. These features are the object of further refinements of the JVM model presented here. For proof details, the instruction refinement, an extensive bibliography and the discussion of related work we refer the interested reader to an extended version of this paper [1].

## 2 $JVM_I$ and the Compilation of $Java_I$ Programs

For the specification of Java, the JVM and the proof machinery, we use Abstract State Machines (ASMs). ASM specifications have a simple mathematical foundation [5], which justifies their intuitive understanding as “pseudo code” over abstract data. We define the basic JVM, called  $JVM_I$ , which is used as the target for compiling Java’s statements and expressions over primitive types. We prove that  $JVM_I$  executes the compilation of  $Java_I$  programs correctly.

The following grammars recall the syntax of  $Java_I$  [3] and introduce the corresponding instruction set  $JVM_I$ :

$  \begin{array}{l}  \textit{Exp} ::= \textit{Lit} \\  \quad   \quad \textit{Uop Exp} \\  \quad   \quad \textit{Exp Bop Exp} \\  \quad   \quad \textit{Var} \\  \quad   \quad \textit{Var} = \textit{Exp} \\  \quad   \quad \textit{Exp? Exp: Exp} \\  \\  \textit{Stm} ::= ; \\  \quad   \quad \textit{Exp}; \\  \quad   \quad \textit{Lab} : \textit{Stm} \\  \quad   \quad \textbf{break} \textit{Lab}; \\  \quad   \quad \textbf{continue} \textit{Lab}; \\  \quad   \quad \textbf{if} (\textit{Exp}) \textit{Stm} \textbf{else} \textit{Stm} \\  \quad   \quad \textbf{while} (\textit{Exp}) \textit{Stm} \\  \quad   \quad \{ \textit{Stm}^* \}  \end{array}  $	$  \begin{array}{l}  \textit{Instr} ::= \textbf{const} (\textit{Lit}) \\  \quad   \quad \textbf{uapply} (\textit{Uop}) \\  \quad   \quad \textbf{bapply} (\textit{Bop}) \\  \quad   \quad \textbf{load} (\textit{Varnum} \times \textit{Typ}) \\  \quad   \quad \textbf{store} (\textit{Varnum} \times \textit{Typ}) \\  \quad   \quad \textbf{dup} (\textit{Typ}) \\  \quad   \quad \textbf{pop} (\textit{Typ}) \\  \quad   \quad \textbf{ifZero} (\textit{Lab}) \\  \quad   \quad \textbf{goto} (\textit{Lab}) \\  \quad   \quad \textbf{label} (\textit{Lab}) \\  \\  \textit{Varnum} == \textit{Nat} \\  \textit{Code} == \textit{Instr}^*  \end{array}  $
---	--

The JVM<sub>*I*</sub> instruction set bears a close resemblance to a traditional stack machine like the P-machine. JVM<sub>*I*</sub> provides instructions to load constants, to apply various unary and binary operators, to load and store a variable, to duplicate and to remove values, and to jump unconditionally or conditionally to a label. Variable locations in the JVM are represented by natural numbers. A JVM<sub>*I*</sub> program is a sequence of instructions.

The universes *Lit*, *Uop*, *Bop*, *Var*, *Typ*, *Lab* contain Java literals, unary and binary operators, variables, primitive types and labels, respectively. With the exception of *Var*, these universes are also used in the JVM<sub>*I*</sub>.

## 2.1 The Machine JVM<sub>*I*</sub> for Imperative Code

The JVM is a typed word-oriented stack-machine running the given bytecode *code* : *Code*. As a consequence the central dynamic part of a JVM<sub>*I*</sub> state consists of a program counter *pc*, a local variable environment *loc* and an operand stack *opd*. The following declarations show their formalization: the first column defines the used types, the second column defines the state, and the third column defines the condition on the initial state. (We consider sequences as isomorphic to functions having an interval of natural numbers starting at 0 as their domain.)

$Pc == Nat$	$pc : Pc$	$pc = next_{unlab}(0, code)$
$Loc == Varnum \rightarrow Word$	$loc : Loc$	$loc = \emptyset$
$Opd == Word^*$	$opd : Opd$	$opd = \epsilon$

The close analogy between the abstract and concrete program counters in Java<sub>*I*</sub> and JVM<sub>*I*</sub>, the memories for local variables and for intermediate values, and their initializations reflects the refinement process, which applied to the machine Java<sub>*I*</sub> yields JVM<sub>*I*</sub>. This correspondence will guide the justification of the correctness of this first step towards an implementation of Java on the JVM.

Local variables and the operand stack store values of the abstract universe *Word*. *Words* are supposed to hold at least 32-bit quantities. Java's values, which occupy at most 32-bits, are represented on the level of the JVM as single *Words*. Java's 64-bit values are mapped to multiple consecutive locations in the local environment and on the operand stack in an implementation dependent way. We define JVM values (*Val*) as sequences of *Words*, i.e.  $Val == Word^*$ . A valid word sequence has length one (32-bit) or two (64-bit). The JVM implements values and operations on Java datatypes as follows. Booleans are represented as integers: 0 is used for *false*, and 1 for *true*. Operations working on boolean, byte, short or char are not supported by the JVM. Instead, upon retrieving the value of a boolean, byte, char or short, it is automatically cast into an int. When writing a value to a boolean, byte, char or short variable, an int is passed and the JVM truncates it to the relevant size.

For the JVM<sub>*I*</sub> we use two static code traversing functions *next* and *jump*, which yield the next statement to be executed and the next statement after the given labeled statement, respectively. Both functions are defined using an aux-

iliary function  $next_{unlab}$  that skips `label` instructions. (The expression  $\iota x \mid p(x)$  denotes the uniquely determined object  $x$  that satisfies  $p(x)$ .)

$$\begin{aligned} next(pc, code) &= next_{unlab}(pc + 1, code) \\ jump(l, code) &= next_{unlab}(\iota pc \mid code(pc) = \text{label}(l)) \\ next_{unlab}(pc, code) &= \min\{pc' \mid pc' \geq pc \wedge \forall l \mid code(pc') \neq \text{label}(l)\} \end{aligned}$$

We also use the following JVM<sub>T</sub> macros, where the homonymy with Java<sub>T</sub> macros reflects the refinement relations on which our correctness proof is based.

$$\begin{aligned} proceed &== pc := next(pc, code) \\ goto(l) &== pc := jump(l, code) \\ pc \text{ is } instr &== code(pc) = instr \end{aligned}$$

The following rules define the semantics of the JVM<sub>T</sub> instructions.

<p><b>if</b> <math>pc</math> <b>is</b> <code>const</code> (<math>lit</math>)  <b>then</b>  <math>opd := \tilde{lit} \cdot opd</math>  <math>proceed</math></p> <p><b>if</b> <math>pc</math> <b>is</b> <code>uapply</code> (<math>\odot</math>)<math>\wedge</math>  <math>(v, opd') = split(\mathcal{A}(\odot), opd)</math>  <b>then</b>  <math>opd := \tilde{\odot} v \cdot opd'</math>  <math>proceed</math></p> <p><b>if</b> <math>pc</math> <b>is</b> <code>bapply</code> (<math>\otimes</math>)<math>\wedge</math>  <math>(v_2, v_1, opd') = split(\mathcal{A}(\otimes), opd) \wedge</math>  <math>(\otimes \in DivMods) \Rightarrow (v_2 \neq 0)</math>  <b>then</b>  <math>opd := v_1 \tilde{\otimes} v_2 \cdot opd'</math>  <math>proceed</math></p> <p><b>if</b> <math>pc</math> <b>is</b> <code>dup</code> (<math>t</math>)<math>\wedge</math>  <math>(v, opd') = split(t, opd)</math>  <b>then</b>  <math>opd := v \cdot v \cdot opd'</math>  <math>proceed</math></p> <p><b>if</b> <math>pc</math> <b>is</b> <code>pop</code> (<math>t</math>)<math>\wedge</math>  <math>(v, opd') = split(t, opd)</math>  <b>then</b>  <math>opd := opd'</math>  <math>proceed</math></p>	<p><b>if</b> <math>pc</math> <b>is</b> <code>load</code> (<math>x, t</math>)  <b>then</b>  <b>if</b> <math>sizeof(t) = 1</math> <b>then</b>  <math>opd := loc(x) \cdot opd</math>  <b>else if</b> <math>sizeof(t) = 2</math> <b>then</b>  <math>opd := loc(x) \cdot loc(x + 1) \cdot opd</math>  <math>proceed</math></p> <p><b>if</b> <math>pc</math> <b>is</b> <code>store</code> (<math>x, t</math>)<math>\wedge</math>  <math>(v, opd') = split(t, opd)</math>  <b>then</b>  <math>opd := opd'</math>  <b>if</b> <math>sizeof(t) = 1</math> <b>then</b>  <math>loc(x) := v(0)</math>  <b>else if</b> <math>sizeof(t) = 2</math> <b>then</b>  <math>loc(x + 1) := v(1)</math>  <math>loc(x) := v(0)</math>  <math>proceed</math></p> <p><b>if</b> <math>pc</math> <b>is</b> <code>goto</code> (<math>l</math>) <b>then</b>  <math>goto(l)</math></p> <p><b>if</b> <math>pc</math> <b>is</b> <code>ifZero</code> (<math>l</math>)<math>\wedge</math>  <math>w \cdot opd' = opd</math>  <b>then</b>  <math>opd := opd'</math>  <b>if</b> <math>w = 0</math> <b>then</b> <math>goto(l)</math>  <b>else</b> <math>proceed</math></p>
---	---

A `const` instruction pushes the JVM value  $\tilde{lit}$  (one or two words) on the operand stack. An unary (binary) operator changes the value(s) on top of the operand stack. The unary (binary) operators are assumed to have the same meaning as in Java (i.e.  $\tilde{\odot}$  ( $\tilde{\otimes}$ )), although they may operate on extended domains. In order to abstract from the different value sizes, we use the function  $split : (Typ^*, Opd) \rightarrow (Val^*, Opd)$ , which given a sequence of  $n$  types and the operand stack, takes the top  $n$  values from the operand stack, such that the  $i$ th value has the size

of the  $i$ th type. The function  $\mathcal{A}(op)$  returns the argument types of  $op$ . The instructions **dup** and **pop** duplicate and remove the top stack value, respectively. A **load** instruction loads the value stored under the location  $x$  on top of the stack. If the type of  $x$  is a double or long, the next two locations are pushed on top of the stack. A **store** instruction stores the top (two) word(s) of the operand stack in the local environment at offset  $x$  (and  $x + 1$ ). A **goto** instruction causes execution to jump to the next instruction determined by the label. The **ifzero** instruction is a conditional **goto**. If the value on top of the operand stack is 0, execution continues at the next instruction determined by the label, otherwise execution proceeds.

The abstract nature of the  $\text{JVM}_{\mathcal{I}}$  instructions is reflected in their parameterization by types and operators. It allows us to restrict our attention to a small set of JVM instructions (or better instruction classes) without losing the generality of our model with respect to the JVM specification [6]. The extended version of this paper [1] shows how to refine these parameterized instruction to JVM's real ones.

## 2.2 Compilation of $\text{Java}_{\mathcal{I}}$ Programs to $\text{JVM}_{\mathcal{I}}$ Code

This section defines the compiling function from  $\text{Java}_{\mathcal{I}}$  to  $\text{JVM}_{\mathcal{I}}$  code. More efficient compilation schemes can be introduced but we leave optimizations for further refinement steps.

The compilation  $\mathcal{E} : \text{Exp} \rightarrow \text{Code}$  of (occurrences of)  $\text{Java}_{\mathcal{I}}$  expressions to  $\text{JVM}_{\mathcal{I}}$  instructions is standard. The resulting sequence of instructions has the effect of storing the value of the expression on top of the operand stack.

To improve readability, we use the following conventions for the presentation of the compilation: We suppress the routine machinery for a consistent assignment of (occurrences of) Java variables  $x$  to JVM variable numbers  $\bar{x}$ . Similarly, we suppress the trivial machinery for label generation. Label providing functions  $lab_i$ ,  $i \in \text{Nat}$ , are defined on occurrences of expressions and statements, are supposed to be injective and to have disjoint ranges. Functions  $\mathcal{T}$  defined on occurrences of variables and expressions return their type. We abbreviate: ‘Let  $e$  be an occurrence of  $exp$  in  $\mathcal{E}(e) = \dots$ ’ by ‘ $\mathcal{E}(e \text{ as } exp) = \dots$ ’.

$$\begin{aligned}
 \mathcal{E}(lit) &= \text{const}(lit) \\
 \mathcal{E}(\odot e) &= \mathcal{E}e \cdot \text{uapply}(\odot) \\
 \mathcal{E}(e_1 \otimes e_2) &= \mathcal{E}e_1 \cdot \mathcal{E}e_2 \cdot \text{bapply}(\otimes) \\
 \mathcal{E}(x) &= \text{load}(\bar{x}, \mathcal{T}(x)) \\
 \mathcal{E}(x = e) &= \mathcal{E}e \cdot \text{dup}(\mathcal{T}(e)) \cdot \text{store}(\bar{x}, \mathcal{T}(x)) \\
 \mathcal{E}(e \text{ as } e_1? e_2: e_3:) &= \mathcal{E}e_1 \cdot \text{ifZero}(lab_1(e)) \cdot \\
 &\quad \mathcal{E}e_2 \cdot \text{goto}(lab_2(e)) \cdot \text{label}(lab_1(e)) \cdot \mathcal{E}e_3 \cdot \text{label}(lab_2(e))
 \end{aligned}$$

Also the compilation  $\mathcal{S} : \text{Stm} \rightarrow \text{Code}$  of  $\text{Java}_{\mathcal{I}}$  statements to  $\text{JVM}_{\mathcal{I}}$  instructions is standard. The compilation of **break**  $lab$ ; and **continue**  $lab$ ; uses the auxiliary function  $target : \text{Stm} \times \text{Lab} \rightarrow \text{Stm}$ . This function provides for occur-

rences of statements and labels the occurrence of the enclosing labeled statement in the given program.

$$\begin{aligned}
S(;) &= \epsilon \\
S(e;) &= \mathcal{E}e \cdot \text{pop}(T(e)) \\
S(\{s_1 \dots s_m\}) &= Ss_1 \cdot \dots \cdot Ss_m \\
\\
S(s \text{ as if } (e) s_1 \text{ else } s_2) &= \mathcal{E}e \cdot \text{ifZero}(lab_1(s)) \cdot \\
&\quad Ss_1 \cdot \text{goto}(lab_2(s)) \cdot \text{label}(lab_1(s)) \cdot Ss_2 \cdot \text{label}(lab_2(s)) \\
S(s \text{ as while } (e) s_1) &= \text{label}(lab_1(s)) \cdot \mathcal{E}e \cdot \text{ifZero}(lab_2(s)) \cdot \\
&\quad Ss_1 \cdot \text{goto}(lab_1(s)) \cdot \text{label}(lab_2(s)) \\
S(s \text{ as lab : } s_1) &= \text{label}(lab_1(s)) \cdot Ss_1 \cdot \text{label}(lab_2(s)) \\
S(s \text{ as continue lab;}) &= \text{goto}(lab_1(\text{target}(s, \text{lab}))) \\
S(s \text{ as break lab;}) &= \text{goto}(lab_2(\text{target}(s, \text{lab})))
\end{aligned}$$

*Correctness Theorem for Java<sub>I</sub>/JVM<sub>I</sub>.* Via the refinement relation and under the assumptions stated above, the result of executing any Java<sub>I</sub> program in the machine Java<sub>I</sub> is equivalent to the result of executing the compiled program on the machine JVM<sub>I</sub>.

### 3 JVM<sub>C</sub> and the Compilation of Class Code

In this section we extend the basic JVM<sub>I</sub> machine to the machine JVM<sub>C</sub>, which handles class (also called static) fields, class methods and class initializers. JVM<sub>C</sub> thus stands for a machine that supports modules, module-local variables and procedures. We add the clauses for compiling class field access, class field assignment, class method calls and return statements to the definition of the Java<sub>I</sub> compilation function.

The following grammar shows the extension of the syntax of Java<sub>I</sub> to the syntax of Java<sub>C</sub>. Furthermore, we define the corresponding JVM<sub>C</sub> instructions:

$$\begin{aligned}
Exp &::= \dots & Instr &::= \dots \\
&| \text{FieldSpec} &&| \text{getstatic}(\text{FieldSpec} \times \\
&| \text{FieldSpec} = Exp && \quad Typ) \\
&| \text{MethSpec}(Exp^*) &&| \text{putstatic}(\text{FieldSpec} \times \\
&&& \quad Typ) \\
\\
Stm &::= \dots &&| \text{invokestatic}(\text{MethSpec}) \\
&| \text{return;} &&| \text{return}(Typ) \\
&| \text{return } Exp; && \\
\\
Init &::= \text{static } Stm && Fcty == (Typ^* \times Typ) \\
&&& FieldSpec == (Class \times Field) \\
&&& MethSpec == (Class \times Meth \times Fcty)
\end{aligned}$$

JVM<sub>C</sub> provides instructions to load and store class fields, and to call and to return from class methods. Both grammars are based on the same abstract definition of field and method specifications. Field specifications consist of a class and a field name, because Java and the JVM allow fields in different classes to have the same name. Method specifications additionally have a functionality (a

sequence of argument types and a result type, which can be `void`), because Java and the JVM support classes with methods having the same name but taking different parameter types.

Field and method specifications use the abstract universes *Class*, *Field* and *Method*. *Class* is assumed to stand for fully qualified Java class names, *Field* and *Method* for identifiers.

### 3.1 The Machine JVM<sub>C</sub> for Class Code

JVM and Java programs are structured into classes, which establish the program's execution environment. For a general, high-level definition of a provably correct compilation scheme from Java to JVM Code, we can abstract from many data structure specifics of the particular JVM class format. This format is called class file in the JVM specification [6].

Our abstract class file refines in a natural way the class environment of Java<sub>C</sub>, providing for every class its kind (whether it is a class or an interface), its superclass (if there is any), a list of the interfaces the class implements, and a table for fields and methods. Class files do not include definitions for fields or methods provided by any superclass.

$$\begin{aligned} Env &== Class \rightarrow ClassDec \\ ClassKind &::= AClass \mid AnInterface \\ ClassDec &== (kind : ClassKind \times super : [Class] \times ifaces : Class^* \times \\ &\quad fTab : Field \rightarrow FieldDec \times mTab : (Meth \times Fcty) \rightarrow MethDec) \end{aligned}$$

In JVM<sub>C</sub> fields and methods can only be static. Fields have a type and optionally a constant value. If a method is implemented in the class, the method body defines its code.

$$\begin{aligned} FieldDec &== (fKind : MemberKind \times fTyp : Typ \times fConstVal : [Val]) \\ MethDec &== (mKind : MemberKind \times mBody : [Code]) \\ MemberKind &::= Static \end{aligned}$$

In JVM<sub>C</sub> we have a fixed environment  $env : Env$ , defined by the given program. The following functions operate on this environment. The function  $mCode$  retrieves for a given method specification the method's code to be executed. The function  $fInitVal$  yields for a given field specification the field's constant value, provided it is available; otherwise, the function returns the default value of the field's type (where  $default : Typ \rightarrow Val$ ).

$$\begin{aligned} mCode(c, m, f) &= mBody(mTab(env(c))(m, f)) \\ fInitVal(c, f) &= \text{case } fTab(env(c))(f) \text{ of } (\_, \_, val) : val \\ &\quad (\_, fTyp, []) : default(fTyp) \end{aligned}$$

The function  $supers$  calculates the transitive closure of  $super$ . The function  $cfields$  returns the set of all fields declared by the class.

$$\begin{aligned} supers &: Class \rightarrow Class^* \\ cfields &: Class \rightarrow \mathcal{P} FieldSpec \end{aligned}$$

For these functions the homonymy to Java<sub>C</sub> functions shows the data refinement relation in going from Java<sub>C</sub> to JVM<sub>C</sub>.

Due to the presence of method calls in JVM<sub>C</sub> we have to embed the one single JVM<sub>I</sub> frame  $(pc, loc, opd)$  into the JVM<sub>C</sub> frame stack *frames*, enriched by a fourth component which always holds the dynamic chain of method specifications. This embedding defines the refinement relation between JVM<sub>I</sub> and JVM<sub>C</sub>. We refine the static function *code*, so that it always denotes the code stored in the environment under the current method specification *mspec*. The current class, method and functionality are denoted by *cclass*, *cmeth* and *cfcty*, respectively, where  $mspec = (cclass, cmeth, cfcty)$ .

$$\begin{array}{llll}
 pcs & : & Pc^* & pc == top(pcs) \\
 locs & : & Loc^* & loc == top(locs) \\
 opds & : & Opd^* & opd == top(opds) \\
 mspecs & : & MethSpec^* & mspec == top(mspecs) \\
 frames & == & (pcs, locs, opds, mspecs) & code == mCode(mspec)
 \end{array}$$

Before a class can be used its class initializers must be executed. At the JVM level class initializers appear as class methods with the special name `<clinit>`. Initialization must be done lazily, i.e. when a class is *first used* in Java, and when a reference is *resolved* in the JVM. Resolution is the process of checking symbolic references from the current class to other classes and interfaces. Since Java's notion of class initialization does not correspond to the related class resolution notion of the JVM, we name the initialization related functions and sets differently. A class can be in one of three states. We introduce a dynamic function *res*, which records the current resolution state. A class is *resolved*, if resolution for this class is in progress or done.

$$\begin{array}{l}
 res : Class \rightarrow ResolvedState \\
 ResolvedState ::= Unresolved \mid Resolved \mid InProgress \\
 resolved(state) = state \in \{InProgress, Resolved\}
 \end{array}$$

The JVM specification [6] uses symbolic references, namely field and method specifications, to support binary compatibility, cf. [4]. As a consequence, the calculation of field offsets and of method offsets is implementation dependent. Therefore, we keep the class field access abstract and define the storage function for class fields to be the same in Java<sub>C</sub> and JVM<sub>C</sub>, namely

$$glo : FieldSpec \rightarrow Val.$$

The runs of JVM<sub>C</sub> start with calling the class method *main* of a distinguished class *Main* being part of the environment. However, before *main* is executed, its class *Main* has to be initialized. Therefore, the frame stack initially has two entries: the *main* method at the bottom and the `<clinit>` method on the top. All classes are initially unresolved and all fields are set to their initial values. This initialization also refines the corresponding conditions imposed on Java<sub>C</sub>:

$$\begin{array}{ll}
 pcs & = start(clinit') \cdot start(main') & res = \{(c, Unresolved) \mid c \in dom(env)\} \\
 locs & = \epsilon \cdot \epsilon & glo = \{(fs, fInitVal(fs)) \mid c \in dom(env), \\
 opds & = \epsilon \cdot \epsilon & fs \in cfields(c)\} \\
 mspecs & = clinit' \cdot main'
 \end{array}$$



The method specifications  $clinit'$  and  $main'$  denote the class methods  $\langle clinit \rangle$  and  $main$  of class  $Main$ . The macro  $start$  returns the first instruction of the code of the given method specification.

$$\begin{aligned} clinit' &== \text{proc}(\text{Main}, \langle clinit \rangle) & start(ms) &== \text{next}_{unlab}(0, mCode(ms)) \\ main' &== \text{proc}(\text{Main}, main) & \text{proc}(c, m) &== (c, m, (\epsilon, \text{void})) \end{aligned}$$

The following rules for  $JVM_C$  define the semantics of the new JVM instructions, provided the class of the field or method specification is already resolved. A **getstatic** instruction loads the value (one or two words), stored under the field specification in the global environment, on top of the operand stack. A **putstatic** instruction stores the top (two) word(s) of the operand stack in the global environment at the given field specification. An **invokestatic** instruction pops the arguments from the stack and sets  $pc$  to the next instruction. The arguments of the invoked method are placed in the local variables of the new stack frame, and execution continues at the first instruction of the new method. A **return** instruction is ‘inverse’ to **invokestatic**. It pops a value from the top of the stack and pushes it onto the operand stack of the invoker. All other items in the current stack are discarded. (If the return type is **void**,  $split$  returns the empty sequence as its value.)

$$\begin{aligned} &\text{if } pc \text{ is } \text{getstatic}((c, f), t) \wedge \\ &\quad \text{resolved}(res(c)) \\ &\text{then} \\ &\quad opd := glo(c, f) \cdot opd \\ &\quad \text{proceed} \\ &\text{if } pc \text{ is } \text{putstatic}((c, f), t) \wedge \\ &\quad \text{resolved}(res(c)) \wedge \\ &\quad (v, opd') = \text{split}(t, opd) \\ &\text{then} \\ &\quad opd := opd' \\ &\quad glo(c, f) := v \\ &\quad \text{proceed} \\ &\text{if } pc \text{ is } \text{invokestatic}(c, m, (ts, t)) \wedge \\ &\quad \text{resolved}(res(c)) \wedge (t_1, \dots, t_n) = ts \wedge \\ &\quad (v_n, \dots, v_1, opd') = \text{split}(t_n, \dots, t_1, opd) \\ &\text{then} \\ &\quad \text{call}(\text{next}(pc, code), v_1 \dots v_n, \\ &\quad \quad opd', (c, m, (ts, t))) \\ &\text{if } pc \text{ is } \text{return}(t) \wedge \\ &\quad (v, opd') = \text{split}(t, opd) \\ &\text{then} \\ &\quad \text{return}(v) \end{aligned}$$

The macros  $call$  and  $return$  update the frames as follows:

$$\begin{aligned} \text{call}(pc, loc, opd, mspec) &== & \text{return}(v) &== \\ \text{let } pc_0 \cdot pcs' = pcs & & \text{if } \text{len}(pcs) = 1 \text{ then} & \\ \quad opd_0 \cdot opds' = opds \text{ in} & & \quad pcs(0) := \text{undef} & \\ pcs &:= \text{start}(mspec) \cdot pc \cdot pcs' & \text{else let } opd_0 \cdot opd_1 \cdot opds' = opds \text{ in} & \\ locs &:= loc \cdot locs & \quad pcs &:= \text{pop}(pcs) \\ opd &:= \epsilon \cdot opd \cdot opds' & \quad locs &:= \text{pop}(locs) \\ mspecs &:= mspec \cdot mspecs & \quad opd &:= (v \cdot opd_1) \cdot opds' \\ & & \quad mspecs &:= \text{pop}(mspecs) \end{aligned}$$

Execution starts in a state in which no class is resolved. A class is resolved, when it is first referenced. Before a class is resolved, its superclass is resolved (if any). Interfaces are not resolved at this time, although this is not specified in Java’s language reference manual [4]. On the level of the JVM resolution leads to three rules. First, resolutions starts, i.e. the class method  $\langle clinit \rangle$  is implicitly called, when the class referred to in a **get**-, **put**- or **invokestatic**

instruction is not resolved. Second, the class initializer records the fact that class initialization is in progress and calls the superclass initializer recursively. Third, after having executed the class initializer, it is recorded that the class is resolved.

<pre> <b>if</b> (<math>pc</math> <b>is</b> <b>putstatic</b> <math>((c, -), -) \vee</math>       <math>pc</math> <b>is</b> <b>getstatic</b> <math>((c, -), -) \vee</math>       <math>pc</math> <b>is</b> <b>invokestatic</b> <math>(c, -, -) \wedge</math>       <math>\neg resolved(res(c))</math>     <b>then</b>       <math>call(pc, \emptyset, opd, proc(c, &lt;clinit&gt;))</math> </pre>	<pre> <b>if</b> <math>res(cclass) = Unresolved</math>     <b>then</b>       <math>res(cclass) := InProgress</math>       <b>if</b> <math>supers(cclass) \neq \epsilon \wedge</math>         <math>\neg resolved(res(super(cclass)))</math>         <b>then</b>           <math>call(pc, \emptyset, opd,</math>             <math>proc(super(cclass), &lt;clinit&gt;))</math>       <b>if</b> <math>pc</math> <b>is</b> <b>return</b> <math>(t) \wedge cmeth = &lt;clinit&gt;</math>         <b>then</b>           <math>res(cclass) := Resolved</math> </pre>
---	---

Firing the second rule depends on the condition that the current class is *Unresolved*—this is the reason why we called the initializer in the first rule. To suppress the simultaneous firing of other rules we strengthen the macro ‘is’:

$$pc \text{ is } instr == code(pc) = instr \wedge resolved(res(cclass))$$

This guarantees that an instruction can only be executed, if the current class is resolved. Opposite to the second rule, the third rule fires simultaneously with the previously presented rule for the **return** instruction.

### 3.2 Compilation of Java<sub>C</sub> Programs to JVM<sub>C</sub> Code

The compilation of Java<sub>I</sub> expressions is extended by defining the compilation of class field access, class field assignment, and by the compilation of calls of class methods.

$$\begin{aligned}
 \mathcal{E}(fspec) &= \text{getstatic}(fspec, \mathcal{T}(fspec)) \\
 \mathcal{E}(fspec = e) &= \mathcal{E}e \cdot \text{dup}(\mathcal{T}(e)) \cdot \text{putstatic}(fspec, \mathcal{T}(fspec)) \\
 \mathcal{E}(mspec(e_1, \dots, e_n)) &= \mathcal{E}e_1 \dots \mathcal{E}e_n \cdot \text{invokestatic}(mspec)
 \end{aligned}$$

We add the clause for **return** statements to the Java<sub>I</sub> compilation.

$$\begin{aligned}
 \mathcal{S}(\text{return } e;) &= \mathcal{E}e \cdot \text{return}(\mathcal{T}(e)) \\
 \mathcal{S}(\text{return};) &= \text{return}(\text{void})
 \end{aligned}$$

To compile a class initializer (the *Init* phrase) means to compile its statement as the body of the static `<clinit>` method.

The extension of Java<sub>I</sub>/JVM<sub>I</sub> to Java<sub>C</sub>/JVM<sub>C</sub> is conservative, i.e. purely incremental. For the proof of the *Correctness Theorem for Java<sub>C</sub>/JVM<sub>C</sub>* it therefore suffices to extend the theorem from Java<sub>I</sub>/JVM<sub>I</sub> to the new expressions and statements occurring in Java<sub>C</sub>/JVM<sub>C</sub>.

## 4 JVM<sub>0</sub> and the Compilation of Java<sub>0</sub> Programs

In this section we extend the machine JVM<sub>C</sub> to JVM<sub>0</sub>. This machine handles the object-oriented features of Java programs, namely instances, instance creation, instance field access, instance method calls with late binding, type casts and null pointers. We add the corresponding new phrases to the definition of the compilation function.

We recall the grammar for the new expressions of Java<sub>0</sub> and define the corresponding JVM<sub>0</sub> instructions:

$  \begin{array}{l}  \text{Exp} ::= \dots \\  \quad   \text{ this} \\  \quad   \text{ new } \text{ConstrSpec} (\text{Exp}^*) \\  \quad   \text{ ConstrSpec} (\text{Exp}^*) \\  \quad   \text{ Exp}.\text{FieldSpec} \\  \quad   \text{ Exp}.\text{FieldSpec} = \text{Exp} \\  \quad   \text{ Exp}.\text{MethSpec}\{\text{CallKind}\}(\text{Exp}^*) \\  \quad   \text{ Exp } \text{instanceof } \text{Class} \\  \quad   (\text{Class}) \text{Exp} \\  \text{ConstrSpec} ::= (\text{Class} \times \text{Typ}^*)  \end{array}  $	$  \begin{array}{l}  \text{Instr} ::= \dots \\  \quad   \text{ new } (\text{Class}) \\  \quad   \text{ getfield } (\text{FieldSpec} \times \text{Typ}) \\  \quad   \text{ putfield } (\text{FieldSpec} \times \text{Typ}) \\  \quad   \text{ dup\_} (\text{Typ}^*) \\  \quad   \text{ invokeinstance } (\text{MethSpec} \times \\  \quad \quad \quad \text{CallKind}) \\  \quad   \text{ instanceof } (\text{Class}) \\  \quad   \text{ checkcast } (\text{Class}) \\  \text{CallKind} ::= \text{Constr} \mid \text{Nonvirtual} \\  \quad \quad \quad \mid \text{Virtual} \mid \text{Super}  \end{array}  $
--	---

Java<sub>0</sub> uses constructor specifications to uniquely denote overloaded instance constructors. JVM<sub>0</sub> provides instructions to allocate a new instance, to access or assign its fields, to duplicate values, to invoke instance methods and to check instance types. Java<sub>0</sub> and JVM<sub>0</sub> use the universe *CallKind*, to distinguish the particular way in which instance methods are called.

### 4.1 The Machine JVM<sub>0</sub> for Object-Oriented Code

JVM<sub>0</sub> uses the same abstract class file as JVM<sub>C</sub>. However, instance fields and instance methods—in opposite to class fields and class methods—are not static but dynamic. So we extend the universe *MemberKind* as follows:

$$\text{MemberKind} ::= \dots \mid \text{Dynamic}$$

The JVM specification [6] fixes the class file. However, the specification does not explain how instances are stored or instance methods are accessed. So we extend the signature of JVM<sub>C</sub> in JVM<sub>0</sub> in the same way as the signature of Java<sub>C</sub> is extended in Java<sub>0</sub>. We introduce the following static functions (homonymy with Java<sub>0</sub> functions) that look up information in the global environment:

$$\begin{array}{ll}
 dfields & : \text{Class} \rightarrow \mathcal{P} \text{FieldSpec} \\
 dlookup & : \text{Class} \times \text{MethSpec} \rightarrow \text{Class} \\
 compatible & : \text{Class} \times \text{Class} \rightarrow \text{Bool}
 \end{array}$$

The function *dfields* determines the instance fields of a class and of all its superclasses (if any). The function *dlookup* returns the first (super) class for the given method specification, which implements this method. The expression

$compatible(myType, tarType)$  returns *true* if  $myType$  is assignment compatible with  $tarType$  [4]. Note that at the JVM level, there is no special lookup function for constructors. Instead, Java's constructors appear in the JVM as instance initialization methods with the special name  $\langle init \rangle$ .

JVM<sub>0</sub> and Java<sub>0</sub> have the same dynamic functions for memorizing the class and the instance field values of a reference. In both machines they are initially empty. References can be obtained from the abstract universe  $Ref$ , which is assumed to be a subset of  $Word$ . (Likewise, we also assume that  $null$  is an element of  $Word$ .)

$$\begin{array}{ll} classOf : Ref \rightarrow Class & classOf = \emptyset \\ dyn : Ref \times FieldSpec \rightarrow Val & dyn = \emptyset \end{array}$$

The following rules define the semantics of the new instructions of JVM<sub>0</sub>, provided that the involved class is resolved.

<p><b>if</b> <math>pc</math> is <b>new</b> <math>(c) \wedge</math>  <math>resolved(res(c))</math>  <b>then</b>              <b>extend</b> <math>Ref</math> <b>by</b> <math>r</math>                  <math>classOf(r) := c</math>              <b>vary</b> <math>fs</math> <b>over</b> <math>dfields(c)</math>                  <math>dyn(r, fs) := fInitVal(fs)</math>              <math>opd := r \cdot opd</math>              <b>proceed</b></p> <p><b>if</b> <math>pc</math> is <b>getfield</b> <math>((c, f), t) \wedge</math>  <math>resolved(res(c)) \wedge</math>  <math>r \cdot opd' = opd \wedge</math>  <math>r \neq null</math>  <b>then</b>              <math>opd := dyn(r, (c, f)) \cdot opd'</math>              <b>proceed</b></p> <p><b>if</b> <math>pc</math> is <b>putfield</b> <math>((c, f), t) \wedge</math>  <math>resolved(res(c)) \wedge</math>  <math>(v, r, opd') = split(t, c, opd) \wedge</math>  <math>r \neq null</math>  <b>then</b>              <math>opd := opd'</math>              <math>dyn(r, (c, f)) := v</math>              <b>proceed</b></p> <p><b>if</b> <math>pc</math> is <b>dup</b> <math>-(t_1, t_2) \wedge</math>  <math>(v_2, v_1, opd') = split(t_2, t_1, opd)</math>  <b>then</b>              <math>opd := v_2 \cdot v_1 \cdot v_2 \cdot opd'</math>              <b>proceed</b></p>	<p><b>if</b> <math>pc</math> is <b>invokeinstance</b> <math>((c, m, (ts, t)), k) \wedge</math>  <math>resolved(res(c)) \wedge (t_1, \dots, t_n) = ts \wedge</math>  <math>(v_n, \dots, v_1, r, opd') =</math>  <math>split(t_n, \dots, t_1, c), opd) \wedge</math>  <math>r \neq null</math>  <b>then</b>              <math>call(next(pc, code), r \cdot v_1 \cdot \dots \cdot v_n,</math>                  <math>opd', (c', m, (ts, t)))</math>              <b>where</b>                  <math>c' = \text{case } k \text{ of}</math>                      <math>Constr : c</math>                      <math>Nonvirtual : cclass</math>                      <math>Virtual : dlookup(classOf(r),</math>                          <math>m, (ts, t))</math>                      <math>Super : dlookup(super(cclass),</math>                          <math>m, (ts, t))</math></p> <p><b>if</b> <math>pc</math> is <b>instanceof</b> <math>(c) \wedge</math>  <math>resolved(res(c)) \wedge</math>  <math>r \cdot opd' = opd</math>  <b>then</b>              <math>opd := (r \neq null \wedge</math>                  <math>compatible(classOf(r), c) \cdot opd'</math>              <b>proceed</b></p> <p><b>if</b> <math>pc</math> is <b>checkcast</b> <math>(c) \wedge</math>  <math>resolved(res(c)) \wedge</math>  <math>r \cdot opd' = opd \wedge</math>  <math>(r = null \vee compatible(classOf(r), c))</math>  <b>then</b>              <b>proceed</b></p>
--	--

A new instruction allocates a fresh reference using the domain extension update of ASMs. The  $classOf$  the reference is set to the given class, the class instance fields are set to default values, and the new reference is pushed on the operand stack. A **getfield** instruction pops the target reference from the stack, retrieves

the value of the field identified by the given field specification from the dynamic store and pushes one or two words on the operand stack. A **putfield** instruction pops a value and the target reference from the stack and sets the dynamic store at the point of the target reference and the given field specification to the popped value. A **dup\_** instruction duplicates the top value and inserts the duplicate below the top value on the stack. An **invokeinstance** instruction pops the arguments and the target reference (which denotes the instance whose method is being called) from the stack and sets *pc* to the next instruction. The method's implementing class is being located. If the call kind is

- *Constr*, the method specification denotes a constructor; its code is located in the given class. (The given method *m* must be *<init>*.)
- *Nonvirtual*, the method specification denotes a private method; its code is located in the current class. (The given class *c* must be *cclass*.)
- *Virtual*, the implementing class is looked up dynamically, starting at the class of the target reference.
- *Super*, the method is looked up dynamically, starting at the superclass of the current class. (The given class *c* must be *super(cclass)*.)

Once a method has been located, **invoke** calls the method: The arguments for the invoked method are placed in the local variables of the new stack frame, placing the target reference *r* (denoting **this** in Java) in *loc*(0). Execution continues at the first instruction of the new method. An **instanceof** instruction pops a reference from the operand stack. If the reference is not *null* and assignment compatible with the required class, the integer 1 is pushed on the operand stack, otherwise 0 is pushed. A **checkcast** instruction checks that the top value on the stack is an instance of the given class.

If the class *c* of a field or method specification or if the explicitly given class *c* of a **new**, an **instanceof** or a **checkcast** instruction is not resolved, the JVM first resolves *c*, i.e. calls *c*'s *<clinit>* method, before the instruction is executed.

```

if (pc is new (c)  $\vee$  pc is putfield (((c, -), -))  $\vee$  pc is getfield (((c, -), -))  $\vee$ 
    pc is invokeinstance (((c, -), -))  $\vee$  pc is instanceof (c)  $\vee$  pc is checkcast (c)  $\wedge$ 
     $\neg$ resolved(res(c))
then
    call(pc,  $\emptyset$ , opd, proc(c, <clinit>))

```

## 4.2 Compilation of Java<sub>0</sub> Programs to JVM<sub>0</sub> Code

Since there are no new statements in Java<sub>0</sub>, only the compilation of Java<sub>0</sub> expressions has to be extended to the new Java<sub>0</sub> expressions. The reference **this** is implemented as the distinguished local variable number 0.

$$\begin{aligned}
\mathcal{E}(\text{this}) &= \text{load}(0, \mathcal{T}(\text{this})) \\
\mathcal{E}(\text{new}(c, ts)(e_1, \dots, e_n)) &= \text{new}(c) \cdot \text{dup}(c) \cdot \mathcal{E}e_1 \cdot \dots \cdot \mathcal{E}e_n \cdot \\
&\quad \text{invokeinstance}((c, \langle \text{init} \rangle, (ts, \text{void})), \text{Constr}) \\
\mathcal{E}((c, ts)(e_1, \dots, e_n)) &= \text{load}(0, \mathcal{T}(\text{this})) \cdot \mathcal{E}e_1 \cdot \dots \cdot \mathcal{E}e_n \cdot \\
&\quad \text{invokeinstance}((c, \langle \text{init} \rangle, (ts, \text{void})), \text{Constr}) \\
\mathcal{E}(e.fspec) &= \mathcal{E}e \cdot \text{getfield}(fspec, \mathcal{T}(fspec)) \\
\mathcal{E}(e_1.fspec = e_2) &= \mathcal{E}e_1 \cdot \mathcal{E}e_2 \cdot \\
&\quad \text{dup}(\mathcal{T}(e_1), \mathcal{T}(e_2)) \cdot \text{putfield}(fspec, \mathcal{T}(fspec)) \\
\mathcal{E}(e.mspec\{k\}(e_1, \dots, e_n)) &= \mathcal{E}e \cdot \mathcal{E}e_1 \cdot \dots \cdot \mathcal{E}e_n \cdot \text{invokeinstance}(mspec, k) \\
\mathcal{E}(e \text{ instanceof } c) &= \mathcal{E}e \cdot \text{instanceof}(c) \\
\mathcal{E}((c) e) &= \mathcal{E}e \cdot \text{checkcast}(c)
\end{aligned}$$

Due to the conservativity of the extension of  $\text{Java}_C/\text{JVM}_C$  to  $\text{Java}_O/\text{JVM}_O$ , for the proof of the *Correctness Theorem for Java<sub>O</sub>/JVM<sub>O</sub>* it suffices to extend the theorem from  $\text{Java}_C/\text{JVM}_C$  to the new expressions occurring in  $\text{Java}_O/\text{JVM}_O$ .

The definitions of class initialization for  $\text{Java}_O$  in [4] and resolution for  $\text{JVM}_O$  in [6] do not match because `instanceof` and class cast expressions in Java do not call the initialization of classes. In opposite, the JVM effect is to execute the initialization of the related class if it is not initialized yet. Under the assumption that also in Java these instructions trigger class initialization, these instructions preserve the theorem for  $\text{Java}_O/\text{JVM}_O$ .

## 5 JVM<sub>E</sub> and the Compilation of Exception Treatment

In this section we extend  $\text{JVM}_O$  to  $\text{JVM}_E$  that handles exceptions. We add the compilation of the new  $\text{Java}_E$  statements and refine the compilation of jump and `return` statements.

The following grammars list the new statements of  $\text{Java}_E$  and the new  $\text{JVM}_E$  instructions.  $\text{JVM}_E$  provides instructions to raise an exception, to jump to and to return from subroutines embedded in methods.

$$\begin{array}{ll}
\text{Stm} ::= \dots & \text{Instr} ::= \dots \\
\begin{array}{|l} \text{throw } Exp; \\ \text{try Stm catch (Typ, Var, Stm)}^* \\ \text{try Stm finally Stm} \end{array} & \begin{array}{|l} \text{athrow} \\ \text{jsr (Lab)} \\ \text{ret (Varnum)} \end{array}
\end{array}$$

### 5.1 The JVM<sub>E</sub> Machine for Executing Exceptions

The JVM supports `try/catch` or `try/finally` by exception tables that list the exceptions of a method. When an exception is raised this table is searched for the handler. Exception tables refine the notion of method body as follows:

$$\begin{aligned}
\text{MethDec} &== (mKind : \text{MemberKind} \times mBody : [\text{Code} \times \text{Exception}^*]) \\
\text{Exception} &== (\text{from, to, handle} : \text{Lab} \times \text{catchTyp} : [\text{Class}])
\end{aligned}$$

The labels *from* and *to* define the range of the protected code; *handle* starts the exception handler for the optional type *catchTyp*. If no *catchTyp* is given

(as is the case for **finally** statements), any exception is caught. We refine the function  $mCode$  from  $JVM_C$  and introduce a new function  $mExcs$ , which returns the exceptions of the given method specification.

$$\begin{aligned} mCode(c, m, f) &= fst(mBody(mTab(env(c))(m, f))) \\ mExcs(c, m, f) &= snd(mBody(mTab(env(c))(m, f))) \end{aligned}$$

If a class initializer raised an exception, which is not handled within the method, Java and therefore the JVM require that the method's class must be labeled as erroneous. So we extend the domain of *ResolvedState* in the same way as we did for Java:

$$ResolvedState ::= \dots \mid Error$$

If the thrown exception is not an **Error** or one of its subclasses, then  $Java_E$  and  $JVM_E$  throw an **ExceptionInInitializerError**. If a class should be resolved but is marked as erroneous, Java and therefore implicitly the JVM require that a **NoClassDefFoundError** is reported.

We formalize the run-time system search for a handler of an exception by a recursively defined function *catch*. This function first searches the active method using *catch'*. If no handler is found (the exception handler list is empty), the current method frame is discarded, the invoker frame is reinstated and *catch* is called recursively. A handler is found if the *pc* is protected by some brackets *from* and *to*, and the thrown exception is compatible with the *catchType*. In this case the operand stack is reduced to the exception and execution continues at the address of the exception handler. When *catch'* returns from a **<clinit>** method, the method has thrown an uncaught exception; according to the strategy presented above the method's class must be labeled as erroneous.

$$\begin{aligned} catch(r, ((pc \cdot pcs, loc \cdot locs, opd \cdot opds, mspec \cdot mspecs), res)) &= \\ catch'(mExcs(mspec)) \textbf{ where} & \\ catch'(\epsilon) = & \\ \textbf{if } pcs = \epsilon \textbf{ then} & \\ ((undef \cdot pcs, loc \cdot locs, opd \cdot opds, mspec \cdot mspecs), res) & \\ \textbf{else let } (c, m, \_) = mspec & \\ res' = \textbf{if } m = \textbf{<clinit>} \textbf{ then } res \oplus \{(c, Error)\} \textbf{ else } res \textbf{ in} & \\ catch(r, ((pcs, locs, opds, mspecs), res')) & \\ catch'((from, to, handle, catchTyp) \cdot excs) = & \\ \textbf{if } jump(from, mCode(mspec)) \leq pc < jump(to, mCode(mspec)) \wedge & \\ (catchTyp = [] \vee compatible(classOf(r), catchTyp)) \textbf{ then} & \\ ((jump(handle, mCode(mspec)) \cdot pcs, loc \cdot locs, r \cdot opds, mspec \cdot mspecs), res) & \\ \textbf{else } catch'(excs) & \end{aligned}$$

The following rules define the semantics of  $JVM_E$  instructions. The **athrow** instruction pops a reference from the stack and throws the exception represented by that reference. The **jsr** instruction is used to implement Java's **finally** clause. This instruction pushes the address of the next instruction on the operand stack and jumps to the given label. This requires that the universe *Pc* (called

**ReturnAddress** in the JVM specification) is embedded in *Word*. The address, which is put on top of the stack, is used by **ret** to return from the subroutine, wherefore the return address first has to be stored in a local variable.

<pre> <b>if</b> <math>pc</math> <b>is</b> <b>athrow</b> <math>\wedge</math>   <math>r \cdot opd' = opd \wedge</math>   <math>r \neq null</math> <b>then</b>   <math>(frames, res) := catch(r, (frames, res))</math> </pre>	<pre> <b>if</b> <math>pc</math> <b>is</b> <b>jsr</b> (<math>lab</math>) <b>then</b>   <math>opd := next(pc, code) \cdot opd</math>   <math>pc := goto(lab)</math> </pre>
<pre> <b>if</b> <math>pc</math> <b>is</b> <b>ret</b> (<math>x</math>) <b>then</b>   <math>pc := loc(x)</math> </pre>	<pre> <b>if</b> <math>res(cclass) = Error</math> <b>then</b>   <math>fail(NoClassDefFoundError)</math> </pre>

If the current class is erroneous, the last rule throws a **NoClassDefFoundError** using the macro  $fail(c)$ . This macro replaces the following instruction sequence:

```
new( $c$ ), dup, invokeinstance(( $c$ , <init>, ( $\epsilon$ , void)),  $Constr$ ), athrow
```

Whether or not the constructor is called is semantically irrelevant, as long as the constructors only call superclass constructors.

We refine in the obvious way rules that raise run-time exceptions. A typical representative of this rule kind is the refinement of **bapply**. It throws an **ArithmeticException**, if the operator is an integer or long division or remainder operator and the right operand is 0.

```

if  $pc$  is bapply ( $\otimes$ )  $\wedge (0, v_1, opd') = split(\mathcal{A}(\otimes), opd) \wedge (\otimes \in DivMods)$ 
then
   $fail(ArithmeticException)$ 

```

$JVM_{\mathcal{E}}$  throws a **NullPointerException** if the target reference of a **getfield**, **putfield** or **invokeinstance** instruction is *null*, or if the reference of the **athrow** instruction is *null*. The machine throws a **ClassCastException**, if the reference on top of stack is neither *null* nor assignment compatible with the required type.

## 5.2 Compilation of Java $\mathcal{E}$ Statements to JVM $\mathcal{E}$ Instructions

Since there are no new expression in Java $\mathcal{E}$ , only the compilation of Java $\mathcal{O}$  statements has to be extended to the compilation of the new Java $\mathcal{E}$  statements.

For **try/catch** statements, the compiled **try** clause is followed by a jump to the end of the compiled statement. Next the handlers are generated. Each handler stores the exception into the 'catch' parameter, followed by the code of the catch clause and a jump to the end of the compiled statement. For **try/finally** statements  $s$ , the **try** clause is compiled followed by a call to the embedded subroutine, which is generated for the **finally** clause. The subroutine first stores the return address into a fresh variable  $ret(s)$ , and finally calls **ret** ( $ret(s)$ ). The handler for exceptions that are thrown in the **try** clause starts at  $lab_3(s)$ . The handler saves an exception of class **Throwable**, which is left on the operand stack, into the fresh local variable  $exc(s)$ , calls the subroutine, and rethrows the



exception. Variable providing functions *exc*, *ret* and also *val* that is used below, return for occurrences of statements *fresh* variable numbers. This means that any returned variable number must be unused when the exception, return address or return value is stored, and this variable definition must reach its corresponding use.

$$\begin{aligned}
 S(\text{throw } e; ) &= \mathcal{E}e \cdot \text{athrow} \\
 S(s \text{ as try } s_0 \text{ catch } (c_1, x_1, s_1) \dots (c_m, x_m, s_m)) &= \\
 &\quad \text{label}(lab_1(s)) \cdot S_{s_0} \cdot \text{goto}(lab_3(s)) \cdot \text{label}(lab_2(s)) \cdot \\
 &\quad \text{label}(lab_{3+1}) \cdot \text{store}(\overline{x_1}, c_1) \cdot S_{s_1} \cdot \text{goto}(lab_3(s)) \cdot \dots \cdot \\
 &\quad \text{label}(lab_{3+m}) \cdot \text{store}(\overline{x_m}, c_m) \cdot S_{s_m} \cdot \text{goto}(lab_3(s)) \cdot \\
 &\quad \text{label}(lab_3(s)) \\
 S(s \text{ as try } s_1 \text{ finally } s_2) &= \\
 &\quad \text{label}(lab_1(s)) \cdot S_{s_1} \cdot \text{jsr}(lab_2(s)) \cdot \text{goto}(lab_4(s)) \cdot \\
 &\quad \text{label}(lab_2(s)) \cdot \text{store}(ret(s), \text{ReturnAddress}) \cdot S_{s_2} \cdot \text{ret}(ret(s)) \cdot \\
 &\quad \text{label}(lab_3(s)) \cdot \text{store}(exc(s), \text{Throwable}) \cdot \text{jsr}(lab_2(s)) \cdot \\
 &\quad \quad \text{load}(exc(s), \text{Throwable}) \cdot \text{athrow} \cdot \\
 &\quad \text{label}(lab_4(s))
 \end{aligned}$$

If a jump statement is nested inside a try clause of a try/finally statement and its corresponding target statement contains try/finally statements, then all finally clauses between the jump statement and the target have to be executed in innermost order. The compilation uses the function *takeFinallyUntilTarget* :  $Stm \times Lab \rightarrow Stm^*$ , which given an occurrence of a statement and a label, returns in innermost order all occurrences of try/finally statements up to the target statement. For *return e* the compiler stores the result of the compiled expression *e* in a *fresh* temporary variable *val*. The compiler then generates code to jump to all outer finally statements in this method using the static function *takeFinally* :  $Stm \rightarrow Stm^*$ . Thereafter, the local variable *val* is pushed back onto the operand stack and the intended return instruction is executed.

$$\begin{aligned}
 S(s \text{ as break } lab; ) &= \text{let } (s_1, \dots, s_m) = \text{takeFinallyUntilTarget}(s, lab) \text{ in} \\
 &\quad \text{jsr}(lab_2(s_1)) \cdot \dots \cdot \text{jsr}(lab_2(s_m)) \cdot \text{goto}(lab_2(\text{target}(s, lab))) \\
 S(s \text{ as continue } lab; ) &= \text{let } (s_1, \dots, s_m) = \text{takeFinallyUntilTarget}(s, lab) \text{ in} \\
 &\quad \text{jsr}(lab_2(s_1)) \cdot \dots \cdot \text{jsr}(lab_2(s_m)) \cdot \text{goto}(lab_1(\text{target}(lab, s))) \\
 S(s \text{ as return } e; ) &= \text{let } (s_1, \dots, s_m) = \text{takeFinally}(s) \text{ in} \\
 &\quad \mathcal{E}e \cdot \text{store}(val(s), T(e)) \cdot \\
 &\quad \text{jsr}(lab_2(s_1)) \cdot \dots \cdot \text{jsr}(lab_2(s_m)) \cdot \text{load}(val(s), T(e)) \cdot \text{return}(T(e)) \\
 S(s \text{ as return; }) &= \text{let } (s_1, \dots, s_m) = \text{takeFinally}(s) \text{ in} \\
 &\quad \text{jsr}(lab_2(s_1)) \cdot \dots \cdot \text{jsr}(lab_2(s_m)) \cdot \text{return}(\text{void})
 \end{aligned}$$

In the generation of an exception table inner try phrases are concatenated before the outer ones. This guarantees that exceptions are searched in innermost order.

$$\begin{aligned}
\mathcal{X}(s \text{ as try } s_0 \text{ catch } (c_1, x_1, s_1), \dots (c_m, x_m, s_m)) &= \\
&\quad \mathcal{X}_{s_0} \cdot (lab_1(s), lab_2(s), lab_{3+1}, c_1) \cdot \mathcal{X}_{s_1} \cdot \dots \cdot \\
&\quad (lab_1(s), lab_2(s), lab_{3+m}, c_m) \cdot \mathcal{X}_{s_m} \\
\mathcal{X}(s \text{ as try } s_1 \text{ finally } s_2) &= \mathcal{X}_{s_1} \cdot (lab_1(s), lab_2(s), lab_3(s), []) \cdot \mathcal{X}_{s_2} \\
\mathcal{X}(\{s_1 \dots s_n\}) &= \mathcal{X}_{s_1} \cdot \dots \cdot \mathcal{X}_{s_n} \\
\mathcal{X}(\text{if } (e) s_1 \text{ else } s_2) &= \mathcal{X}_{s_1} \cdot \mathcal{X}_{s_2} \\
\mathcal{X}(\text{while } (e) s) &= \mathcal{X}_s \\
\mathcal{X}(lab : s) &= \mathcal{X}_s \\
\mathcal{X}(-) &= \epsilon
\end{aligned}$$

If during execution of a class initializer an exception is thrown and this is not an `Error` or one of its subclasses, then  $\text{Java}_{\mathcal{E}}$  and  $\text{JVM}_{\mathcal{E}}$  throw an `ExceptionInInitializerError`. We refine the compilation of the phrase *Init* as follows:

$$\begin{aligned}
S(\text{static } s) &= \\
S(\text{try } s \text{ catch } (\text{Exception}, x, & \\
&\quad \text{throw new } (\text{ExceptionInInitializerError}, (\epsilon, \text{void})) ());)
\end{aligned}$$

Due to the conservativity of the extension of  $\text{Java}_{\mathcal{O}}/\text{JVM}_{\mathcal{O}}$  to  $\text{Java}_{\mathcal{E}}/\text{JVM}_{\mathcal{E}}$ , for the proof of the *Correctness Theorem for  $\text{Java}_{\mathcal{E}}/\text{JVM}_{\mathcal{E}}$*  it suffices to extend the theorem from  $\text{Java}_{\mathcal{O}}/\text{JVM}_{\mathcal{O}}$  to expression and statement execution in finally and error handling code, and to prove the following

*Exception Lemma.* The execution of code in  $\text{Java}_{\mathcal{E}}$  and the execution of the corresponding compiled code in  $\text{JVM}_{\mathcal{E}}$  produce exceptions at corresponding values of the program counters in  $\text{Java}_{\mathcal{E}}$  and  $\text{JVM}_{\mathcal{E}}$ , for the same reasons, with the same failure classes (if any) and trigger the same exception handling.

## 6 Conclusion

We have presented implementation independent, rigorous yet easy to understand abstract code for the JVM as target machine for compilation of Java programs. Our definition captures faithfully the corresponding explanations of the Java Virtual Machine specification [6] and provides a practical basis for the mathematical analysis and comparison of different implementations of the machine. In particular it allowed us to prove the correctness of a general scheme for compiling Java programs into JVM code. Additionally, we have validated our work by a successful implementation in the functional programming language Haskell. The extended version of this paper [1] includes the proof details, the instruction refinement, an extensive bibliography and the discussion of related work. In an accompanying study [2] we refine the present JVM model to a defensive JVM, where we also isolate the bytecode verifier and the resolution component (including dynamic loading) of the JVM. This JVM can be used to execute compiled Java code as well as any bytecode that is loaded from the net.

*Acknowledgment.* We thank Ton Vullings for comments on this work. The first author thanks the IRIN (Institut de Recherche en Informatique de Nantes, Université de Nantes & École Centrale), in particular the *Équipe Génie logiciel, Méthodes et Spécifications formelles* for the good working environment offered during the last stage of the work on this paper.

## References

- [1] E. Börger and W. Schulte. Defining the Java Virtual Machine as platform for provably correct Java compilation. Technical report, Universität Ulm, Fakultät für Informatik. Ulm, Germany, 1998.
- [2] E. Börger and W. Schulte. A modular design for the Java VM architecture. In E. Börger, editor, *Architecture Design and Validation Methods*. Springer LNCS, to appear, 1998.
- [3] E. Börger and W. Schulte. A programmer friendly modular definition of the semantics of Java. In J. Alves-Foss, editor, *Formal Syntax and Semantics of Java(tm)*, Springer LNCS, to appear. 1998.
- [4] J. Gosling, B. Joy, and G. Steele. *The Java(tm) Language Specification*. Addison Wesley, 1996.
- [5] Y. Gurevich. Evolving algebras 1993: Lipari guide. In E. Börger, editor, *Specification and Validation Methods*. Oxford University Press, 1995.
- [6] T. Lindholm and F. Yellin. *The Java(tm) Virtual Machine Specification*. Addison Wesley, 1996.