

# A Formal Framework for Java Separate Compilation

Davide Ancona, Giovanni Lagorio, and Elena Zucca\*

DISI - Università di Genova  
Via Dodecaneso, 35, 16146 Genova (Italy)  
email: {davide,lagorio,zucca}@disi.unige.it

**Abstract.** We define a formal notion, called *compilation schema*, suitable for specifying different possibilities for performing the overall process of Java compilation, which includes typechecking of source fragments with generation of corresponding binary code, typechecking of binary fragments, extraction of type information from fragments and definition of dependencies among them. We consider three compilation schemata of interest for Java, that is, *minimal*, *SDK* and *safe*, which correspond to a minimal set of checks, the checks performed by the SDK implementation, and all the checks needed to prevent run-time linkage errors, respectively. In order to demonstrate our approach, we define a kernel model for Java separate compilation and execution, consisting in a small Java subset, and a simple corresponding binary language for which we provide an operational semantics including run-time verification. We define a safe compilation schema for this language and formally prove type safety.

## 1 Introduction

In modern programming languages, the notion of “program” as a whole has become more and more obsolete. Now, the process of developing software typically consists in writing separate pieces of code, which we will call *fragments*, following [8], each one implementing some basic functionality and relying on other functionalities provided by other fragments. The language should provide facilities which allow the development of fragments to be as much modular as possible. In particular, a highly desirable property is *separate compilation*, which means, in its strongest formulation, that a single source fragment  $f$  can be typechecked in isolation, generating a corresponding binary fragment, in a context where only type information is available on the fragments it depends on, say  $f_1, \dots, f_n$ , but no code. This phase is called *intra-checking* in [8]. Then, an executable application can be constructed linking together a collection of binary fragments s.t., for each fragment  $f$ , all the fragments  $f_1, \dots, f_n$  it depends on are available

---

\* Partially supported by DART - Dynamic Assembly, Reconfiguration and Type-checking, Murst NAPOLI - Network Aware Programming: Oggetti, Linguaggi, Implementazioni, and APPLIED SEMantics - Esprit Working Group 26142.

and satisfy the required type assumptions. This phase is called *inter-checking* in [8]. Inter-checking should guarantee that the resulting application does not raise linkage errors (for instance, a fragment needed during execution either does not exist or does not provide some expected functionality).

Though Java is a widely known paradigmatic example of language supporting separate compilation, it does not match the above schema in many respects.

1. A fragment  $f$  which depends on  $f_1, \dots, f_n$  can be typechecked in absence of their sources, but at least the corresponding binaries must be available, since in Java class files play the dual roles of interfaces and object files.
2. If (some of)  $f_1, \dots, f_n$  are only available in source form, then Java compilers enforce their typechecking, hence generation of the corresponding code<sup>1</sup>, too.
3. Linking is *dynamic*, in the sense that starting an application corresponds to running just one fragment  $f$  and then, during the execution, other needed fragments are loaded and linked on demand.

The fact (3) that linking is dynamic has a very strong impact on changing the schema we outlined above based on two distinct phases of intra-checking (separate compilation) and inter-checking (linking).

First of all, run-time inter-checks become necessary in order to guarantee that execution does not crash. For instance, the Java Virtual Machine (JVM from now on) has a bytecode verifier which finds linkage errors and throws corresponding exceptions.

On the other hand, since run-time checks are performed, inter-checks at compile-time are in a sense redundant, since in any case the fact that the execution does not crash is guaranteed by the bytecode verifier. Nevertheless, Java compilers try to anticipate at compile-time the detection of some linkage errors, performing *some* (but not all) inter-checks at compile-time. In other words, the overall compilation process is not solely specified by the typechecking rules for the source code, but some additional aspects must be modeled in order to fully capture the behavior of a compiler.

For instance, if the fragment  $f$  to be typechecked depends on a source fragment  $f'$ , Java compilers act like a typical static linker, by enforcing the typechecking of  $f'$  too, in order to ensure that, if  $f'$  either is missing or does not satisfy the type requirements of  $f$ , then the error is detected before execution.

However, Java compilers take this approach only partially, and do not perform *all* the checks which would be actually possible (for instance, no checks are performed if  $f'$  is in binary form). As a consequence, as well-known to experienced Java programmers, standard Java compilers, as SDK and Jikes, are *not type safe* in the sense that they can generate binary fragments whose execution throws linkage errors. This seems in contradiction with the fact that type safety results have been proved for the Java language [19, 10, 18]; the explanation is that these formal type systems, and the related type safety results, are only related to

---

<sup>1</sup> In this paper, we will always use the terminology *typechecking* a fragment meaning the process that in the case of a source fragment also includes generation of binary code.

the special case when a closed set of source fragments is typechecked, while (1) and (2) are not taken into account (apart from [12], see Section 5). Even worse, there is *no* specification of separate compilation in [14], hence the outcome of compilations may strongly depend on the particular compiler implementation. Moreover, as will be shown in detail in the following, rules adopted by existing compilers are quite complex and cannot be easily explained informally.

For all these reasons, we believe necessary the definition of a formal framework for Java separate compilation, providing a rigorous basis for:

- defining and investigating different possibilities for the overall compilation process (for instance: a minimal set of checks, the checks performed by standard Java compilers, as many checks as possible)
- proving desirable properties, like type safety, for a compilation process.

This is what we achieve in this paper, by modeling the overall compilation process via the formal notion of *compilation schema*. A compilation schema defines two different judgments, a *source type judgment*  $\Gamma \vdash \mathbf{S} \rightsquigarrow \mathbf{B}$  and a *binary type judgment*  $\Gamma \vdash \mathbf{B} \diamond$  modeling typechecking, in a given type environment  $\Gamma$ , of source code  $\mathbf{S}$  and binary code  $\mathbf{B}$ , respectively; in the former case the corresponding binary code  $\mathbf{B}$  is also generated. These two components model the part of Java compilation which corresponds to truly separate compilation in the sense of [8]. The fact (2) that in Java typechecking of a fragment may enforce typechecking of other fragments is modeled by another component of a compilation schema, that is, a *dependency function*. Finally, the fact (1) that in Java type information for a fragment cannot be provided separately from code is modeled by a *type extraction function* which extracts from a compilation environment  $ce$  (collection of source and binary fragments) a type environment  $\Gamma$  providing the type information needed for typechecking fragments in  $ce$ .

We consider three different compilation schemata for Java. The first, which we call *minimal*, corresponds to true separate typechecking, in the sense that no other fragment is typechecked when compilation is invoked on a fragment  $f$ . In this case, all inter-checks are left to the run-time verifier. However, note that some of the fragments  $f$  depends on *must* be available, since some of the type information needed for typechecking  $f$  has to be extracted from them (see Section 5 for some more comments).

The second compilation schema is that used by standard Java compilers (at least for what we have been able to understand by experiments, since no specification is available). In this case only some inter-checks are performed: invoking compilation on  $f$  enforces typechecking of other fragments, but not of all those which could be possibly loaded at run-time; moreover no checks are performed on binary fragments, that is, the type judgment  $\Gamma \vdash \mathbf{B} \diamond$  is always trivially valid. As a consequence, binaries obtained as result of the compilation are not guaranteed to safely link at run-time, as we will show by the examples in Section 2. Finally, we propose a compilation schema which is type safe, that is, guarantees safe linking at run-time. For this last schema we provide a full definition of the four components for a small Java subset and we prove type safety.

The paper is organized as follows. In Section 2 we illustrate by means of some examples the three different schemata. In Section 3 we introduce our framework, formally defining the notion of compilation schema and the related type safety property. In Section 4 we define, for a small subset of Java, a type safe compilation schema. Note that this implies that we not only define the syntax for the source level, but also for the binary level, for which we also provide, in order to express and prove type safety, a simple execution model, partly inspired by those defined in [11, 9], but much simpler, since our aim here is not to provide a realistic model of Java dynamic loading and linking, but to focus on type safety. Hence a contribution of the paper is also to define a kernel calculus for the Java Virtual Machine, in the spirit of other calculi that have been defined for modeling in isolation relevant aspects of the Java language [15, 3]. In Section 5 we discuss related work. Finally, Section 6 summarizes the contribution of this paper and outlines further work.

## 2 Some motivating examples

In this section we illustrate, by means of some examples, three different Java compilation schemata, called *minimal*, *SDK* and *safe*, respectively.

As already explained, the minimal compilation schema requires the minimal amount of checks over fragments: typechecking is performed only for those source fragments on which the compiler has been explicitly invoked and no checks (except those strictly necessary for compiling the sources) are performed on binary fragments. This schema fits well in open environments where source fragments to be compiled are expected to be later dynamically linked with fragments that are not available at compile-time. For instance, assume that the class `C1` we want to compile depends on a class `C2`. Even in the case the source of `C2` is available, it could be sensible avoiding typechecking of `C2` if there is a high probability that it does not correspond to the actual code that will be linked with `C1` at runtime.

The SDK schema simply corresponds to the SDK implementation of Java<sup>2</sup>. As already said, this schema falls in between the minimal and the safe schemata: it enforces more checks than the former but less than the latter. For instance, the compilation of a class `C` requires that all source fragments<sup>3</sup> directly used by `C` must be typechecked, while for all binary fragments<sup>4</sup> directly used by `C`, only their existence and format is checked (but no real typecheck is performed).

Finally, the safe schema can be sensibly applied when we expect that the fragments that will be linked at run-time are those available after the compilation; under this assumption, it makes sense to typecheck all fragments (either source or binary) used (either directly or indirectly) by a class `C`, in order to ensure that no execution of `C` will throw a linkage error (like, for instance, `NoClassDefFoundError` or `NoSuchMethodError`). To this aim, the compilation

---

<sup>2</sup> The examples in this paper are based on version 1.4 beta 2 of SDK.

<sup>3</sup> Whose corresponding binary fragment is either unavailable or older.

<sup>4</sup> Which either are more recent than the corresponding source fragment or do not have a corresponding source fragment.

schema must include all those checks on binary fragments that a safe linker would perform if Java classes were statically linked<sup>5</sup>.

While all these three schemata share the same source type judgment and type extraction function (corresponding to the Java type system defined in [14] and formalized in, e.g., [10]), they remarkably differ in the other two components (that is, dependency function and binary type judgment) as described in the following examples.

Consider the following class declarations, that we assume each one contained in a single `.java` file:

```
class Main {
    static void main(String[] args){new Used().m();}
    void g (UsedAsType x) {}
}
class Used extends UsedParent {
    int m(){return new TransUsed().m();}
}
class UsedParent{
    int m(){return 1;}
}
class TransUsed {
    int m(){return 1;}
}
class UsedAsType { ... }
```

As already stated, the same type extraction function is shared by the three schemata. The definition is straightforward: the type environment is a mapping associating to each available class `C` the type information which can be extracted from its code, that is, a pair consisting of the direct superclass of `C` and the list of method signatures<sup>6</sup> directly declared in `C`, which we will call the *direct type* of `C`.

Assume now that we want to compile `Main`. In the minimal schema, our aim is just to perform the separate typechecking (intra-checking), hence we only need the type information necessary to typecheck the source code of `Main`. In particular, for each class used in `Main`, there are two possible situations: either the class is used only as a type, like in the method `g`, and in this case we only need the existence of the class in the type environment, or it is used as the receiver's type in a method call, like in `new Used().m()`. In this case we need to know which are *all* the method signatures of the class (either directly declared or inherited). This information, which we will call the *full type* of a class, can be safely constructed having all the direct types of the classes in its ancestor

---

<sup>5</sup> Of course, such checks are usually performed at run-time by the JVM.

<sup>6</sup> In these examples we will consider for simplicity only instance method declarations, as in the Java subset defined in Section 4; in full Java the direct type of a class would also include other declared members. We assume that the method `main` is just used for starting execution.

hierarchy and provided that this hierarchy is acyclic<sup>7</sup> (this will be formalized in Figure 5 later on). In summary, `Main.java` can be successfully typechecked, producing a corresponding `Main.class`, in the type environment  $\Gamma$  defined by

$$\begin{aligned} \{\text{Main} \mapsto \langle \text{Object}, \text{void } g(\text{UsedAsType}) \rangle, \text{Used} \mapsto \langle \text{UsedParent}, \text{int } m() \rangle, \\ \text{UsedParent} \mapsto \langle \text{Object}, \text{int } m() \rangle, \text{UsedAsType} \mapsto \dots \}, \end{aligned}$$

which can be extracted from a compilation environment  $ce_1$  where only the source files `Main.java`, `Used.java`, `UsedParent.java` and `UsedAsType.java` are available. This is formalized by the validity of the judgment  $\Gamma \vdash S \rightsquigarrow B$  that will be defined in Figure 3. Note that, since, as already said, in Java type information on fragments cannot be provided separately from their code, either the `.java` or `.class` files for `Used`, `UsedParent` and `UsedAsType` must be available, even though no typechecking is performed on their code (for `UsedAsType` not even the type information is used). In the minimal schema, indeed, the set of dependencies of `Main.java` is  $\{\text{Main}\}$ , reflecting the fact that we are only interested in typechecking `Main.java`.

In both the SDK and safe schema, the set of dependencies of `Main` in  $ce_1$  includes also `Used`, `UsedParent`, `UsedAsType` and `TransUsed`. As a consequence, compilation of `Main.java` in  $ce_1$  fails for both the SDK and safe schema.

Let us consider now two compilation environments able to discriminate between the SDK and the safe schemata. First consider  $ce_2$  which contains the source files `Main.java`, `Used.class`, `UsedParent.class`, `UsedAsType.class`<sup>8</sup>, and a changed version of `TransUsed.java` which does not satisfy intra-checks (for instance, the body of method `m` returns a boolean). The type environment extracted from  $ce_2$  is still  $\Gamma$ . However, the set of dependencies of `Main` in  $ce_2$  is  $\{\text{Main}, \text{Used}, \text{UsedParent}, \text{UsedAsType}, \text{TransUsed}\}$  in the safe schema, while it is  $\{\text{Main}, \text{Used}, \text{UsedParent}, \text{UsedAsType}\}$  in the SDK schema.

Hence, in SDK the source of class `TransUsed` is not typechecked; then, a new binary fragment `Main.class` is produced. However, in the new environment obtained by enriching  $ce_2$  with the fragment `Main.class`, the execution of class `Main` throws the error `NoClassDefFoundError` (note that this error is raised instead of a type error since the class has not been compiled, hence there is no corresponding binary), whereas this error is detected at compile-time by the safe schema which performs the typechecking of the source code of `TransUsed`.

In this case the difference between the Java and the safe schema is given by the dependency function. However, even in the case dependencies are the same, the two schemata can still behave differently due to the fact that the safe schema also performs a significant binary typechecking (formalized by the validity of the judgment  $\Gamma \vdash B \diamond$  that will be defined in Figure 6). For instance, invoking the compilation of both `Main` and `TransUsed` in the compilation environment  $ce_3$  which contains `Main.java`, `Used.class`, `UsedParent.class`, `UsedAsType.class` and a changed version of `TransUsed.java` which does not satisfy type requirements in `Used` (for instance, declaring `boolean m() {return`

<sup>7</sup> However, in SDK 1.4 beta 2 the compiler loops if a used binary class has a cyclic inheritance hierarchy!

<sup>8</sup> Obtained, e.g., by compiling the whole program in the example.

`true;}`), in SDK no checks are performed on the the binary code of `Used`. Hence, again, in the binary environment obtained after the compilation, the execution of class `Main` throws the exception `NoSuchMethodError`, whereas this error is detected at compile-time by the safe schema.

### 3 Framework

We now formally define our framework for modeling the Java overall compilation process. Consistently with this aim, we will everywhere use a Java-related terminology. However, most of the notions presented here could be generalized to model the compilation process of other languages, as we will briefly discuss in Section 5.

Let us denote by  $\mathbb{C}$  the set of fragment names, that is, in Java, class/interface names<sup>9</sup>, ranged over by  $\mathcal{C}$ , and by  $\mathbb{S}$  and  $\mathbb{B}$  the set of source and binary fragments, respectively. We assume that source fragments are `.java` files containing (for simplicity) exactly one class/interface declaration, and binary fragments are `.class` files.

A *compilation environment*  $ce$  is a pair<sup>10</sup>

$$\langle ce_b, ce_s \rangle \in CE = [\mathbb{C} \rightarrow_{fin} \mathbb{B}] \times [\mathbb{C} \rightarrow_{fin} \mathbb{S}]$$

s.t.  $Def(ce_b) \cap Def(ce_s) = \emptyset$ . We will call  $ce_b$  and  $ce_s$  a *binary* and a *source environment*, respectively. Note that the assumption  $Def(ce_b) \cap Def(ce_s) = \emptyset$  means that, even in the case a class has both a binary and a source definition, the compiler considers only one of them, according to some rule.

The results of (successful) compilations are binary environments. Hence, we can model the compilation process by a (partial) function, called *compilation function*:

$$\mathcal{C} : CE \times \wp(\mathbb{C}) \rightarrow [\mathbb{C} \rightarrow_{fin} \mathbb{B}]$$

where  $\mathcal{C}(\langle ce_b, ce_s \rangle, \mathbf{CS}) = ce'_b$  intuitively means that the compilation, invoked on fragments with names in  $\mathbf{CS}$ , in the compilation environment consisting of binary fragments  $ce_b$  and source fragments  $ce_s$ , generates binary fragments  $ce'_b$ .

We introduce now the formal notion of *compilation schema*, meant to express different Java compilation processes.

A compilation schema consists of the following four components.

- A *dependency function*  $\mathcal{D}$  which gives, for any compilation environment  $ce$  and set of fragment names  $\mathbf{CS}$ , the set  $\mathbf{CS}'$  of all the fragment names on which typechecking is enforced when the compilation is invoked on  $\mathbf{CS}$ .
- A *type extraction function*  $\mathcal{T}$  which extracts from a compilation environment  $ce$  a *type environment*  $\Gamma$  providing the type information necessary to typecheck fragments in  $ce$ .

<sup>9</sup> We will consider only classes in the Java subset in Section 4.

<sup>10</sup> We denote by  $[A \rightarrow_{fin} B]$  the set of *finite partial functions* from  $A$  into  $B$ , that is, functions  $f$  from  $A$  into  $B$  which are defined on a finite subset of  $A$ , denoted  $Def(f)$ .

- A *source type judgment*  $\Gamma \vdash S \rightsquigarrow B$  expressing that in the type environment  $\Gamma$  the source fragment  $S$  is successfully typechecked generating the binary fragment  $B$ .
- A *binary type judgment*  $\Gamma \vdash B \diamond$  expressing that in the type environment  $\Gamma$  the binary fragment  $B$  is successfully typechecked.

To compile a set of fragments  $\mathbf{CS}$  in a compilation environment  $ce$ , first the needed type environment  $\Gamma$  is extracted applying  $\mathcal{T}$  to  $ce$ . Then, all the fragments in the set  $\mathbf{CS}_d$  computed from  $\mathbf{CS}$  using  $\mathcal{D}$  are typechecked generating corresponding binaries for those which were in source form. This can be formalized by the inference rule in Figure 3 which defines a compilation function  $\mathcal{C}$  in terms of the four components of a compilation schema. The second side condition,  $\mathbf{CS} \subseteq \mathit{Def}(ce_s)$ , states that a compilation can be only invoked on a set of existing sources.

$$\frac{\begin{array}{l} \forall \mathbf{C} \in \mathbf{CS}_b \ \Gamma \vdash ce_b(\mathbf{C}) \diamond \\ \forall \mathbf{C} \in \mathbf{CS}_s \ \Gamma \vdash ce_s(\mathbf{C}) \rightsquigarrow B_{\mathbf{C}} \end{array}}{\mathcal{C}(\langle ce_b, ce_s \rangle, \mathbf{CS}) = \{\mathbf{C} \mapsto B_{\mathbf{C}} \mid \mathbf{C} \in \mathbf{CS}_s\}}$$

$$\begin{array}{l} \Gamma = \mathcal{T}(\langle ce_b, ce_s \rangle) \\ \mathbf{CS} \subseteq \mathit{Def}(ce_s) \\ \mathbf{CS}_d = \mathcal{D}(\langle ce_b, ce_s \rangle, \mathbf{CS}) \\ \mathbf{CS}_b = \mathbf{CS}_d \cap \mathit{Def}(ce_b) \\ \mathbf{CS}_s = \mathbf{CS}_d \cap \mathit{Def}(ce_s) \end{array}$$

**Fig. 1.** Definition of compilation function

Let us now apply the above definitions for specifying the three different compilation schemata informally introduced in Section 2.

The type extraction function is the same for the three schemata: the type environment extracted from a compilation environment  $ce = \langle ce_b, ce_s \rangle$  is a finite partial function which associates to each  $\mathbf{C} \in \mathit{Def}(ce_b) \cup \mathit{Def}(ce_s)$  its *direct type*, that is, a pair consisting of the superclass of  $\mathbf{C}$  and the list of the method signatures declared in  $\mathbf{C}$  (see the formal definition in Figure 15 in the Appendix).

Again, the source type judgment is the same for the three schemata, and corresponds to the Java type system defined in [14] and formalized, e.g., in [10]. The formalization for the small Java subset for which we define a type safe compilation schema is given in Figure 3 in the Appendix.

On the contrary, the three schemata remarkably differ in the other two components.

For what concerns the dependency function,  $\mathcal{D}(ce, \{\mathbf{C}\})$  contains only  $\mathbf{C}$  in the minimal schema; in the safe schema  $\mathcal{D}(ce, \{\mathbf{C}\})$  contains  $\mathbf{C}$  and all the classes transitively used by  $\mathbf{C}$ , regardless that  $\mathbf{C}$  is in source or binary form (see the formal definition in Figure 14 later on). In the SDK schema the definition is much more involved. First of all,  $\mathcal{D}(ce, \{\mathbf{C}\})$  contains  $\mathbf{C}$  and all the classes directly used by  $\mathbf{C}$ . For each of these classes, say  $\mathbf{C}'$ ,  $\mathcal{D}(ce, \{\mathbf{C}\})$  also (recursively) contains  $\mathcal{D}(ce, \{\mathbf{C}'\})$  if  $\mathbf{C}'$  is in source form. If  $\mathbf{C}'$  is in binary form, then the behavior is different depending whether  $\mathbf{C}'$  is only used in  $\mathbf{C}$  as “abstract” type (for instance, field

type, parameter type, method return type) or information on the components provided by  $\mathbf{C}$  is also needed (for instance, there is a method call with receiver's type  $\mathbf{C}$ ). In the former case  $\mathcal{D}(ce, \{\mathbf{C}\})$  contains only  $\mathbf{C}'$ , in the latter it contains  $\mathbf{C}'$ , all the ancestor classes of  $\mathbf{C}'$  and (recursively)  $\mathcal{D}(ce, \{\mathbf{C}''\})$  for each ancestor  $\mathbf{C}''$  which is in source form.

This rule is quite complex, and has been extrapolated only by performing a number of compilation tests, since it is hard to deduce it from the compiler source and no other form of documentation seems to be available.

The binary type judgment also differs from schema to schema. In the minimal schema no typechecks are performed (that is, the judgment  $\Gamma \vdash \mathbf{B} \diamond$  trivially holds). In the safe schema the checks performed on a binary fragment are similar to those performed on a source fragment. A difference is, for instance, the way a method invocation is checked. In the source case the method must be found searching in all the ancestor classes, and overloading must be resolved, while in the binary case a method call is already annotated with the class where the method should be found together with its signature. The formalization for the small Java subset for which we define a type safe compilation schema is given in Figure 3 (last rule) and Figure 4 for the source case and in Figure 6 (last rule) for the binary case.

In the SDK schema no typechecks are performed on binary code. The only checks which are performed, together with the existence check on  $\mathbf{C}'$ , when type-checking a class  $\mathbf{C}$  which uses  $\mathbf{C}'$ , are on the format of the binary (analogous to the fact that Java grammar is respected in the source case) and on the correspondence between the fragment name and the name of the class defined inside. In the formal model in Section 4 we assume for simplicity that fragments are well-formed in this sense.

We introduce now the formal property of type safety for separate compilation. We assume a judgment of the form  $\mathbf{C} \rightsquigarrow_{ce_b} V$  which is valid if and only if execution of class  $\mathbf{C}$  in the binary environment  $ce_b$  terminates producing a value  $V$  which can be either a normal value or a linkage exception. Intuitively, this judgment corresponds to start the execution from class  $\mathbf{C}$  in an environment  $ce_b$  corresponding to the set of all available binaries that can be dynamically linked during the execution.

The formal definition of this judgment for the small Java subset for which we define a type safe compilation schema is given in Figure 13. Let us denote with  $ce_b[ce'_b]$  the partial function  $f$  s.t.  $Def(f) = Def(ce_b) \cup Def(ce'_b)$  and for any  $\mathbf{C} \in Def(f)$   $f(\mathbf{C}) = ce'_b(\mathbf{C})$  if  $\mathbf{C} \in Def(ce'_b)$  and  $f(\mathbf{C}) = ce_b(\mathbf{C})$  otherwise.

**Definition 1.** *A compilation function  $\mathcal{C}$  is type safe iff for any compilation environment  $\langle ce_b, ce_s \rangle$  and set of class names  $\mathcal{CS}$ , if  $\mathcal{C}(\mathcal{CS}, \langle ce_b, ce_s \rangle) = ce'_b$ , then, for all class names  $\mathbf{C} \in Def(ce'_b)$  and values  $V$ , if  $\mathbf{C} \rightsquigarrow_{ce_b[ce'_b]} V$ , then  $V$  is not a (linkage) exception.*

Note that type safety requires that execution does not throw linkage errors only when started from classes that are the product of the compilation. Indeed, an error raised by an execution started from a class  $\mathbf{C}$  present in the original binary

environment  $ce_b$  can be either an error which was already present, hence not due to the compilation, or an error due to the fact that some binary used by  $C$  has been modified. In this case we will say that the compilation function does not satisfy a different property, which we call *contextual binary compatibility* (see the Conclusion).

## 4 A safe compilation schema for Java

$ \begin{aligned} S &::= \text{class } C \text{ extends } C' \{ MDS^s \} \text{ main } E^s \\ MDS^s &::= MD_1^s \dots MD_n^s \\ MD^s &::= MH \{ \text{return } E^s; \} \\ MH &::= T_0 \text{ m}(T_1 \ x_1, \dots, T_n \ x_n) \\ E^s &::= \text{new } C \mid x \mid N \mid \\ &\quad E_0^s.m(E_1^s, \dots, E_n^s) \\ \\ T &::= C \mid \text{int} \\ \\ B &::= \langle C, C', MDS^b, E^b \rangle \\ MDS^b &::= MD_1^b \dots MD_n^b \\ MD^b &::= MH \{ \text{return } E^b; \} \\ V &::= \text{new } C \mid N \mid \varepsilon \\ E^b &::= V \mid x \mid E_0^b.m \prec \boxed{C}, T_1 \dots T_n, T \succ (E_1^b, \dots, E_n^b) \\ &\quad E_0^b.m \prec C, T_1 \dots T_n, T \succ (E_1^b, \dots, E_n^b) \mid \text{new } \boxed{C} \\ \varepsilon &::= \text{ClassNotFound} \mid \text{ClassCircularityError} \mid \text{VerifyError} \mid \text{NoSuchMethod} \\ \\ MS &::= T \text{ m}(T_1 \dots T_n) \\ AMS &::= C \ MS \\ AMSS &::= AMS_1 \dots AMS_n \end{aligned} $
--

Fig. 2. Syntax and types

The language we consider is reminiscent of Featherweight Java [15], in the sense that it is a small functional subset of Java (see Figure 2); however, since here we are mainly interested in code generation and bytecode execution, we present a simple binary language as well, together with its reduction semantics. The dynamic semantics of our Java subset is indirectly defined by a compilation function mapping well-typed source fragments into well-typed binary fragments.

Metavariables  $C$ ,  $m$ ,  $x$  and  $N$  range over sets of class, method and parameter names, and integer literals, respectively.

A source fragment  $S$  is a class declaration consisting of the class name, the name of the superclass, a sequence of method declarations  $MDS^s$ , and an expression  $E^s$  playing the role of the (static) `main` method, that for simplicity we assume present in all classes. A method declaration  $MD^s$  consists of a method

header and a method body (an expression). A method header  $MH$  consists of a (return) type, a method name and a sequence of parameter types and names. There are four kinds of expression: instance creation, parameter name, integer literal and method invocation. A type is either a class name or `int`. We will use the abbreviation  $\bar{T}$  for  $T_1 \dots T_n$  in the following.

Our description of bytecode is rather abstract: we basically enrich the source language with two kinds of annotation. Each method invocation is annotated with a method descriptor, which is a triple describing the method which has been statically selected for the call: the first component corresponds to the static type of the receiver<sup>11</sup>, the second to the type of the parameters and the third to the return type. Moreover, each class name mentioned either in class creation or as the first component of method descriptors in method call is (initially) boxed. The idea is that a reference to a class is “sealed” in a box until it has been verified (at runtime). Such a reference cannot be used until it has been unboxed.

A binary fragment  $B$  consists of the name of the class, the name of the superclass, a set of binary method declarations  $MDS^b$  and the binary expression corresponding to the method `main`. This expression is used as the starting point of the execution when the class is run. A binary method declaration  $MD^b$  is structurally equivalent to a source method declaration except that the body is a binary expression.

Binary expressions can be either *values*, or parameters, or (either boxed or unboxed) method calls, or a boxed creation expression. Values correspond to the normal forms of the reduction semantics of binary fragments (defined in Figure 13), and can be either unboxed creation expressions, or integer literals, or exceptions (in case of abnormal termination).

Note that exceptions and unboxed method invocation and creation are only needed for defining the rewriting rules for bytecode execution (see Figure 13), but they are not considered valid binary formats, even though for sake of simplicity we do not have introduced two separate syntactic categories corresponding to valid binary format and valid run-time expressions, respectively.

In the last part of Figure 2 we define *method signatures* and *annotated method signatures*, which are not part of the syntax but will be used in the type judgments. A method signature  $MS$  is a method header without the argument names; an annotated method signature  $AMS$  is a method signature prefixed by an annotation indicating the class which contains the method declaration.

We start now the formal definition of the four components of our safe compilation schema, which will be used (as shown in Figure 3 in Section 3 to define the corresponding compilation function).

The dependency function  $\mathcal{D}$  is defined as follows:

$$\mathcal{D}(ce, CS) = \{C' \mid \exists C \in CS \text{ s.t. } C \xrightarrow{*}_{ce} C'\}$$

---

<sup>11</sup> This is a change introduced in SDK 1.4, since in the previous versions the first component of method descriptors corresponded to the class where the method was statically found.

where  $\xrightarrow{*}_{ce}$  is the reflexive and transitive closure of the relation  $\rightarrow_{ce}$  defined by  $C \rightarrow_{ce} C'$  iff  $C' \in \text{refClasses}(C, ce)$ . This latter function, defined in Figure 14 in the Appendix, gives the set of all classes explicitly mentioned in the code of  $C$ .

A type environment  $\Gamma$  is a finite (partial) function from class names into *direct types* of classes, that is, pairs  $\langle C', MS_1 \dots MS_n \rangle$  where  $C'$  is the (direct) superclass of  $C$  and  $MS_1 \dots MS_n$  is the list of the signatures of methods declared in  $C$ .

The type extraction function  $\mathcal{T}$ , defined in Figure 15 in the Appendix, simply throws away method bodies and parameter names, retaining type information from all classes in the compilation environment.

$\frac{\Gamma \vdash C ::_{\diamond} \_ \quad \Gamma \vdash \text{MDS}^s \rightsquigarrow \text{MDS}^b \quad \Gamma; \emptyset \vdash E^s : \_ \rightsquigarrow E^b}{\Gamma \vdash \text{class } C \text{ extends } C' \{ \text{MDS}^s \} \text{ main } E^s \rightsquigarrow \langle C, C', \text{MDS}^b, E^b \rangle}$	
$\frac{\forall i \in 1..n \quad \Gamma \vdash \text{MD}_i^s \rightsquigarrow \text{MD}_i^b}{\Gamma \vdash \text{MD}_1^s \dots \text{MD}_n^s \rightsquigarrow \text{MD}_1^b \dots \text{MD}_n^b}$	
$\frac{\Gamma; \{x_1 \mapsto T_1, \dots, x_n \mapsto T_n\} \vdash E^s : T \rightsquigarrow E^b \quad \Gamma \vdash T \leq T_0}{\Gamma \vdash T_0 \text{ m}(T_1 \ x_1, \dots, T_n \ x_n) \{ \text{return } E^s; \} \rightsquigarrow T_0 \text{ m}(T_1 \ x_1, \dots, T_n \ x_n) \{ \text{return } E^b; \}}$	$\forall i \in 1..n \ T_i \in \text{Def}(\Gamma) \cup \{\text{int}\}$
$\frac{}{\Gamma; H \vdash \text{new } C : C \rightsquigarrow \text{new } \boxed{C}} \quad C \in \text{Def}(\Gamma)$	
$\frac{}{\Gamma; H \vdash N : \text{int} \rightsquigarrow N}$	
$\frac{}{\Gamma; H \vdash x : T \rightsquigarrow x} \quad H(x) = T$	
$\frac{\Gamma; H \vdash E_0^s : C \rightsquigarrow E_0^b \quad \forall i \in 1..n \ \Gamma; H \vdash E_i^s : T_i \rightsquigarrow E_i^b}{\Gamma; H \vdash E_0^s.m(E_1^s, \dots, E_n^s) : T' \rightsquigarrow E_0^b.m \prec \boxed{C}, \bar{T}', T' \succ (E_1^b, \dots, E_n^b)} \quad \langle \bar{T}', T' \rangle = \text{methRes}(\Gamma, C, m, T_1 \dots T_n)$	

**Fig. 3.** Source type-judgment

In Figure 3 are defined the rules for typechecking, in a given type environment, source fragments (with generation of the corresponding binary code).

The first rule defines the typechecking of a class declaration.

First, it is checked that  $C$  has a well-formed *full type* in  $\Gamma$ . The full type of a class is the list of all the annotated method signatures of the class, either directly declared or inherited, and it can be safely constructed if there are no cycles in the inheritance hierarchy of  $C$  and the Java rules on method overriding are respected. This is formalized by the judgment  $\Gamma \vdash C ::_{\diamond} \_$  defined in Figure 5.

Then, the method bodies and the main expression are checked and compiled.

The second rule defines the typechecking of a sequence of method declarations. Each method declaration is correct (third rule) if the type of the expression body is a subtype of the declared return type (premises) and all argument types exist in  $\Gamma$  (note that *no* check is performed on the fact that they have well-formed full types). The judgment  $\Gamma \vdash T \leq T'$  is valid whenever  $T$  is a subtype of  $T'$  in the type environment  $\Gamma$ , and is defined in Figure 7.

Other rules define the typechecking of expressions, which also needs a local type environment  $\Pi$  which is a (partial) function from parameters into types.

An instance creation expression, `new C`, is well-typed, and has type  $C$ , in  $\Gamma$  and  $\Pi$  if  $C$  exists in  $\Gamma$ . An integer literal is trivially well-typed, and has type `int`, in every  $\Gamma$  and  $\Pi$ .

A parameter is well-typed in  $\Gamma$  and  $\Pi$  if it belongs to the domain of the local type environment, and it has the corresponding type.

A method call expression is typechecked in two steps: first, the receiver expression and all the argument expressions are typechecked finding their types  $T_i$  (and generating the corresponding binary expressions  $E_i^b$ ). Then, using this information, the most specific among the applicable methods is selected, as formally defined by the function *methRes*, defined in Figure 4, which returns a pair consisting of the type of parameters and returned value used to annotate the binary method call produced as the result of the compilation; the annotation is used at runtime by the JVM (see Figure 13). Recall that in SDK 1.4 the first component of the descriptor which annotates the method call is the static type of the receiver ( $C$  in the rule).

$$\begin{aligned}
\text{methRes}(\Gamma, C, m, \bar{T}) &= \begin{cases} \langle \bar{T}', T' \rangle & \text{if } \Gamma \vdash C :: \text{AMSS} \wedge \\ & \text{mostSpec}(\Gamma, \text{appMeth}(\Gamma, \text{AMSS}, m, \bar{T})) = C' T' m(\bar{T}') \\ \perp & \text{otherwise} \end{cases} \\
\text{appMeth}(\Gamma, A, m, \bar{T}) &= A \\
\text{appMeth}(\Gamma, C T m(\bar{T}'), \text{AMSS}, m, \bar{T}) &= \\
\begin{cases} C T m'(\bar{T}') \text{ appMeth}(\Gamma, \text{AMSS}, m, \bar{T}) & \text{if } m = m' \wedge \Gamma \vdash \bar{T} \leq \bar{T}' \\ \text{appMeth}(\Gamma, \text{AMSS}, m, \bar{T}) & \text{otherwise} \end{cases} \\
\text{mostSpec}(\Gamma, A) &= \perp \\
\text{mostSpec}(\Gamma, C T m(\bar{T})) &= C T m(\bar{T}) \\
\text{mostSpec}(\Gamma, C T m(\bar{T}) \text{ AMSS}) &= \text{min}(\Gamma, C T m(\bar{T}), \text{mostSpec}(\Gamma, \text{AMSS})) \text{ with } \text{AMSS} \neq A \\
\text{min}(\Gamma, C T m(\bar{T}), C' T' m(\bar{T}')) &= \begin{cases} C T m(\bar{T}) & \text{if } \Gamma \vdash C \bar{T} \leq C' \bar{T}' \\ C' T' m(\bar{T}') & \text{if } \Gamma \vdash C' \bar{T}' \leq C \bar{T} \\ \perp & \text{otherwise} \end{cases}
\end{aligned}$$

**Fig. 4.** Definition of *methRes*

In Figure 5 we define the judgment  $\Gamma \vdash C :: \text{AMSS}$ , associating to a class its full type AMSS, and the judgment  $\Gamma \vdash C ::_{\diamond} \text{AMSS}$ , which is valid only if AMSS is well-formed.

$\overline{\Gamma \vdash \text{Object} :: A}$	
$\frac{\Gamma \vdash C' :: \text{AMSS}'}{\Gamma \vdash C :: \text{AMSS}' \ C \ \text{MS}_1 \ \dots \ C \ \text{MS}_n}$	$\Gamma(C) = \langle C', \text{MS}_1 \dots \text{MS}_n \rangle$
$\frac{\Gamma \vdash C :: \text{AMSS}}{\Gamma \vdash C ::_{\diamond} \text{AMSS}}$	$\text{noRep}(\text{AMSS}) \wedge \text{okOverride}(\text{AMSS})$
$\text{noRep}(\text{AMS}_1 \dots \text{AMS}_n) \iff \forall i, j \in 1..n \ \text{AMS}_i = \text{AMS}_j \implies i = j$	
$\text{okOverride}(\text{AMSS}) \iff$ $\forall \text{AMS}, \text{AMS}' \in \text{AMSS} \ \text{AMS} = C \ \text{T} \ m(\bar{\text{T}}) \wedge \text{AMS}' = C' \ \text{T}' \ m(\bar{\text{T}}) \implies \text{T} = \text{T}'$	

**Fig. 5.** Full type of a class

As already said, the full type of a class consists of the sequence of the annotated signatures of the methods either directly declared in  $C$  or inherited.

A full class type AMSS is well-formed if it does not contain duplicate method signatures (predicate *noRep*) and if the Java rules on overriding are satisfied (predicate *okOverride*).

In Figure 5 the notation  $\text{AMS} \in \text{AMSS}$  is a shortcut for  $\exists \text{AMSS}_0, \text{AMSS}_1 : \text{AMSS} = \text{AMSS}_0 \ \text{AMS} \ \text{AMSS}_1$ .

Note that neither  $\Gamma \vdash C :: \text{AMSS}$  nor  $\Gamma \vdash C ::_{\diamond} \text{AMSS}$  can be deduced for  $C$  if it has a cyclic inheritance hierarchy in  $\Gamma$ .

Figure 6 shows the rules for typechecking binary fragments, which are analogous to those for source fragments shown in Figure 3, except for the last rule, concerning method calls.

Indeed, as already mentioned, in a binary method call the descriptor annotation indicates exactly which method to look for, and we only have to check that the types of the receiver expression and of the parameters are subtypes of those specified in the descriptor, and the class of the actual receiver  $C$  still implements such a method (premise  $\Gamma \vdash C \triangleleft C' \ \text{T}' \ m(\text{T}'_1 \dots \text{T}'_n)$ , see Figure 7 above for the definition of this judgment); this informally means that class  $C$  must inherit method  $\text{T}' \ m(\text{T}'_1 \dots \text{T}'_n)$  from  $C'$  or any of  $C'$  superclasses (of course, if  $C' = C$ , then the method can also be defined in  $C$  itself). Note that this corresponds to the run-time check performed by the JVM when invoking methods, therefore requiring method  $\text{T}' \ m(\text{T}'_1 \dots \text{T}'_n)$  to be exactly defined in  $C'$  would be too strong. The judgment  $\Gamma \vdash \text{int} \leq \text{int}$  is trivially valid in every  $\Gamma$ .

Every class  $C$  defined in  $\Gamma$  is a subtype of itself (second rule) and of its (direct) superclass (third rule). Note that every class in  $\Gamma$  is considered subclass

$$\begin{array}{c}
\frac{\Gamma \vdash \mathbf{C} ::_{\diamond} \_ \quad \Gamma \vdash \text{MDS}^b \diamond \quad \Gamma; \emptyset \vdash \mathbf{E}^b : \_ \diamond}{\Gamma \vdash \langle \mathbf{C}, \mathbf{C}', \text{MDS}^b, \mathbf{E}^b \rangle \diamond} \\
\\
\frac{\forall i \in 1..n \quad \Gamma \vdash \text{MD}_i^b \diamond}{\Gamma \vdash \text{MD}_1^b \dots \text{MD}_n^b \diamond} \\
\\
\frac{\Gamma; \{x_1 \mapsto \mathbf{T}_1, \dots, x_n \mapsto \mathbf{T}_n\} \vdash \mathbf{E}^b : \mathbf{T} \diamond \quad \Gamma \vdash \mathbf{T} \leq \mathbf{T}_0}{\Gamma \vdash \mathbf{T}_0 \text{ m}(\mathbf{T}_1 \ x_1, \dots, \mathbf{T}_n \ x_n) \{ \text{return } \mathbf{E}^b; \} \diamond} \quad \forall i \in 1..n \ \mathbf{T}_i \in \text{Def}(\Gamma) \cup \{\text{int}\} \\
\\
\frac{}{\Gamma; \Pi \vdash \text{new } \boxed{\mathbf{C}} : \mathbf{C} \diamond} \quad \mathbf{C} \in \text{Def}(\Gamma) \\
\\
\frac{}{\Gamma; \Pi \vdash \mathbf{N} : \text{int} \diamond} \\
\\
\frac{}{\Gamma; \Pi \vdash \mathbf{x} : \mathbf{T} \diamond} \quad \Pi(\mathbf{x}) = \mathbf{T} \\
\\
\frac{\Gamma; \Pi \vdash \mathbf{E}_0^b : \mathbf{C} \diamond \quad \forall i \in 1..n \ \Gamma; \Pi \vdash \mathbf{E}_i^b : \mathbf{T}_i \diamond \quad \forall i \in 1..n \ \Gamma \vdash \mathbf{T}_i \leq \mathbf{T}'_i \quad \Gamma \vdash \mathbf{C} \triangleleft \mathbf{C}' \ \mathbf{T}' \text{ m}(\mathbf{T}'_1 \dots \mathbf{T}'_n)}{\Gamma; \Pi \vdash \mathbf{E}_0^b . \text{m} \langle \boxed{\mathbf{C}'}, \mathbf{T}'_1 \dots \mathbf{T}'_n, \mathbf{T}' \rangle (\mathbf{E}_1^b, \dots, \mathbf{E}_n^b) : \mathbf{T}' \diamond}
\end{array}$$

Fig. 6. Binary type-judgment

$$\begin{array}{c}
\frac{}{\Gamma \vdash \text{int} \leq \text{int}} \\
\\
\frac{}{\Gamma \vdash \mathbf{C} \leq \mathbf{C}} \quad \mathbf{C} \in \text{Def}(\Gamma) \\
\\
\frac{\Gamma \vdash \mathbf{C}' :: \_ \quad \Gamma(\mathbf{C}) = \langle \mathbf{C}', \_ \rangle}{\Gamma \vdash \mathbf{C} \leq \mathbf{C}'} \\
\\
\frac{\Gamma \vdash \mathbf{C} \leq \mathbf{C}' \quad \Gamma \vdash \mathbf{C}' \leq \mathbf{C}''}{\Gamma \vdash \mathbf{C} \leq \mathbf{C}''} \\
\\
\frac{\forall i \in 1..n \ \Gamma \vdash \mathbf{T}_i \leq \mathbf{T}'_i}{\Gamma \vdash \mathbf{T}_1 \dots \mathbf{T}_n \leq \mathbf{T}'_1 \dots \mathbf{T}'_n} \\
\\
\frac{\Gamma \vdash \mathbf{C} \leq \mathbf{C}_k}{\Gamma \vdash \{ \mathbf{C}_1 \ \text{MS}_1, \dots, \mathbf{C}_n \ \text{MS}_n \} \triangleleft \mathbf{C} \ \text{MS}} \quad k \in 1..n \wedge \text{MS} = \text{MS}_k \\
\\
\frac{\Gamma \vdash \mathbf{C} :: \text{AMSS} \quad \Gamma \vdash \text{AMSS} \triangleleft \text{AMS}}{\Gamma \vdash \mathbf{C} \triangleleft \text{AMS}}
\end{array}$$

Fig. 7. Implementation and widening

of itself, even if its inheritance hierarchy is cyclical, because this does not lead to wrong type assumption. Vice versa, a class  $C$  is considered subclass of  $C'$  only if its inheritance relation is acyclic.

Subtyping relation is transitive, fourth rule. The fifth rule extends the sub-type relation to tuples of types.

The judgment  $\Gamma \vdash \{C_1 MS_1, \dots, C_n MS_n\} \triangleleft C MS$  is valid whenever one of the annotated method signatures  $C_i MS_i$  “implements” the annotated method signature  $C MS$ , that is, there is a method with the same signature in  $C$  or any of its superclasses. If a class has full type  $AMSS$  and the latter implements the method signature  $AMS'$ , then the class is said to implement  $AMS'$ .

The last rule defines the judgment  $\Gamma \vdash C \triangleleft AMS'$ , that is valid whenever class  $C$  implements  $AMS'$ .

We define now execution and verification of binary fragments. We anticipate some auxiliary definitions.

$$\boxed{
 \begin{array}{l}
 S_1 \cup^E S_2 =_{\text{def}} \begin{cases} S_1 & \text{if } S_1 \in \varepsilon \\ S_2 & \text{if } S_2 \in \varepsilon, S_1 \notin \varepsilon \\ S_1 \cup S_2 & \text{otherwise} \end{cases} \\
 b_1 \wedge^E b_2 =_{\text{def}} \begin{cases} b_1 & \text{if } b_1 \in \varepsilon \\ b_2 & \text{if } b_2 \in \varepsilon, b_1 \notin \varepsilon \\ Ok & \text{otherwise} \end{cases}
 \end{array}
 }$$

**Fig. 8.** Exception-aware union and disjunction

Figure 8 shows the definition of the two operations  $\cup^E$  and  $\wedge^E$ , whose arguments are either set of class names or exceptions.

The former is just set union when both arguments are not exceptions, whereas it returns one of its arguments when it is an exception (giving priority to the left argument).

The latter is similar to a boolean conjunction; it returns *Ok* when both arguments are *Ok* and one of its argument when it is an exception (again, giving priority to the left argument).

The function  $MBody(ce_b, C, m, T_1 \dots T_n, T)$ , defined in Figure 9, models method look-up at run-time; it searches in the binary environment  $ce_b$  for the body of a method named  $m$  whose argument types are  $T_1 \dots T_n$  and return type is  $T$ , starting from class  $C$ . If such a method is not found in  $C$ , then it is searched, recursively, in its superclass.

The result of  $MBody$  can be either the body of the method (if found) or an exception if it is not found. There are two kinds of error that can happen during the method look-up: if  $C = \text{Object}$  then this means that the method cannot be found, so the exception *NoSuchMethod* is returned. Otherwise, if  $C$  cannot be found in the binary environment, then the exception *ClassNotFound* is returned.

$$\begin{array}{l}
MBody(ce_b, \mathbf{C}, \mathbf{m}, T_1 \dots T_n, T) = \\
\left\{ \begin{array}{l}
NoSuchMethod \text{ if } \mathbf{C} = \mathbf{Object} \\
ClassNotFound \text{ if } \mathbf{C} \neq \mathbf{Object}, \mathbf{C} \notin Def(ce_b) \\
\langle \mathbf{E}^b, \mathbf{x}_1 \dots \mathbf{x}_n \rangle \text{ if } T \mathbf{m}(T_1 \mathbf{x}_1, \dots, T_n \mathbf{x}_n) \{ \mathbf{E}^b \} \in code(ce_b(\mathbf{C})) \\
MBody(ce_b, superclass(ce_b(\mathbf{C})), \mathbf{m}, T_1 \dots T_n, T) \text{ otherwise}
\end{array} \right. \\
\\
code(\langle -, -, MDS^b, - \rangle) = MDS^b \\
\\
superclass(\langle -, \mathbf{C}', -, - \rangle) = \mathbf{C}'
\end{array}$$

**Fig. 9.** Definition of *MBody*

$$\begin{array}{l}
[\cdot]^{Exp} ::= [\cdot].m \prec \langle \mathbf{C}, T_1 \dots T_n, T \rangle (\mathbf{E}_1^b, \dots, \mathbf{E}_n^b) \mid \\
\quad \mathbf{new} \mathbf{C}.m \prec \langle \mathbf{C}', T_1 \dots T_n, T \rangle (v_1, \dots, v_{i-1}, [\cdot], \mathbf{E}_{i+1}^b, \dots, \mathbf{E}_n^b) \\
[\cdot]^{Type} ::= \mathbf{new} \mathbf{C}.m \prec [\cdot], T_1 \dots T_n, T \rangle (v_1, \dots, v_n) \mid \\
\quad \mathbf{new} [\cdot]
\end{array}$$

**Fig. 10.** Contexts

In Figure 10 we introduce two kinds of contexts: expression contexts  $[\cdot]^{Exp}$  and type contexts  $[\cdot]^{Type}$ . In rewrite semantics, given in Figure 13, the former are used to propagate execution to sub-expressions and the latter to verify class references in order to unbox them (making it possible to “normal” execution to proceed). The function *Weak\_subtype* is used, at verification time, to check whether a type is subtype of another;  $Weak\_subtype(T, T', ce_b)$ , defined in Figure 11, returns *Ok* when the type  $T$  is a subtype of  $T'$  in the binary environment  $ce_b$  or an appropriate exception when it is not. The subtype is “weak” because of the special case  $T = T'$ : any type  $T$  is *always* considered subtype of itself; it does not even matter whether  $T$  exists or not in the binary environment  $ce_b$ . In all the other cases both  $T$  and  $T'$  must exist, otherwise an exception *ClassNotFound* is returned. When both exist in  $ce_b$ , the auxiliary function *Supertypes* is used to check the relationship between  $T$  and  $T'$ ; indeed,  $T$  is subtype of  $T'$  iff  $T'$  is a supertype of  $T$ . The function  $Supertypes(T, ce_b)$  can either return the set of supertypes of  $T$  in  $ce_b$ , when they can be computed, or an exception in case of error. There are two possible error situations: when a class (directly or indirectly) extends itself and when a parent class is not found in  $ce_b$ ; in these cases the exceptions *ClassCircularityError* and *ClassNotFound* are, respectively, returned. Figure 12 shows the verification judgments. The top-level judgment  $\vdash_{ce_b} \mathbf{C} : Ok$  is valid whenever the class  $\mathbf{C}$  can be verified in the binary environment  $ce_b$ , otherwise it is valid a judgment  $\vdash_{ce_b} \mathbf{C} : \varepsilon$ , where  $\varepsilon$  indicates the error occurred in the verification steps. Indeed, it can be proved that the verification process is deterministic and always terminates (either with *Ok* or with an exception).

$$\begin{array}{l}
\text{Weak\_subtype}(\mathbf{T}, \mathbf{T}', ce_b) = \begin{cases} Ok & \text{if } \mathbf{T} = \mathbf{T}' \\ \varepsilon & \text{if } \mathbf{T} \neq \mathbf{T}', \text{Supertypes}(\mathbf{T}, ce_b) = \varepsilon \\ Ok & \text{if } \mathbf{T} \neq \mathbf{T}', \mathbf{T}' \in \text{Def}(ce_b), \\ & \mathbf{T}' \in \text{Supertypes}(\mathbf{T}, ce_b) \\ \text{VerifyError} & \text{if } \mathbf{T} \neq \mathbf{T}', \mathbf{T}' \in \text{Def}(ce_b), \\ & \mathbf{T}' \notin \text{Supertypes}(\mathbf{T}, ce_b) \\ \text{ClassNotFound} & \text{otherwise} \end{cases} \\
\\
\text{Supertypes}(\mathbf{T}, ce_b) = \text{Supertypes}_{aux}(\mathbf{T}, \emptyset, ce_b) \\
\text{Supertypes}_{aux}(\mathbf{T}, LC, ce_b) = \begin{cases} \text{ClassCircularityError} & \text{if } \mathbf{T} \in LC \\ \text{ClassNotFound} & \text{if } \mathbf{T} \notin \text{Def}(ce_b) \cup \{\text{int}, \text{Object}\} \\ \{\text{Object}\} \cup LC & \text{if } \mathbf{T} = \text{Object} \\ \{\text{int}\} & \text{if } \mathbf{T} = \text{int} \\ \{\mathbf{T}\} \cup^E \text{Supertypes}_{aux}(\text{superclass}(ce_b(\mathbf{T})), \{\mathbf{T}\} \cup LC, ce_b) & \text{if } \mathbf{T} \notin LC, \mathbf{T} \in \text{Def}(ce_b) \end{cases}
\end{array}$$

**Fig. 11.** *Weak\_subtype* and *Supertypes* definitions

Note the interesting relation between verification and typechecking of binaries as defined in Figure 6; the former corresponds to dynamic typechecking of binaries, whereas the latter to static typechecking, and, hence, is more conservative than the former. This relation is formalized by Theorem 4 below.

When a class exists in the current binary environment  $ce_b$  its verification consists of: checking that there are no different method declarations with the same signature in the code of the class (predicate *noDup*) and verifying that all method declarations ( $\vdash_{ce_b} \text{MD}_i^b : b_i$ ) and the main expression ( $\emptyset \vdash_{ce_b} \mathbf{E}^b : \langle -, b \rangle$ ) are *Ok*.

When a class does not exist in a binary environment  $ce_b$  its verification simply gives *ClassNotFound* (second metarule).

The judgment  $\sigma \vdash_{ce_b} \mathbf{E}^b : \langle \mathbf{T}, b \rangle$  is valid whenever the expression  $\mathbf{E}^b$  in a binary context  $ce_b$  and local type environment  $\sigma$  has type  $\mathbf{T}$  and the result of its verification is  $b$ . The value  $b$  can be either *Ok*, when the verification succeeds, or an exception, indicating the problem, when the verification fails. In this latter case the value of  $\mathbf{T}$  is immaterial. A local type environment  $\sigma$  is a (partial) function from argument names to types.

The verification of a method declaration (third rule) succeeds when the verification of its body succeeds and the type of the body is a subtype of the declared return type.

The verification of a method invocation succeeds when the number of arguments coincides with the number of parameter types in the method descriptor, the verification of the receiver and of each argument type succeeds and the type of the receiver and of each argument is a (weak) subtype of the corresponding type contained in the method descriptor (fourth rule).

The fifth rule covers the case when the numbers of arguments differs (side condition  $k \neq n$ ). Note that if a binary fragment is the result of the compilation

$\frac{\forall i \in 1..n \vdash_{ce_b} \mathbf{MD}_i^b : b_i \quad \emptyset \vdash_{ce_b} \mathbf{E}^b : \langle \neg, b \rangle}{\vdash_{ce_b} \mathbf{C} : noDup(\mathbf{MD}_1^b \dots \mathbf{MD}_n^b) \wedge_{i \in 1..n} b_i \wedge^E b} \quad ce_b(\mathbf{C}) = \langle \neg, \neg, \mathbf{MD}_1^b \dots \mathbf{MD}_n^b, \mathbf{E}^b \rangle$
$\frac{}{\vdash_{ce_b} \mathbf{C} : ClassNotFound} \quad \mathbf{C} \notin Def(ce_b)$
$\frac{\{x_1 \mapsto T_1, \dots, x_n \mapsto T_n\} \vdash_{ce_b} \mathbf{E}^b : \langle T, b \rangle}{\vdash_{ce_b} T_0 \ m(T_1 \ x_1, \dots, T_n \ x_n) \ \{\mathbf{return} \ \mathbf{E}^b;\} : b \wedge^E Weak\_subtype(T, T_0, ce_b)}$
$\frac{\forall i \in 0..n \ \sigma \vdash_{ce_b} \mathbf{E}_i^b : \langle T_i, b_i \rangle}{\sigma \vdash_{ce_b} \mathbf{E}_{0..m}^b \prec \boxed{T'_0}, \bar{T}', T \succ (\mathbf{E}^b_1, \dots, \mathbf{E}^b_n) : \langle T, \bigwedge_{i \in 0..n} s_i \rangle} \quad \begin{array}{l} \bar{T}' = T'_1 \dots T'_n \\ s_i = b_i \wedge^E \\ Weak\_subtype(T_i, T'_i, ce_b) \end{array}$
$\frac{}{\sigma \vdash_{ce_b} \mathbf{E}_{0..m}^b \prec \boxed{T'_0}, T'_1 \dots T'_n, T \succ (\mathbf{E}^b_1, \dots, \mathbf{E}^b_k) : \langle T, VerifyError \rangle} \quad k \neq n$
$\frac{}{\sigma \vdash_{ce_b} n : \langle \mathbf{int}, Ok \rangle} \quad n = 0, 1, -1, 2, -2, \dots$
$\frac{}{\sigma \vdash_{ce_b} \mathbf{new} \ \boxed{\mathbf{C}} : \langle \mathbf{C}, Ok \rangle}$
$\frac{}{\sigma \vdash_{ce_b} \mathbf{x} : \langle \sigma(\mathbf{x}), Ok \rangle} \quad \mathbf{x} \in Def(\sigma)$
$noDup(\mathbf{MD}_1^b \dots \mathbf{MD}_n^b) = \begin{cases} Ok & \text{if } \forall i, j \in 1..n \ \mathit{methSig}(\mathbf{MD}_i^b) = \mathit{methSig}(\mathbf{MD}_j^b) \implies i = j \\ VerifyError & \text{otherwise} \end{cases}$
$\mathit{methSig}(T \ m(T_1 \ x_1, \dots, T_n \ x_n) \ \{\mathbf{E}^b\}) = T \ m(T_1 \dots T_n)$

**Fig. 12.** Verification

of a source fragment, the number of arguments is indeed equal to the number of parameter types in the descriptor; such a mismatch may only be found in “malicious” binary fragments.

$\frac{\vdash_{ce_b} C : Ok \quad E^b \overset{*}{\rightsquigarrow}_{ce_b} V}{C \rightsquigarrow_{ce_b} V} \quad ce_b(C) = \langle -, \rightarrow, \rightarrow, E^b \rangle$	
$\frac{\vdash_{ce_b} C : \varepsilon}{C \rightsquigarrow_{ce_b} \varepsilon}$	
$\frac{\vdash_{ce_b} C : \varepsilon}{\boxed{C}^{Type} \rightsquigarrow_{ce_b} \varepsilon}$	
$\frac{\vdash_{ce_b} C : Ok}{\boxed{C}^{Type} \rightsquigarrow_{ce_b} \boxed{C}^{Type}}$	
$\frac{E^b \rightsquigarrow_{ce_b} E_1^b}{\boxed{E^b}^{Exp} \rightsquigarrow_{ce_b} \boxed{E_1^b}^{Exp}} \quad E_1^b \neq \varepsilon$	
$\frac{E^b \rightsquigarrow_{ce_b} \varepsilon}{\boxed{E^b}^{Exp} \rightsquigarrow_{ce_b} \varepsilon}$	
$\frac{}{\mathbf{new} C.m \langle C', \bar{T}, T \rangle (v_1, \dots, v_n) \rightsquigarrow_{ce_b} \varepsilon}$	$\bar{T} = T_1 \dots T_n$ $MBody(ce_b, C', m, \bar{T}, T) = \varepsilon$
$\frac{E^b[v_1/x_1, \dots, v_n/x_n] \rightsquigarrow_{ce_b} E_1^b}{\mathbf{new} C.m \langle C', \bar{T}, T \rangle (v_1, \dots, v_n) \rightsquigarrow_{ce_b} E_1^b}$	$\bar{T} = T_1 \dots T_n$ $MBody(ce_b, C', m, \bar{T}, T) \neq \varepsilon$ $MBody(ce_b, C, m, \bar{T}, T) = \langle E^b, x_1 \dots x_n \rangle$

**Fig. 13.** Rewriting

Figure 13 shows the rewriting rules for the program execution. The first two rules deal with the execution of the main method of a class  $C$ ; the former covers the case when class  $C$  is verified, whereas the latter considers the case when  $C$  does not pass verification.

The third and fourth rules cover the loading/verification process. The former is used in case of error: the whole term is rewritten in the exception thrown by the verifier. The latter is used when the verification is carried out successfully; in this case the term is rewritten in itself except for the reference to the class  $C$  that is unboxed.

The third rule is just the standard closure.

The fourth rule propagates an exception rewriting an entire term containing an exception  $\varepsilon$  in the exception itself.

The fifth and the sixth rules deal with method invocation. When the method cannot be found starting the search from the class contained in the method

descriptor the entire expression is rewritten in the exception; otherwise a second call to  $MBody$ , passing as starting class the dynamic type of the receiver, returns the method body and the name of the arguments. These information are used to expand the method call.

#### 4.1 Main Results

We prove three main theorems claiming the safety of source and binary type-checking and of the safe compilation schema, respectively; the former two theorems are necessary for proving the latter.

**Theorem 1 (Safe Source Typechecking).** *For all type environments  $\Gamma$ , sources  $S$  and binaries  $B$ , if  $\Gamma \vdash S \rightsquigarrow B$ , then  $\Gamma \vdash B \diamond$ .*

**Theorem 2 (Safe Binary Typechecking).** *Let  $\langle ce_b, ce_s \rangle$  and  $\mathcal{C}$  be a compilation environment and a class name, respectively. For all values  $V$ , if  $\forall \mathcal{C}' \in \mathcal{D}(\langle ce_b, ce_s \rangle, \{\mathcal{C}\}) \quad \mathcal{T}(\langle ce_b, ce_s \rangle) \vdash ce_b(\mathcal{C}') \diamond$  and  $\mathcal{C} \rightsquigarrow_{ce_b} V$ , then  $V$  is not an exception.*

We can state now the main property of the safe compilation schema: if a set of classes is successfully compiled w.r.t. the safe schema, then the execution of any binary produced by such compilation in the updated binary environment never throws a linkage exception.

**Theorem 3 (Safety).** *Let  $\langle ce_b, ce_s \rangle$ ,  $\mathcal{CS}$  and  $ce'_b$  be a compilation environment, a set of class names and a binary environment, respectively. For all class names  $\mathcal{C} \in \text{Def}(ce'_b)$  and values  $V$ , if  $\mathcal{C}(\langle ce_b, ce_s \rangle, \mathcal{CS}) = ce'_b$  and  $\mathcal{C} \rightsquigarrow_{ce_b[ce'_b]} V$ , then  $V$  is not an exception.*

#### 4.2 Formal Proofs (Sketched)

*Safe Source Typechecking:* For reasons of space we omit the proof, which is by induction over the rules defining the source typechecking judgment  $\Gamma \vdash S \rightsquigarrow B$ . The proof uses a number of lemmas claiming the same property for all kinds of subcomponents of a source; the most interesting one concerns expressions.

**Lemma 1.** *For all type and variable environments  $\Gamma$  and  $\Pi$ , source and binary expressions  $E^s$  and  $E^b$  and types  $T$ , if  $\Gamma; \Pi \vdash E^s : T \rightsquigarrow E^b$  then  $\Gamma; \Pi \vdash E^b : T \diamond$ .*

*Safe Binary Typechecking:* Safety of binary typechecking comes from the following two theorems, the former connecting static with dynamic binary typechecking (that is, the binary typechecking judgment  $\Gamma \vdash B \diamond$  with the verification judgment  $\vdash_{ce_b} \mathcal{C} : b$ ), the latter expressing subject reduction for binary expressions. For reasons of space we omit the proofs of these two theorems; the former can be proved by induction over the rules for binary typechecking, while the latter can be proved by induction over the rewriting rules for binary expressions.

**Theorem 4 (Binary Typechecking Implies Verification).** *Let  $\langle ce_b, ce_s \rangle$  and  $C$  be a compilation environment and a class name, respectively, s.t. the following condition holds:  $\mathcal{D}(\langle ce_b, ce_s \rangle, \{C\}) \subseteq \text{Def}(ce_b)$ . If  $\mathcal{T}(\langle ce_b, ce_s \rangle) \vdash ce_b(C) \diamond$ , then  $\vdash_{ce_b} C : Ok$ .*

Note that the converse implication does not hold, since typechecking at (dynamic) load time is more accurate than that at compile time. For instance, typechecking of a binary declaration of a class  $C$  requires the check of all classes explicitly mentioned in  $C$ , whereas the JVM only checks those classes that are actually needed by that particular execution.

**Theorem 5 (Binary Subject Reduction).** *Let  $\langle ce_b, ce_s \rangle$  and  $E^b$  be a compilation environment and a binary expression, respectively. If  $\mathcal{T}(\langle ce_b, ce_s \rangle); \emptyset \vdash E^b : T \diamond$  and  $E^b \rightsquigarrow_{ce_b} E^{b'}$ , then there exists a type  $T'$  s.t.  $\mathcal{T}(\langle ce_b, ce_s \rangle); \emptyset \vdash E^{b'} : T' \diamond$  and  $\mathcal{T}(\langle ce_b, ce_s \rangle) \vdash T' \leq T$ .*

We are now able to prove safety of binary typechecking. Let us assume that  $\forall C' \in \mathcal{D}(\langle ce_b, ce_s \rangle, \{C\}) \quad \mathcal{T}(\langle ce_b, ce_s \rangle) \vdash ce_b(C') \diamond$  and  $C \rightsquigarrow_{ce_b} V$ . From the first assumption we easily deduce  $\mathcal{D}(\langle ce_b, ce_s \rangle, \{C\}) \subseteq \text{Def}(ce_b)$  and  $\mathcal{T}(\langle ce_b, ce_s \rangle) \vdash ce_b(C) \diamond$  (since, trivially,  $C \in \mathcal{D}(\langle ce_b, ce_s \rangle, \{C\})$ ).

As a consequence, Theorem 4 can be applied, therefore  $\vdash_{ce_b} C : Ok$  holds. This means that  $C \rightsquigarrow_{ce_b} V$  has been deduced by instantiating the first (and not the second) meta-rule in Figure 13, so  $E^b \rightsquigarrow_{ce_b}^* V$  must hold. Furthermore, the validity of  $\mathcal{T}(\langle ce_b, ce_s \rangle) \vdash ce_b(C) \diamond$  implies the validity of  $\mathcal{T}(\langle ce_b, ce_s \rangle); \emptyset \vdash E^b : T \diamond$ , since there is only one meta-rule that can be instantiated in Figure 6. Therefore we can apply Theorem 5 and deduce the validity of  $\mathcal{T}(\langle ce_b, ce_s \rangle); \emptyset \vdash V : T' \diamond$ , with  $T'$  subtype of  $T$ . Since exceptions do not typecheck (see rules in Figure 6), we can conclude that  $V$  is not an exception.

*Safety:* To prove safety we need two lemmas claiming that both the dependency and type extraction functions are invariant w.r.t. source typechecking. These lemmas can be proved by induction over the definition of  $\mathcal{D}$  and  $\mathcal{T}$ , respectively; for reasons of space we omit their proofs.

In what follows, let  $ce_s \setminus C$  denotes the partial function obtained by restricting the definition domain of  $ce_s$  to the set  $\text{Def}(ce_s) \setminus C$ .

**Lemma 2.** *Let  $\langle ce_b, ce_s \rangle$ ,  $C$  and  $B$  be a compilation environment, a class name, and a binary fragment, respectively. If  $\mathcal{T}(\langle ce_b, ce_s \rangle) \vdash ce_b(C) \rightsquigarrow B$ , then for all class name  $C'$  the following equality holds:*

$$\mathcal{D}(\langle ce_b, ce_s \rangle, \{C'\}) = \mathcal{D}(\langle ce_b[C \mapsto B], ce_s \setminus C \rangle, \{C'\}).$$

**Lemma 3.** *Let  $\langle ce_b, ce_s \rangle$ ,  $C$  and  $B$  be a compilation environment, a class name, and a binary fragment, respectively. If  $\mathcal{T}(\langle ce_b, ce_s \rangle) \vdash ce_b(C) \rightsquigarrow B$ , then the following equality holds:*

$$\mathcal{T}(\langle ce_b, ce_s \rangle) = \mathcal{T}(\langle ce_b[C \mapsto B], ce_s \setminus C \rangle).$$

Now assume that  $\mathcal{C}(\langle ce_b, ce_s \rangle, \mathbf{CS}) = ce'_b$ . By virtue of the top-level rule in Figure 3, the following judgments are valid:

$$\begin{aligned} \forall \mathbf{C} \in \mathbf{CS}_b \quad \Gamma \vdash ce_b(\mathbf{C}) \diamond \\ \forall \mathbf{C} \in \mathbf{CS}_s \quad \Gamma \vdash ce_s(\mathbf{C}) \rightsquigarrow \mathbf{B}_{\mathbf{C}} \end{aligned}$$

where  $\mathbf{CS}_b = \mathbf{CS}_d \cap \text{Def}(ce_b)$ ,  $\mathbf{CS}_s = \mathbf{CS}_d \cap \text{Def}(ce_s)$ ,  $\mathbf{CS}_d = \mathcal{D}(\langle ce_b, ce_s \rangle, \mathbf{CS})$  and  $\Gamma = \mathcal{T}(\langle ce_b, ce_s \rangle)$ . Furthermore,  $ce'_b(\mathbf{C}) = \{\mathbf{C} \mapsto \mathbf{B}_{\mathbf{C}} \mid \mathbf{C} \in \mathbf{CS}_s\}$ .

By Theorem 1,  $\Gamma \vdash ce'_b(\mathbf{C}) \diamond$  for all  $\mathbf{C} \in \mathbf{CS}_s$ , therefore we can easily deduce  $\Gamma \vdash ce_b[ce'_b](\mathbf{C}) \diamond$  for all  $\mathbf{C} \in \mathbf{CS}_d$ .

Let us now prove the main theorem by assuming that  $\mathbf{C}$  is a class name in  $\mathbf{CS}_b$  (recall that  $\text{Def}(ce'_b) = \mathbf{CS}_b$ ) and that  $\mathbf{C} \rightsquigarrow_{ce_b[ce'_b]} V$  for a certain value  $V$ . By lemmas 2 and 3 and by induction on the cardinality of  $\mathbf{CS}$ ,  $\mathcal{D}(\langle ce_b, ce_s \rangle, \{\mathbf{C}\}) = \mathcal{D}(\langle ce_b[ce'_b], ce_{s \setminus \mathbf{CS}} \rangle, \{\mathbf{C}\})$  and  $\mathcal{T}(\langle ce_b, ce_s \rangle) = \mathcal{T}(\langle ce_b[ce'_b], ce_{s \setminus \mathbf{CS}} \rangle)$ . Therefore we can apply Theorem 2 and conclude that  $V$  cannot be an exception.

## 5 Related Work

We already mentioned in the Introduction that the seminal paper on separate typecheck of fragments is [8]. There, the basic idea is to distinguish a phase of *intra-checking*, which models separate compilation, in which a single fragment is type-checked w.r.t. a typing environment (which expresses the interface of the fragment in terms of both imported and exported services), and a phase of *inter-checking* which models (static) linking, in which it is checked that all the fragments we want to link have been type-checked w.r.t. compatible type environments.

Formally<sup>12</sup>, intra-checking is modeled by a judgment  $\Gamma \vdash f : T$  (in [8] issues of code generation are avoided by always working at the source level, even when discussing linking), expressing that the fragment  $f$  has type  $T$  in the type environment  $\Gamma$ . Inter-checking takes place on *linksets* which are, roughly, collections of named fragment  $x_i \mapsto \Gamma_i \vdash f_i : T_i^{i \in 1..n}$ , and succeeds if and only if intra-checking succeeds (that is, each  $\Gamma_i \vdash f_i : T_i$  holds) and, moreover, for each  $j, k \in 1..n$ ,  $x_j$  has type  $T_j$  in  $\Gamma_k$ . This corresponds to require *exact* agreement among the actual interface of a fragment and that assumed in another: in realistic systems, this condition should be weakened, for instance requiring some subtyping relation (see below for the Java case).

As already discussed in the Introduction, Java has many features which make this view not immediately applicable: class files play the dual roles of interfaces (type environments) and object files; there is no separate linking phase, since linking takes place at run-time; compilers usually incorporate *some* inter-checks, but not enough to guarantee safe run-time linking.

In this paper, we present a framework which models separate compilation in the Java sense. Though a detailed formal comparison of our framework with

<sup>12</sup> We use slightly different notations from [8] in order to conform to those used in this paper.

that in [8] is matter of further work, we can list the main contact points and differences.

- The intra-checking phase is modeled in the same way, apart from the fact that here we are interested in distinguishing source and code fragments, hence in modeling code generation.
- Instead of starting from a fixed linkset, here we assume that the set of fragments to be linked is determined by the dependency function.
- Most importantly, instead of having that each fragment  $f_i$  is equipped with its own type environment  $T_i$  (expressing its interface), here we assume a *global* type environment. This reflects the fact that in Java there is no separate notion of interface<sup>13</sup> of a fragment (class  $C$ ) describing both imported and exported services, but this interface must in some way be extracted from the code. Here, since our aim was to model separate compilation as it happens in Java rather than to compare with truly separate compilation in the sense of [8], we have taken the simpler approach to type-check all the fragments in the same type environment, which trivially consists in the compilation environment where we drop method bodies. Hence inter-checking is trivial. An approach more in the spirit of [8] consists in extracting the interface  $I_i$  for each class  $C_i$  containing the minimal assumptions needed for successfully typecheck  $C_i$ , and then to check that different interfaces are compatible. This approach poses non trivial problems both on how to perform the type extraction and how to define the “right” subtyping relation among interfaces. However, the investigation is very interesting because it could lead to innovative techniques for Java compilation supporting truly separate compilation and then the possibility of performing static checks on binaries. We refer to [6, 2] for a more complete treatment in this direction.

Another important stream of research related to this paper is that devoted to the formal definition of Java (see [1] for a survey). As already mentioned, the type judgment which we consider at the source level is based (though much simpler, since we consider a small Java subset) on the many existing formal Java type systems, in particular those in [10]. For what concerns an integrated formal model covering all Java aspects, the most remarkable amount of work in this direction is that of Sophia Drossopoulou and her group. The already cited [10] provides a formal type system at the source level for a substantial subset Java<sup>s</sup> of Java and a translation of this language into a binary language Java<sup>b</sup>, which is in turn a subset of a language Java<sup>r</sup> of run-time terms for which an operational semantics is given. This allows to prove type safety of the Java subset. In [13] the focus is on binary compatibility. In [9] a model is defined for dynamic loading and linking, distinguishing five components in a Java implementation: evaluation, resolution, loading, verification, and preparation, with their associated checks. These five together are proved to guarantee type soundness. This paper is the most important reference for the execution model of our small binary language; however, in our case the main aim is not to define a realistic model of the JVM,

---

<sup>13</sup> Not to be confused with a Java interface.

taking into account all features, but to show how absence of linkage errors can be guaranteed by a compilation schema, so we take a much more abstract view. Finally, [12] enhances the previous formal description of Java in [10], introducing, among other improvements, an account of separate compilation. Indeed, type information used in typechecking Java<sup>s</sup> can also be extracted from the binary language Java<sup>r</sup>, analogously to what we do in this paper by means of the  $\mathcal{T}$  function. However, the judgment for typechecking source classes defined in [12] do not correspond to separate compilation as happens in our framework, simply because its validity requires the type environment extracted from the compilation environment to be well-formed.

Finally, several interesting papers can be found in literature on separate compilation for ML (see among many others [17, 16, 7]). All these papers clearly show that separate compilation in ML is not a simple issue, and for this reason, needs to be properly formalized. However, ML separate compilation is based on traditional static linking, therefore many problems arising in Java disappear in ML; for instance, the static type-checks proposed in [17] are sensible for a static linker, but cannot be performed at run-time by a virtual machine without seriously compromising efficiency. Furthermore, it seems that no unifying frameworks have been defined for investigating ML separate compilation, and in fact, this would be useful to compare all the technical results and to understand how they can contribute all together to the design of a better compiler/linker for ML. For instance, using the terminology used in our paper to model the overall compilation process, [17] is mainly concerned with the definition of the type extraction function, while [16] with the typechecking of sources and [7] with the definition of the dependency function.

## 6 Conclusion

We have introduced a formal framework modeling Java separate compilation. The overall compilation process is modeled by the formal notion of compilation schema, in which the aspects which concern truly separate typechecking of fragments (source type judgment and binary type judgment) are isolated from the definition of dependencies and extraction of the type information from the fragments. We have considered three compilation schemata of interest for Java, that is, minimal, SDK, and safe, correspondingly to perform, when a single fragment is typechecked, no typechecks on other fragments, only the typechecks performed by SDK, and enough typechecks in order to ensure absence of linkage errors at run-time.

In order to demonstrate our approach, we have defined a kernel model for Java separate compilation and execution, consisting in a small Java subset, and a simple corresponding binary language for which we provide an operational semantics including run-time verification. We have defined a safe compilation schema for this language and formally proved type safety.

In this paper we have focused on the safety property (a preliminary work pointing out that Java compilers are not safe in the context of separate com-

pilation was [4]); however, there are other interesting properties of compilation schemata we want to investigate, like *monotonicity* and *contextual binary compatibility*. By monotonicity we mean the fact that, when a subset of the source fragments composing a program is changed, re-compiling only this set gives the same result as re-compiling the whole program (this property is mentioned as desirable in [8], and a preliminary formalization is given in [5]). By *contextual binary compatibility* we mean a property analogous to Java binary compatibility, but related to compilation: we say that a compilation schema respects contextual binary compatibility if all the binary fragments which could be safely linked before compilation still safely link after. Note that, even though a formal analysis of their relation is still matter of further work (see below), these three properties seem at a first sight to be independent.

The work presented in this paper is a first step, and many interesting research directions are open. On the more theoretical side, we plan to formally analyze the difference between separate compilation in the Java sense, modeled in this paper, and truly separate compilation in the sense of, e.g., [8] (see [6, 2]). This should lead to the definition of a very abstract framework, like that in [8], but including dynamic linking and verification, in which we will express formal properties like type safety, monotonicity and contextual binary compatibility and analyze their relations. Furthermore, the interesting relation between verification and binary typechecking deserves further investigation in order to better understand the deep issue of binary compatibility.

On the side of application to Java, we plan to extend the safe type system defined here to more substantial Java subsets and to develop extended compilers which satisfy good properties like type safety.

*Acknowledgment* We are extremely grateful to Sophia Drossopoulou for her stimulating discussions and for her precious suggestions on preliminary drafts of this paper.

## References

1. J. Alves-Foss, editor. *Formal Syntax and Semantics of Java*. Number 1523 in Lecture Notes in Computer Science. Springer, 1999.
2. D. Ancona and G. Lagorio. Supporting true separate compilation in Java: A modular approach. Technical report, Dipartimento di Informatica e Scienze dell'Informazione, Università di Genova, 2002. Submitted for publication.
3. D. Ancona, G. Lagorio, and E. Zucca. A core calculus for Java exceptions. In *ACM Symp. on Object-Oriented Programming: Systems, Languages and Applications 2001*. ACM Press, October 2001. SIGPLAN Notices.
4. D. Ancona, G. Lagorio, and E. Zucca. Java separate type checking is not safe. In *3th Intl. Workshop on Formal Techniques for Java Programs 2001*, June 2001.
5. D. Ancona, G. Lagorio, and E. Zucca. Separate compilation in Java: Avoiding ambiguity via monotonicity. Technical Report, DISI, July 2001.
6. D. Ancona, G. Lagorio, and E. Zucca. True separate compilation of Java classes. Technical report, Dipartimento di Informatica e Scienze dell'Informazione, Università di Genova, 2002. Submitted for publication.

7. M. Blume. Dependency analysis for standard ML. *ACM Transactions on Programming Languages and Systems*, 21(4):790–812, 1999.
8. L. Cardelli. Program fragments, linking, and modularization. In *ACM Symp. on Principles of Programming Languages 1997*, pages 266–277. ACM Press, 1997.
9. S. Drossopoulou. Towards an abstract model of Java dynamic linking and verification. In R. Harper, editor, *TIC'00 - Third Workshop on Types in Compilation (Selected Papers)*, volume 2071 of *Lecture Notes in Computer Science*, pages 53–84. Springer, 2001.
10. S. Drossopoulou and S. Eisenbach. Describing the semantics of Java and proving type soundness. In J. Alves-Foss, editor, *Formal Syntax and Semantics of Java*, number 1523 in *Lecture Notes in Computer Science*, pages 41–82. Springer, 1999.
11. S. Drossopoulou, S. Eisenbach, and D. Wragg. A fragment calculus - towards a model of separate compilation, linking and binary compatibility. In *Proc. 14th Ann. IEEE Symp. on Logic in Computer Science*, July 1999.
12. S. Drossopoulou, T. Valkevych, and S. Eisenbach. Java type soundness revisited. Technical report, Dept. of Computing - Imperial College of Science, Technology and Medicine, September 2000.
13. S. Drossopoulou, D. Wragg, and S. Eisenbach. What is Java Binary Compatibility? In *ACM Symp. on Object-Oriented Programming: Systems, Languages and Applications 1998*, volume 33(10) of *ACM SIGPLAN Notices*, pages 341–358, October 1998.
14. J. Gosling, B. Joy, G. Steele, and G. Bracha. *The Java™ Language Specification, Second Edition*. Addison-Wesley, 2000.
15. A. Igarashi, B. Pierce, and P. Wadler. Featherweight Java: A minimal core calculus for Java and GJ. In *ACM Symp. on Object-Oriented Programming: Systems, Languages and Applications 1999*, pages 132–146, November 1999.
16. X. Leroy. Manifest types, modules and separate compilation. In *ACM Symp. on Principles of Programming Languages 1994*, pages 109–122. ACM Press, 1994.
17. Z. Shao and A.W. Appel. Smartest recompilation. In *ACM Symp. on Principles of Programming Languages 1993*, pages 439–450. ACM Press, 1993.
18. D. Syme. Proving Java type sound. In Jim Alves-Foss, editor, *Formal Syntax and Semantics of Java*, number 1523 in *Lecture Notes in Computer Science*, pages 83–118. Springer, 1999.
19. D. von Oheimb and T. Nipkow. Machine-checking the Java specification: Proving type-safety. In Jim Alves-Foss, editor, *Formal Syntax and Semantics of Java*, number 1523 in *Lecture Notes in Computer Science*, pages 119–156. Springer, 1999.

## A Appendix

$$\begin{aligned}
refClasses(\mathbf{C}, \langle ce_b, ce_s \rangle) &= \begin{cases} refClasses(ce_b(\mathbf{C})) & \text{if } \mathbf{C} \in Def(ce_b) \\ refClasses(ce_s(\mathbf{C})) & \text{if } \mathbf{C} \in Def(ce_s) \\ \emptyset & \text{otherwise} \end{cases} \\
refClasses(\mathbf{class } \mathbf{C} \text{ extends } \mathbf{C}' \{ \mathbf{MDS}^s \} \text{ main } \mathbf{E}^s) &= \\
&\quad \{ \mathbf{C}, \mathbf{C}' \} \cup refClasses(\mathbf{MDS}^s) \cup refClasses(\mathbf{E}^s) \\
refClasses(\mathbf{MD}_1^s \dots \mathbf{MD}_n^s) &= \bigcup_{i \in 1..n} refClasses(\mathbf{MD}_i^s) \\
refClasses(\mathbf{MH} \{ \text{return } \mathbf{E}^s; \}) &= refClasses(\mathbf{MH}) \cup refClasses(\mathbf{E}^s) \\
refClasses(\mathbf{T}_0 \text{ m}(\mathbf{T}_1 \mathbf{x}_1, \dots, \mathbf{T}_n \mathbf{x}_n)) &= \{ \mathbf{T}_0, \dots, \mathbf{T}_n \} \\
refClasses(\mathbf{new } \mathbf{C}) &= \{ \mathbf{C} \} \\
refClasses(\mathbf{x}) = refClasses(\mathbf{N}) &= \emptyset \\
refClasses(\mathbf{E}_0^s \text{ m}(\mathbf{E}_1^s, \dots, \mathbf{E}_n^s)) &= \bigcup_{i \in 0..n} refClasses(\mathbf{E}_i^s) \\
\\
refClasses(\langle \mathbf{C}, \mathbf{C}', \mathbf{MDS}^b, \mathbf{E}^b \rangle) &= \{ \mathbf{C}, \mathbf{C}' \} \cup refClasses(\mathbf{MDS}^b) \cup refClasses(\mathbf{E}^b) \\
refClasses(\mathbf{MD}_1^b \dots \mathbf{MD}_n^b) &= \bigcup_{i \in 1..n} refClasses(\mathbf{MD}_i^b) \\
refClasses(\mathbf{MH} \{ \text{return } \mathbf{E}^b; \}) &= refClasses(\mathbf{MH}) \cup refClasses(\mathbf{E}^b) \\
refClasses(\mathbf{E}_0^b \text{ m} \prec \boxed{\mathbf{C}}, \mathbf{T}_1 \dots \mathbf{T}_n, \mathbf{T}_0 \succ (\mathbf{E}_1^b, \dots, \mathbf{E}_n^b)) &= \{ \mathbf{C} \} \cup \bigcup_{i \in 0..n} (refClasses(\mathbf{E}_i^b) \cup \{ \mathbf{T}_i \}) \\
refClasses(\mathbf{new } \boxed{\mathbf{C}}) &= \{ \mathbf{C} \}
\end{aligned}$$

**Fig. 14.** Definition of the dependency function

$$\begin{aligned}
\forall \mathbf{C} \mathcal{T}(\langle ce_b, ce_s \rangle)(\mathbf{C}) &= \begin{cases} \langle \mathbf{C}', \mathcal{T}(\mathbf{MDS}^b) \rangle & \text{if } ce_b(\mathbf{C}) = \langle \mathbf{C}, \mathbf{C}', \mathbf{MDS}^b, \mathbf{E}^b \rangle \\ \langle \mathbf{C}', \mathcal{T}(\mathbf{MDS}^s) \rangle & \text{if } ce_s(\mathbf{C}) = \mathbf{class } \mathbf{C} \text{ extends } \mathbf{C}' \{ \mathbf{MDS}^s \} \text{ main } \mathbf{E}^s \\ \perp & \text{otherwise} \end{cases} \\
\mathcal{T}(\mathbf{MD}_1^s \dots \mathbf{MD}_n^s) &= \mathcal{T}(\mathbf{MD}_1^s) \dots \mathcal{T}(\mathbf{MD}_n^s) \\
\mathcal{T}(\mathbf{MD}_1^b \dots \mathbf{MD}_n^b) &= \mathcal{T}(\mathbf{MD}_1^b) \dots \mathcal{T}(\mathbf{MD}_n^b) \\
\mathcal{T}(\mathbf{MH} \{ \text{return } \mathbf{E}^s; \}) &= \mathcal{T}(\mathbf{MH}) \\
\mathcal{T}(\mathbf{MH} \{ \text{return } \mathbf{E}^b; \}) &= \mathcal{T}(\mathbf{MH}) \\
\mathcal{T}(\mathbf{T}_0 \text{ m}(\mathbf{T}_1 \mathbf{x}_1, \dots, \mathbf{T}_n \mathbf{x}_n)) &= \mathbf{T}_0 \text{ m}(\mathbf{T}_1 \dots \mathbf{T}_n)
\end{aligned}$$

**Fig. 15.** Definition of the type extraction function