

# CIL to Java-bytecode Translation for Static Analysis Leveraging

Pietro Ferrara  
Julia SRL, Italy  
pietro.ferrara@juliasoft.com

Agostino Cortesi  
Università Ca' Foscari Venezia, Italy  
cortesi@unive.it

Fausto Spoto  
Università di Verona, Italy  
spoto@univr.it

## ABSTRACT

A formal translation of CIL (*i.e.*, .Net) bytecode into Java bytecode is introduced and proved sound with respect to the language semantics. The resulting code is then analyzed with Julia, an industrial static analyzer of Java bytecode. The overall process of translation and analysis is fast, scales up to industrial programs, and introduces a negligible number of false alarms. The main result of this work is to leverage existing, mature, and sound analyzers for Java bytecode by applying them to the (translated) CIL bytecode.

## KEYWORDS

Static Analysis, bytecode, JVM, .Net, CIL.

## 1 INTRODUCTION

Static analysis infers, at compile-time, properties about the runtime behavior of computer programs. It allows one to verify, for instance, the absence of runtime errors or security breaches. Static analysis applies also to compiled code in assembly or bytecode format. This is particularly interesting for applications distributed on the Internet, or downloaded from public (and possibly unsafe) application repositories (*e.g.*, the Google Play Store), when the source code is not available, but the user would like to statically check some safety or security properties.

The analysis of Java bytecode for the Java Virtual Machine (from now on, JB) has a long research tradition and many analyzers exist [24]. Some analyses build on formal mathematical roots, such as abstract interpretation [13, 20, 23]. Moreover, JB makes the design of static analysis easier by requiring bytecode to be type-checkable [19] and without unsafe operations such as free pointer operations. On the contrary, CIL bytecode, that is, the compiled bytecode used for the .Net platform (from now on, just CIL), has not received much attention from the static analysis community yet. Moreover, CIL can be used in an *unsafe* way, that is, allowing free pointer operations, which makes its static analysis harder. However, these operations are very often used in very controlled contexts, hence, in most cases, a static analyzer could possibly capture their actual behavior anyway.

Despite clear differences, JB and CIL share strong similarities, being both low-level object-oriented languages where objects are stored and shared in the heap. Hence, it is tempting to leverage mature existing static analyses and tools for JB by translating CIL into *equivalent* JB and running the tools on the latter. Obviously, this introduces issues about the exact meaning of *equivalence* between CIL and its translation into JB. Moreover, the translation should not introduce code artifacts that confuse the analyzer and should work on industrial-size CIL applications, supporting as many unsafe pointer operations as possible.

The main contribution of this work is the introduction of a translation of CIL to JB that is (i) theoretically sound, and (ii) effective in practice, so that an industrial static analyzer for JB can be applied to

.Net (and in particular C#) programs. More languages compile into JB (*e.g.*, Scala) and CIL (*e.g.*, VB.Net and F#), with distinct features and code structures. Here, we focus on Java and C#, that have similar structure and compile into comparable bytecode.

We start by formalizing the concrete semantics of a representative subset of CIL and JB, and the translation of CIL into JB. Then, we prove this translation sound, that is, the concrete semantics of the initial CIL program is equivalent to that of the translated JB program. This guarantees that, if we prove a property of the JB program, then such property holds also for the original CIL program. Then we present a deep experimental evaluation over industrial-size open source popular programs, by applying the Julia static analyzer [23] to the translated JB.

We focus on three main research issues about scalability, precision, and coverage of the overall approach:

**RESEARCH QUESTION 1 (SCALABILITY).** *Does the CIL to JB translation scale up, that is, (i) can it deal with libraries of industrial size (100KLOCs) in a few minutes, and (ii) is its computational time comparable to that required by the static analysis phase?*

**RESEARCH QUESTION 2 (PRECISION).** *Does the CIL to JB translation introduce less than 10% of the false alarms produced by the static analyzer?*

**RESEARCH QUESTION 3 (LIBRARIES).** *Despite supporting only a subset of CIL, does the CIL to JB translation succeed on at least 95% of the system libraries?*

Notice that it is crucial that a static analyzer understands the behaviour of system libraries, as otherwise it could only rely on manual annotations or on (possibly unsound) assumptions on their execution. System libraries need to access memory through unsafe pointer. Java allows such behaviour through native methods (written in languages other than Java and bound through the Java Native Interface), while .Net allows unsafe pointers in its code. In these cases, our translation produces Java native methods. We ensure that the effort of manually annotating .Net libraries is comparable to that needed for Java.

## 1.1 Related Work

Few attempts have been made in the past to translate CIL to JB. Grasshopper is probably the most popular one. However, it is not available any more<sup>1</sup>. As far as we can see, it was abandoned about a decade ago, and we cannot make any comparison with our translation. A similar tool was CLR2JVM [2]: it translates CIL to an intermediate  $XML_{CLR}$  representation, that can be then translated into  $XML_{JVM}$ , and finally to JB. As far as we can see<sup>2</sup>,

<sup>1</sup>We were unable to access the website <http://dev.mainsoft.com>, that seems to be the website of the tool from past forum discussions (<http://stackoverflow.com/questions/95163/differences-between-msil-and-java-bytecode>)

<sup>2</sup><http://xmlvm.org/documentation/>

the tool should read .NET executables, but it failed parsing all the executable files of our experiments (see Sec. 5 for the complete list). This probably happened because CLR2JVM is not maintained any more (the last commit to the repository <https://sourceforge.net/p/xmlvm/code/HEAD/tree/trunk/xmlvm/src/clr2jvm/> was more than six years ago), and it does not support the last CIL versions. Neither Grasshopper nor CLR2JVM has any documentation or discussion about how the translation is performed (in particular, how they handle instructions that are different from CIL to JB such as direct references). Therefore, as far as we can see our translation from CIL to JB is the only one that (i) works on recent releases of CIL and JB, and (ii) is formalized and proved sound.

Other translations between low-level languages exist, justified by the need of applying verification tools that work on a specific language only. For instance, [9] defines a translation from Boogie into WhyML and proves it sound, as we have done from CIL to JB. Similar translations work also at runtime, in particular inside a just-in-time compiler, as in [12]. However, we did not find any literature on the translation of CIL into JB for industrial-size software.

Many other static analysis tools for the .Net platform exist, in particular for C# code. There are tools that verify compliance to some guideline, such as Fxcop [22] and Coverity Prevent [15]. Other tools, such as NDepend [5] and CodeMetrics [6], provide metrics about the code under analysis. ReSharper [18] applies syntactical code inspections, finds code smells and guarantees compliance to coding standards. As far as we can see, there exist only two main fundamental tools with scientific base: (i) Spec# [11], an extension of C# with static checking of various kinds of manual specifications, and (ii) CodeContracts [21], an abstract interpretation-based static analyzer for CIL. In the Java world, the number of static analyzers based on syntactic reasoning, such as Checkstyle [1], FindBugs [4] and PMD [7], is comparable to that for .Net. However, Java attracted much more attention from the scientific community, and more semantic analyzers have been introduced during the last decade, such as CodeSonar [3], ThreadSafe [10] and Julia [23]. Few semantic analyzers, such as WALA [8], have been applied to various languages (e.g., Java and JavaScript), but with ad-hoc source translations.

Our approach lets us apply all the Java analyzers on .Net programs (almost) *for free*, that is, by translating CIL into JB and using the analyzers as they are (we expect that few manual annotations are needed to improve the precision of the analysis, in particular when dealing with library calls). We have also studied performance and results with Julia. As far as we know, our work is the first translation of CIL into JB for static analysis that is proven to be *sound* and comes with evidence that this translation applies to *industrial-size software* with results that are comparable in terms of precision and efficiency to those obtained on JB.

## 2 BACKGROUND

We provide here a basic introduction to CIL and JB, with a running example and discussion on the architecture of the Julia static analyzer. For an exhaustive definition of JB and CIL, see [19] and [16], respectively.

### 2.1 CIL and JB

Bytecode is a machine-independent low-level programming language, used as target of the compilation of high-level languages,

JB	CIL	
iadd ladd	add	(arith. op.)
iload i lload i aload i istore i lstore i astore i	ldloc i stloc i ldarg i	(local vars)
invokevirtual invokestatic	call	(meth. call)
new T getfield f putfield f	newobj T(...) ldfld f stfld f	(objects)
if.icmpgt	bgt	(cond. branch)
dup dup2	dup	(stack)
	ldloca i stind ldind	(pointers)

Figure 1: JB and CIL minimal bytecode languages.

that hence becomes machine-independent. Bytecode languages are interpreted by their corresponding virtual machine, specific to each execution architecture. Both .Net and Java compile into bytecode. However, they use distinct instructions and virtual machines. .Net compiles into CIL, while Java compiles into JB. These have strong similarities: both use an operand stack for temporary values and an array of local variables standing for source code variables; both are object-oriented, with instructions for object creation, field access and virtual method dispatch. Despite these undeniable similarities, CIL and JB differ for the way of performing parameter passing (CIL uses a specific array of variables for the formal parameters, while JB merges them into the array of local variables); they handle object creation differently (CIL creates and initializes the object at the same time, while these are distinct operations in JB); they allocate memory slots differently (in CIL each value uses a slot, while JB uses 1 or 2 slots per 32- or 64- bit values, respectively); finally, CIL uses pointers explicitly, also in type-unsafe ways, while JB has no notion of pointer.

We focus our formalization on a minimal representative subset of JB and CIL, as defined in Fig. 1. That figure presents bytecode instructions for:

**arithmetic:** JB has type-specific operations, such as `iadd` and `ladd` to add two integer or long values, respectively. Instead, CIL has generic operations, such as `add` to add two numerical values of the same type;

**local variables access:** JB has a single array of variables for both local variables and method arguments, and reads and writes values from this array through `xload` and `xstore`, where `x` is `i` for integer values, `l` for long values and `a` for references, respectively. In this array, 64 bits values use two subsequent slots. CIL, instead, uses two arrays: one for method’s arguments (`ldarg i` loads the value of the `i`-th argument) and one for local variables (`ldloc i` and `stloc i` read and write the `i`-th local variable, respectively). In addition, it uses one slot both for 32- and 64-bit variables;

**method call:** JB has several kinds of method call instructions, such as `invokevirtual` and `invokestatic`. Instead, CIL has a unique `call` instruction;

**object manipulation:** in JB, instructions `new`, `getField`, and `putField` allocate a new object, read, and write its fields, respectively. CIL has similar instructions `newobj`, that also calls the constructor, `ldfld` and `stfld`;

**conditional branch:** JB has type-specific conditional branch instructions such as `if_icmpgt` (to branch if the greater than operator returns true on the topmost two integer values of the operand stack); CIL has generic instructions such as `bgt`;

**stack:** CIL duplicates the top value of the stack through the `dup` instruction. JB does the same with `dup` for 32 bits values and `dup2` for 64 bits values;

**pointers:** CIL contains some instructions to load the address of a local variable (`ldloca i`), and to store and load a value into the memory cell pointed by a reference (`stind` and `ldind`, respectively). Instead, JB has no direct pointer manipulation.

In the rest of this article,  $St_{CIL}$  and  $St_{JB}$  denote CIL and JB instructions or statements, respectively. A method (both in CIL and JB) is represented by (i) a sequence of (possibly conditional and branching) statements, and (ii) the number and static types of arguments and local variables.

## 2.2 Running Example

Fig. 2 shows the running example that Sec. 3 and 4 use to clarify the formalization. The C# code in Fig. 2a defines a class `Wrapper` that wraps an integer value, and a static method `WrappersCollection` that, given an integer `n`, returns a collection of `n` wrappers containing values from 0 to `n - 1`. Fig. 2b presents the (simplified) CIL obtained from its compilation: as usual with CIL, code is unstructured (e.g., there are branches at lines 7 and 20), and each source code statement could be translated into many bytecode statements (e.g., line 7 in Fig. 2a compiles into lines 3 and 4 in Fig. 2b). Fig. 2c presents the results of our translation of CIL into JB. Next sections explain the steps of the translation. First of all, notice some of the differences highlighted in Sec. 2.1. Namely, the type-generic CIL statement `stloc.1` at line 6 of Fig. 2b is translated into the type-specific `istore.2` at line 7 of Fig. 2c. Similarly, `newobj` (line 11) is translated into multiple JB statements (line 12-16). The running example contains few instructions that are not part of the minimal language defined in Sec. 2.1, and in particular (i) `ldc` and `i_const` that load constant (integer) values, (ii) `blt` and `if_icmplt` for conditional branching when an (integer) value is strictly less than another, (iii) `invokespecial` to invoke a specific method in JB, and (iv) `ret` and `areturn` to return a (reference) value.

**2.2.1 Example with Direct References.** Safe C# code adopts direct pointers only for `out` and `ref` method parameters. These parameters can be assigned (and read as well in case of `ref`) inside the method, and “any change to the parameter in the called method is reflected in the calling method”<sup>3</sup>. In our translation, we build up wrapper objects to soundly represent their semantics.

Consider for instance the C# code in Fig. 3a. Method `init` receives a `ref` parameter and it increments it by one. This is compiled (Fig. 3b) into a method reading the value pointed by the direct reference (line 5), and writing it (line 8). Our goal is to translate this code into the Java code in Fig. 3c: we simulate the direct reference by constructing a wrapper object (line 7), assigning

<sup>3</sup><https://msdn.microsoft.com/en-us/library/14akc2c7.aspx>

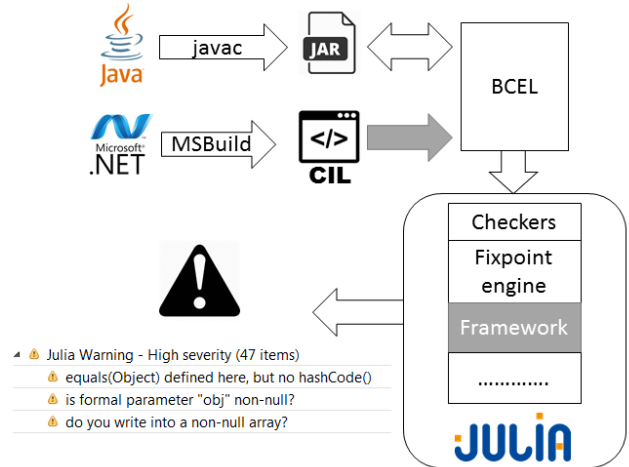


Figure 4: Julia’s architecture.

the value of the local variable to field `value` of the wrapper (line 8), and then propagating back the results of the call to `init` (line 9) by assigning the value of the local variable with the one stored in the wrapper (line 10). In addition, the `ref` parameter is replaced by the type of the wrapper object.

## 2.3 Julia

Julia [23] is a static analyzer for JB, based on abstract interpretation [14]. It transforms JB to basic blocks of code, that analyzes through a fixpoint algorithm. Analyses are constraint-based or denotational. Currently, Julia features around 70 checkers, including nullness, termination, synchronization and taint analysis. Since Julia works on JB, it may analyse any programming language that compiles into that bytecode. In particular, we apply Julia to the compilation of CIL into JB.

Fig. 4 is a high level view of Julia’s architecture. Java code is compiled by `javac` into a `jar` file, then parsed through the BCEL library [17]. Julia receives this latter format, applies its analysis (by using many components such as the checkers, that define what properties to check, a fixpoint engine, and the framework specifying the semantics of some specific components of the programming language), and outputs a list of warnings. The added component in our approach is the translation of CIL into BCEL format (grey arrow in the upper part of Fig. 4). Since a program in BCEL format can be dumped into a `jar` file, we can dump a `.dll` file in this format. We added few annotations in a new framework of Julia to specify the main components of the .Net framework (e.g., the signatures in the `Object` class).

## 3 CONCRETE SEMANTICS

We formalize here the concrete semantics of CIL and JB.

### 3.1 Notation

Let `Ref` and `Num` be the set of reference and numerical values, respectively, and let `Val = Ref ∪ Num`. A stack `s` of elements in `T` is a function in  $\mathbb{N} \rightarrow T$  such that  $\exists i \in \mathbb{N} : \forall i_1 \leq i : i_1 \in \text{dom}(s) \wedge \forall i_2 > i : i_2 \notin \text{dom}(s)$ . We will refer to `i` as the height of `s` (`height(s)`). Given a stack `s` and an element `e`, `s :: e` denotes a

```

1 public class Wrapper
2 {
3     int f;
4     Wrapper(int f) { this.f = f; }
5     static ICollection<Wrapper> WrappersCollection(int n)
6     {
7         ICollection<Wrapper> result = new List<Wrapper>();
8         for (int i = 0; i < n; i++)
9             result.Add(new Wrapper(i));
10        return result;
11    }
12 }

```

(a) The C# source code of the running example.

```

1 static ICollection<Wrapper> WrappersCollection(int n)
2 {
3     newobj List<Wrapper>::ctor()
4     stloc.0
5     ldc.0
6     stloc.1
7     br #18
8
9     ldloc.0
10    ldloc.1
11    newobj Wrapper::ctor(int32)
12    call Generic.ICollection<Wrapper>::Add
13    ldloc.1
14    ldc.1
15    add
16    stloc.1
17
18    ldloc.1
19    ldarg.0
20    blt #9
21
22    ldloc.0
23    ret
24 }

```

(b) The compilation into CIL of the code in (a).

```

1 static WrappersCollection(IListSystem/Collections/Generic/ICollection {
2     new <System/Collections/Generic/List>
3     dup
4     invokespecial <System/Collections/Generic/List.<init>>
5     astore.1
6     iconst.0
7     istore.2
8     goto 23
9
10    aload.1
11    iload.2
12    istore.3
13    new <Wrapper>
14    dup
15    iload.3
16    invokespecial <Wrapper.<init>>
17    invokevirtual <System/Collections/Generic/ICollection.Add>
18    iload.2
19    iconst.1
20    iadd
21    istore.2
22
23    iload.2
24    iload.0
25    if.icmplt 10
26
27    aload.1
28    areturn
29 }

```

(c) The translation into JB of the CIL in (b).

Figure 2: The C# code, CIL, and JB of the running example.

```

1 void init(ref int i)
2 {
3     i++;
4 }
5
6 void run()
7 {
8     int i=0;
9     init(ref i);
10 }

```

(a) C# code

```

1 static void init(ref A& a)
2 {
3     ldarg.0
4     ldarg.0
5     lldind.i4
6     ldc.i4.1
7     add
8     stind.i4
9 }
10
11 public static int run()
12 {
13     ldc.i4.0
14     stloc.0
15     ldloca.s 0
16     call void Temporary.Foo::'init'(int32&)
17 }

```

(b) CIL

```

1 public static final void init(WrapRef i) {
2     i.value = i.value + 1;
3 }
4
5 public static final int run() {
6     int i = 0;
7     WrapRef wrapper = new WrapRef();
8     wrapper.value = i;
9     init(wrapper);
10    i = wrapper.value;
11 }

```

(c) Java code

Figure 3: CIL code using `ref` parameters.

stack whose top element (that is, the one with the highest index) is  $e$  followed by the stack  $s$ .

As usual for object-oriented programming languages, an object is a map from field names to values, and a heap is a map from references to objects. Formally,  $\text{Heap} : \text{Ref} \rightarrow \text{Field} \rightarrow \text{Val}$ , where

Field is the set containing all field names.  $\text{fresh}(T, h) = (r, h')$  allocates an object of type  $T$  in heap  $h$  and returns (i) the reference  $r$  of the freshly allocated object, and (ii) the heap  $h'$  resulting from the allocation of memory on  $h$ .

For simplicity, we consider only integer (Int) and long (Long) numerical types ( $\text{NumTypes} = \{\text{Int}, \text{Long}\}$ ), and references (Ref). Given a value  $v$ ,  $\text{typeOf}(v)$  returns its type (Int, Long, or Ref). Since JB instructions often append a prefix to distinguish instructions dealing with different types (e.g.,  $\text{iadd}$  and  $\text{ladd}$ ), we define a support function  $\text{JVMprefix}$  that, given a type  $t$ , returns the prefix of the given type (i.e.,  $\text{i}$  if  $t = \text{Int}$ ,  $\text{l}$  if  $t = \text{Long}$ , and  $\text{a}$  if  $t = \text{Ref}$ ). We define by  $\text{WRef}$  an object type with a unique field value.

Given a method signature  $m$  and a list of arguments  $L$  (with the receiver in the first argument if  $m$  is not static),  $\text{body}(m, L) : \mathbb{N} \rightarrow \text{St}$  returns the body of the method resolving the call, that is, a sequence of statements (represented by a function mapping indexes to statements). Similarly, each statement belongs to a method; hence  $\text{getBody}(\text{st}) = b$  is the body of the method where  $\text{st}$  occurs. Finally,  $\text{isStatic}(m)$  means that  $m$  is static.

### 3.2 CIL

First, we define the concrete semantics of our CIL fragment (Sec. 2.1).

**Concrete State** A local state in CIL is composed by a stack of values or reference to local variables  $\text{Stack} : \mathbb{N} \rightarrow \text{Val} \cup \text{Ref}_{\text{Loc}}$  (where  $r_i \in \text{Ref}_{\text{Loc}}$  represents the cell's reference of the  $i$ -th local variable), an array of local variables  $\text{Loc} : \mathbb{N} \rightarrow \text{Val}$ , and an array of method arguments  $\text{Arg} : \mathbb{N} \rightarrow \text{Val}$ . A concrete CIL state consists of a local state and a heap, that is,  $\Sigma_{\text{CIL}} = \text{Stack} \times \text{Loc} \times \text{Arg} \times \text{Heap}$ .

**Concrete Semantics** Fig. 5 shows the concrete CIL semantics  $(\text{st}, \sigma) \rightarrow_{\text{CIL}} \sigma'$ . For a statement  $\text{st}$  and an entry state  $\sigma$ , it yields the state  $\sigma'$  resulting from the execution of  $\text{st}$  over  $\sigma$ ; or a program label  $l$ , meaning that the next instruction to execute is that at  $l$ . Otherwise, the next instruction to execute is implicitly assumed to be the subsequent one, sequentially (if any).

For the most part, the concrete semantics just formalizes the runtime semantics of the CIL ECMA Standard [16]. For instance, rule  $\text{add}$  pops the two topmost values of the operand stack and replaces them with their addition. However, its semantics is defined iff the two values have the same type. Instead,  $\text{ldloc } i$  pushes to the operand stack the value of the  $i$ -th local variables, while  $\text{stloc } i$  stores the top of the operand stack into the  $i$ -th local variable. Statements working with objects, such as  $\text{ldfld}$  and  $\text{stfld}$ , read from and write into the heap, if their receiver is not null.  $\text{call}$  and  $\text{callstatic}$  create a frame (i.e., an array of arguments, an empty array of local variables and an empty operand stack), execute the callee and leave its returned value on the stack, if any. For simplicity, the formalization assumes that there is no returned value. Finally,  $\text{ldloca } i$  loads the reference to the  $i$ -th local variable to the stack (represented by  $r_i$ ),  $\text{stind}$  stores the given value to the given reference, and  $\text{ldind}$  loads the value pointed by the given reference.

**Running example:** Consider the running example in Fig. 2b and apply the concrete semantics when  $n$  is 1, that is, when the entry state consists of an empty operand stack and of an array of local variables, while the value of the arguments is  $[0 \mapsto 1]$ . Assume that  $\text{newobj}$  allocates the object at address #1. Then after the first block (lines 3-7) the address of the object is stored in local variable 0, local variable 1 (representing variable  $i$  of the source program) holds 0, and address #1 in the heap holds an object of type  $\text{List}(\text{Collection})$ . Formally, the concrete state at line 7 is  $(\emptyset, [0 \mapsto \#1, 1 \mapsto 0], [0 \mapsto 1], [\#1 \mapsto \langle \rangle])$ , where  $\langle \rangle$  stands for the

empty list. Then the body of the  $\text{for}$  loop (lines 9-16) is executed once and creates a new  $\text{Wrapper}$  object (assume at address #2) wrapping 0, adds it to the list at address #1 and increments counter  $i$  (i.e., local variable 0) by 1. Hence the execution of the concrete semantics of the body of the loop leads to the concrete state  $\sigma = (\emptyset, [0 \mapsto \#1, 1 \mapsto 1], [0 \mapsto 1], [\#1 \mapsto \langle \#2 \rangle, \#2 \mapsto \{f \mapsto 0\}])$  at line 16. The condition at line 20 will then route the program to line 22 (since both argument 0 and local variable 1 hold 1). In conclusion, the concrete semantics will reach statement  $\text{ret}$  at line 23 with a state  $\sigma_{\text{CIL}}$  equal to  $\sigma$ , but where the operand stack contains reference #1.

### 3.3 JB

Let us turn to the concrete semantics of the JB fragment (Sec. 2.1).

**Concrete State** A local state in JB is composed by a stack of values  $\text{Stack} : \mathbb{N} \rightarrow \text{Val}$  and an array of local variables  $\text{Loc} : \mathbb{N} \rightarrow \text{Val}$ . A concrete JB state consists of a local state and a heap, that is,  $\Sigma_{\text{JB}} = \text{Stack} \times \text{Loc} \times \text{Heap}$ .

**Concrete Semantics** Fig. 6 reports the concrete JB semantics  $\rightarrow_{\text{JB}}$ . For a statement  $\text{st}$  and an entry state  $\sigma$ , it yields the state  $\sigma'$  resulting from the execution of  $\text{st}$  over  $\sigma$ .

For the most part, the behavior of this semantics is identical to that for CIL. The main differences are that JB instructions work on specific types (e.g., while CIL  $\text{add}$  statement adds two values of the same type, JB  $\text{iadd}$  and  $\text{ladd}$  statements add the values iff they are both  $\text{int}$  or  $\text{long}$ , respectively), and  $\text{new}$  only allocates a new object, while CIL  $\text{newobj}$  statements also calls a constructor.

**Running example:** The application of the JB concrete semantics is similar to that for CIL, but there are two minor differences: (i) there is only one array of local variables representing both CIL arguments and local variables (e.g., CIL local variable 0 is represented by JB local variable 1, since the first local variable holds the argument of the method); and (ii) there is an *instrumentation* local variable at index 3. Therefore, after we apply the JB concrete semantics from the entry state that maps the argument to 1, we obtain the concrete state  $([0 \mapsto \#1], [0 \mapsto 1, 1 \mapsto \#1, 2 \mapsto 1, 3 \mapsto 0], [\#1 \mapsto \langle \#2 \rangle, \#2 \mapsto \{f \mapsto 0\}])$ .

## 4 FROM CIL TO JB

This section formalizes and proves correct the translation of CIL statements into JB statements.

### 4.1 Concrete States

Function  $\mathbb{T}_\sigma \llbracket \cdot \rrbracket : \Sigma_{\text{CIL}} \rightarrow \Sigma_{\text{JB}}$  translates CIL concrete states into JB concrete states:  $\mathbb{T}_\sigma \llbracket (s, l, a, h) \rrbracket = (s', \text{cnvrtLoc}(l, a), h')$

This function (i) replaces direct reference with wrapper objects, and (ii) merges the array of local variables and arguments, adjusting variable indexes for 64 bits values. Formally:

$$i \in \text{dom}(s), l = \text{height}(s)$$

$$s' = \left[ i \mapsto \begin{cases} s(i) & \text{if } s(i) \in \text{Val} \\ r_i & \text{if } s(i) \in \text{Ref}_{\text{Loc}} \end{cases} \right]$$

$$\text{where } h'_{-1} = h, \text{ and } (h'_i, r_i) = \text{allocWrp}(h'_{i-1}, i)$$

$$\text{allocWrp}(h, j) = \begin{cases} (h, \text{null}) & \text{if } s(j) \in \text{Val} \\ (h'[r \mapsto h(r)[\text{value} \mapsto l(j)]]) & \text{if } s(i) \in \text{Ref}_{\text{Loc}} \\ \quad \text{where } (r, h') = \text{fresh}(\text{WRef}, h) \end{cases}$$

$$\begin{array}{c}
\frac{\text{typeOf}(v_1) = \text{typeOf}(v_2)}{\langle \text{add}, (s :: v_1 :: v_2, l, a, h) \rangle \rightarrow_{\text{CIL}} (s :: (v_1 + v_2), l, a, h)} \quad (\text{add}) \qquad \frac{}{\langle \text{ldloc } i, (s, l, a, h) \rangle \rightarrow_{\text{CIL}} (s :: l(i), l, a, h)} \quad (\text{ldloc}) \\
\frac{}{\langle \text{stloc } i, (s :: v, l, a, h) \rangle \rightarrow_{\text{CIL}} (s :: l[i \mapsto v], a, h)} \quad (\text{stloc}) \qquad \frac{}{\langle \text{ldarg } i, (s :: l, a, h) \rangle \rightarrow_{\text{CIL}} (s :: a(i), l, a, h)} \quad (\text{ldarg}) \\
\frac{\text{isStatic}(m(\text{arg}_0, \dots, \text{arg}_i)) = \text{false} \wedge t \neq \text{null} \wedge \langle \text{body}(m(\text{arg}_0, \dots, \text{arg}_i), (t, v_1, \dots, v_i)), (\varepsilon, \emptyset, [0 \mapsto t, j \mapsto v_j : j \in [1..i]], h) \rangle \rightarrow_{\text{CIL}} (s', l', a', h')}{\langle \text{call } m(\text{arg}_1, \dots, \text{arg}_i), (s :: t :: v_1 :: \dots :: v_i, l, a, h) \rangle \rightarrow_{\text{CIL}} (s, l, a, h')} \quad (\text{call}) \\
\frac{\text{isStatic}(m(\text{arg}_0, \dots, \text{arg}_i)) = \text{true} \wedge \langle \text{body}(m(\text{arg}_0, \dots, \text{arg}_i), (v_1, \dots, v_i)), (\varepsilon, \emptyset, [j - 1 \mapsto v_j : j \in [1..i]], h) \rangle \rightarrow_{\text{CIL}} (s', l', a', h')}{\langle \text{call } m(\text{arg}_1, \dots, \text{arg}_i), (s :: v_1 :: \dots :: v_i, l, a, h) \rangle \rightarrow_{\text{CIL}} (s, l, a, h')} \quad (\text{callstatic}) \\
\frac{\text{fresh}(T, h) = (r, h_1) \wedge \langle \text{body}(\text{ctor}(\text{arg}_1, \dots, \text{arg}_i), (v_1, \dots, v_i)), (\varepsilon, \emptyset, [0 \mapsto r, j \mapsto v_j : j \in [1..i]], h_1) \rangle \rightarrow_{\text{CIL}} (s', l', a', h')}{\langle \text{newobj } T(a_1, \dots, a_i), (s :: v_1 :: \dots :: v_i, l, a, h) \rangle \rightarrow_{\text{CIL}} (s :: r, l, a, h')} \quad (\text{newobj}) \\
\frac{o \neq \text{null}}{\langle \text{ldfld } f, (s :: o, l, a, h) \rangle \rightarrow_{\text{CIL}} (s :: h(o)(f), l, a, h)} \quad (\text{ldfld}) \qquad \frac{o \neq \text{null} \quad s' = h(o)[f \mapsto v]}{\langle \text{stfld } f, (s :: o :: v, l, a, h) \rangle \rightarrow_{\text{CIL}} (s, l, a, h[o \mapsto s'])} \quad (\text{stfld}) \\
\frac{\text{typeOf}(v_1) = \text{typeOf}(v_2) \wedge v_1 > v_2}{\langle \text{bgt } l, (s :: v_1 :: v_2, l, a, h) \rangle \rightarrow_{\text{CIL}} \langle l, (s, l, a, h) \rangle} \quad (\text{bgt true}) \qquad \frac{\text{typeOf}(v_1) = \text{typeOf}(v_2) \wedge v_1 \leq v_2}{\langle \text{bgt } l, (s :: v_1 :: v_2, l, a, h) \rangle \rightarrow_{\text{CIL}} (s, l, a, h)} \quad (\text{bgt false}) \\
\frac{}{\langle \text{ldloca } i, (s, l, a, h) \rangle \rightarrow_{\text{CIL}} (s :: r_i, l, a, h)} \quad (\text{ldloca}) \qquad \frac{}{\langle \text{stind}, (s :: r_i :: v, l, a, h) \rangle \rightarrow_{\text{CIL}} (s, l, a, h[r_i \mapsto v])} \quad (\text{stind}) \\
\frac{}{\langle \text{dup}, (s :: v, l, a, h) \rangle \rightarrow_{\text{CIL}} (s :: v :: v, l, a, h)} \quad (\text{dup}) \qquad \frac{}{\langle \text{ldind}, (s :: r_i, l, a, h) \rangle \rightarrow_{\text{CIL}} (s :: h(r_i), l, a, h)} \quad (\text{ldind})
\end{array}$$

Figure 5: Concrete CIL semantics.

$$\begin{array}{c}
\frac{\text{typeOf}(v) \neq \text{Long}}{\langle \text{dup}, (s :: v, l, h) \rangle \rightarrow_{\text{JB}} (s :: v :: v, l, h)} \quad (\text{dup}) \\
\frac{\text{typeOf}(v_1) \neq \text{Long}}{\langle \text{dup2}, (s :: v_1 :: v_2, l, h) \rangle \rightarrow_{\text{JB}} (s :: v_1 :: v_2 :: v_1 :: v_2, l, h)} \quad (\text{dup2 } 32) \qquad \frac{\text{typeOf}(v) = \text{Long}}{\langle \text{dup2}, (s :: v, l, h) \rangle \rightarrow_{\text{JB}} (s :: v :: v, l, h)} \quad (\text{dup2 } 64) \\
\frac{\text{typeOf}(v_1) = \text{Int} \wedge \text{typeOf}(v_2) = \text{Int}}{\langle \text{iadd}, (s :: v_1 :: v_2, l, h) \rangle \rightarrow_{\text{JB}} (s :: (v_1 + v_2), l, h)} \quad (\text{iadd}) \qquad \frac{\text{typeOf}(v_1) = \text{Long} \wedge \text{typeOf}(v_2) = \text{Long}}{\langle \text{ladd}, (s :: v_1 :: v_2, l, h) \rangle \rightarrow_{\text{JB}} (s :: (v_1 + v_2), l, h)} \quad (\text{ladd}) \\
\frac{x = \text{JVMprefix}(\text{typeOf}(l(i)))}{\langle \text{xload } i, (s, l, h) \rangle \rightarrow_{\text{JB}} (s :: l(i), l, h)} \quad (\text{xload}) \qquad \frac{x = \text{JVMprefix}(\text{typeOf}(v))}{\langle \text{xstore } i, (s :: v, l, h) \rangle \rightarrow_{\text{JB}} (s, l[i \mapsto v], h)} \quad (\text{xstore}) \\
\frac{\text{isStatic}(m(\text{arg}_0, \dots, \text{arg}_i)) = \text{false} \wedge t \neq \text{null} \wedge \langle \text{body}(m(\text{arg}_0, \dots, \text{arg}_i), (t, v_1, \dots, v_i)), ([], [0 \mapsto t, j \mapsto v_j : j \in [1..i]], h) \rangle \rightarrow_{\text{JB}} (s', l', h')}{\langle \text{invokevirtual } m(\text{arg}_1, \dots, \text{arg}_i), (s :: t :: v_1 :: \dots :: v_i, l, h) \rangle \rightarrow_{\text{JB}} (s, l, h')} \quad (\text{invokevirtual}) \\
\frac{\text{isStatic}(m(\text{arg}_0, \dots, \text{arg}_i)) = \text{true} \wedge \langle \text{body}(m(\text{arg}_0, \dots, \text{arg}_i), (v_1, \dots, v_i)), ([], [j - 1 \mapsto v_j : j \in [1..i]], h) \rangle \rightarrow_{\text{JB}} (s', l', h')}{\langle \text{invokestatic } m(\text{arg}_1, \dots, \text{arg}_i), (s :: v_1 :: \dots :: v_i, l, h) \rangle \rightarrow_{\text{JB}} (s, l, h')} \quad (\text{invokestatic}) \\
\frac{\text{fresh}(T, h) = (r, h')}{\langle \text{new } T, (s, l, h) \rangle \rightarrow_{\text{JB}} (s :: r, l, h')} \quad (\text{new}) \qquad \frac{o \neq \text{null}}{\langle \text{getfield } f, (s :: o, l, h) \rangle \rightarrow_{\text{JB}} (s :: h(o)(f), l, h)} \quad (\text{getfield}) \\
\frac{o \neq \text{null} \quad s' = h(o)[f \mapsto v]}{\langle \text{putfield } f, (s :: o :: v, l, h) \rangle \rightarrow_{\text{JB}} (s, l, h[o \mapsto s'])} \quad (\text{putfield}) \qquad \frac{\text{typeOf}(v_1) = \text{Int} \wedge \text{typeOf}(v_2) = \text{Int} \wedge v_1 > v_2}{\langle \text{if\_icmpgt } l, (s :: v_1 :: v_2, l, h) \rangle \rightarrow_{\text{JB}} \langle l, (s, l, h) \rangle} \quad (\text{if\_icmpgt true}) \\
\frac{\text{typeOf}(v_1) = \text{Int} \wedge \text{typeOf}(v_2) = \text{Int} \wedge v_1 \leq v_2}{\langle \text{if\_icmpgt } l, (s :: v_1 :: v_2, l, h) \rangle \rightarrow_{\text{JB}} (s, l, h)} \quad (\text{if\_icmpgt false})
\end{array}$$

Figure 6: Concrete JB semantics.

Intuitively, each direct reference in the operand stack (that is,  $s(i) \in \text{Ref}_{\text{Loc}}$ ) is replaced by another reference pointing to a wrapper object freshly allocated and containing in its field the value pointed by the original direct reference. 6

Then, for an array of values  $b$  and an index  $i$ , the following function counts the 64 bits types among the first  $i$ :

$$64_b^i = |\{j : 0 \leq j < i \text{ and } b[j] \text{ is a 64 bit value}\}|$$

Then the array  $cnvrtLoc(l, a)$  is defined as follows:

$$\forall 0 \leq i < |a| : cnvrtLoc(l, a)[i + 64_a^i] = a[i]$$

$$\forall 0 \leq i < || : cnvrtLoc(l, a)[|a| + 64_a^{||} + i + 64_a^i] = l[i]$$

**Running example:** Consider the CIL exit state computed in Sec. 3.2, that is,  $\sigma_{CIL} = ([0 \mapsto \#1], [0 \mapsto \#1, 1 \mapsto 1], [0 \mapsto 1], [\#1 \mapsto < \#2 >, \#2 \mapsto \{f \mapsto 0\}])$ . The CIL to JB translation computes  $\mathbb{T}_\sigma[\sigma_{CIL}] = ([0 \mapsto \#1], [0 \mapsto 1, 1 \mapsto \#1, 2 \mapsto 1], [\#1 \mapsto < \#2 >, \#2 \mapsto \{f \mapsto 0\}])$  (it merges CIL arguments and local variables). This state is almost identical to the exit state of the JB concrete semantics applied to the running example (Sec. 3.3) but for the instrumentation local variable at index 3.

## 4.2 Statements

Fig. 7 formalizes the translation  $\mathbb{T}[\text{st}_{CIL}, K] = \text{st}_{JB}$  of a single CIL statement into a sequence of JB statements.  $K$  is the static type information about locals, arguments and stack elements, computed at  $\text{st}_{CIL}$  by a standard algorithm [16]. In particular, types and height of the stack are fixed and statistically known at each bytecode. In addition, the forth component  $\bar{w}$  is a stack of elements in  $\perp \cup (\mathbb{N} \times \mathbb{N})$  that, for each element in the operand stack, tells (i)  $\perp$  if it is not a direct reference, or (ii)  $(i, j)$  where  $i$  is the index of the local variables pointed by the direct reference<sup>4</sup>, and  $j$  the index of the local instrumentation variable containing a pointer to the wrapper simulating the direct reference.

Few CIL statements (namely,  $ldfld$  and  $stfld$ ) have a one-to-one translation into a JB statement ( $getField$  and  $putfield$ ). The statements reading and writing local variables and arguments ( $ldarg$ ,  $ldloc$ , and  $stloc$ ) are translated into their JB counterpart ( $xload$ ,  $xstore$ , respectively) taking into account the type of the value at the top of the stack, and adjusting the index of the variable taking into account arguments and 64 bit variables. Some CIL statements ( $dup$ ) get translated into different JB statements on the basis of contextual information such as the type of values in the operand stack ( $dup$  and  $dup2$ ). Other CIL statements can be translated only if the type of the values in the operand stack is numeric: (i)  $add$  can be translated into  $ladd$  and  $iadd$ , and (ii)  $bgt$  to  $is.icmpgt$ ; if they are applied to references (as in generic CIL code), then the code is unsafe and we do not support its translation.

$call$  requires to (i) translate the method call to the corresponding static or dynamic invocation statement in JB, and (ii) to propagate the side effects on direct pointers passed to the method as  $out/ref$  parameters to the local variables of the callee.

The translation of  $newobj$  is tricky because of the different patterns used in CIL and JB for object creation<sup>5</sup>. While CIL creates and initializes the object (*i.e.*, calls its constructor) with a single instruction, JB splits these operations and requires the newly created object to occur below the arguments on the stack, before calling the constructor. Hence, the translation relies on a function  $freshIdx$  to store and load the values of the constructor arguments through instrumentation local variables. In particular, given a CIL method  $m$ , the number and types of arguments and local variables of the

<sup>4</sup>Since the language we introduced in Fig. 1 supports only  $ldloca$  to get a direct pointer, we need to track only this information in the formalization.

<sup>5</sup>For sake of simplicity, we assume the constructor does not have  $out/ref$  parameters. In the implementation, they are treated as for  $call$  statements

method are known (Sec. II.15.4 of [16]). Therefore, function  $cnvrtLoc$  tells which local variables the translated JB method already uses. Then, for each argument of each  $newobj$  statement in  $m$ , it is possible to allocate a fresh local variable to store and load its value. In this way, the translation allocates a new object and puts its address below the constructor arguments.

Instructions dealing with direct pointers (namely,  $ldloca$ ,  $stind$ , and  $ldind$  in our minimal language) are translated through equivalent CIL instructions dealing with wrapper objects (and their field value). Therefore,  $stind$  and  $ldind$  are simply translated through equivalent write and read of field value, respectively.  $ldloca$  instead requires to allocate a wrapper object  $newobj$ , stores a reference to the wrapper ( $stloc$ ) in an instrumentation variable obtained through  $freshIdx$ , stores the value pointed by the direct reference in the local variable to its field value ( $ldloc$  and  $stfld$ ), and leaves a reference to the wrapper object in the operand stack ( $dup$ ).

In general, each CIL statement is translated into one or more JB statements, hence offsets are not preserved. Thus, function  $statementIdx : St \rightarrow \mathbb{N}$  yields the JB offset of the first statement in the translation of the given CIL statement. In addition, since direct references are replaced by wrapper objects, when a method parameter has a direct reference type  $\&T$  (and this happens when it is a  $ref$  or  $out$  parameter in safe C#), this is replaced by a wrapper object  $WRef$ .

**Running example:** Consider the running example in Fig. 2. Most CIL statements are translated into a single JB statement (*e.g.*, lines 18-20 and 22-23 of Fig. 2b are translated into lines 23-25 and 27-28 of Fig. 2c), with the noticeable exception of the CIL  $newobj$  statements at line 3 and 11, translated into lines 2-4 and 12-16, respectively. The former passes no argument to the constructor; the latter (that instantiates a  $Wrapper$ ) calls a constructor with an argument, hence requiring an instrumentation variable at index 3. **Direct references** As sketched in Sec. 2.2.1, we model the semantics of pointers in safe C# code through wrapper objects. In particular,  $ldind$  (line 5 of Fig. 3b) is translated into the field access  $i.value$  (right side of the assignment at line 2 of Fig. 3c), while  $stind$  (line 8 of Fig. 3b) is translated into the assignment of  $i.value$  (left side of the assignment at line 2 of Fig. 3c). In addition,  $ldloca$  (line 15 of Figure 3b) leads to the construction and assignment of a wrapper object (lines 7 and 8 of Figure 3c), while after the method call the value contained in the wrapper object is written into the local variable (line 10 of Figure 3c).

## 4.3 Correctness

We prove the translation from CIL to JB being correct. Namely, given a concrete CIL state  $\sigma_{CIL}$ , by applying the operational semantics for a statement  $st$ , we get a state that, when translated into JB, is exactly the one resulting from the translation of  $\sigma_{CIL}$  into JB and the application of the JB semantics to it:

$$\begin{aligned} \forall st \in St_{CIL}, \sigma_{CIL} \in \Sigma_{CIL} : \\ \langle st, \sigma_{CIL} \rangle \rightarrow_{CIL} \sigma'_{CIL} \text{ and } \langle \mathbb{T}[st, K], \mathbb{T}_\sigma[\sigma_{CIL}] \rangle \rightarrow_{JB} \sigma'_{JB} \\ \Downarrow \\ \mathbb{T}_\sigma[\sigma'_{CIL}] = \bullet \sigma'_{JB} \end{aligned}$$

where  $\sigma_1 = \bullet \sigma_2$  means that the two states are equal up to instrumentation variables introduced by the translation process. Formally, let

$$\begin{aligned}
\mathbb{T}[\text{dup}, \bar{s} :: t, \bar{l}, \bar{a}, \bar{w}] &= \begin{cases} \text{dup} & \text{if } t \neq \text{Long} \\ \text{dup2} & \text{if } t = \text{Long} \end{cases} \\
\mathbb{T}[\text{add}, \bar{s} :: t_1 :: t_2, \bar{l}, \bar{a}, \bar{w}] &= \begin{cases} \text{iadd} & \text{if } t_1 = t_2 = \text{Int} \\ \text{ladd} & \text{if } t_1 = t_2 = \text{Long} \end{cases} \\
\mathbb{T}[\text{ldloc } i, \bar{s}, \bar{l}, \bar{a}, \bar{w}] &= \text{xload } j \text{ where } j = |a| + 64_a^{|a|} + i + 64_j^i \wedge x = \text{JVMprefix}(\text{typeOf}(\bar{l}(i))) \\
\mathbb{T}[\text{stloc } i, \bar{s} :: t, \bar{l}, \bar{a}, \bar{w}] &= \text{xstore } j \text{ where } j = |a| + 64_a^{|a|} + i + 64_j^i \wedge x = \text{JVMprefix}(\text{typeOf}(\bar{l}(i))) \\
\mathbb{T}[\text{ldarg } i, \bar{s}, \bar{l}, \bar{a}, \bar{w}] &= \text{xload } j \text{ where } j = i + 64_a^i \wedge x = \text{JVMprefix}(\text{typeOf}(\bar{a}(i))) \\
\mathbb{T}[\text{call } m(\text{arg}_1, \dots, \text{arg}_i), \\ \bar{s} :: t_1 :: \dots :: t_i, \bar{l}, \bar{a}, \bar{w} :: p_1 :: \dots :: p_i] &= \text{invoke}; \text{aload } p_{idx_1}^1; \text{getfield value}; x_{idx_1} \text{store } p_{idx_1}^2; \dots \\ &\quad \dots \text{aload } p_{idx_j}^1; \text{getfield value}; x_{idx_j} \text{store } p_{idx_j}^2; \\ &\quad \text{where invoke} = \begin{cases} \text{invokestatic } m(\text{arg}_1, \dots, \text{arg}_i) & \text{if } \text{isStatic}(m(\text{arg}_1, \dots, \text{arg}_i)) \\ \text{invokevirtual } m(\text{arg}_1, \dots, \text{arg}_i) & \text{otherwise} \end{cases} \\ &\quad \{idx_1, \dots, idx_j\} = \{k : \text{arg}_k \in \text{Ref}_{\text{Loc}}\} \\ &\quad \forall k \in [1..j] : x_{idx_k} = \text{JVMprefix}(\text{typeOf}(\bar{l}(p_{idx_k}^2))) \wedge \forall r \in [1..i] : p_i = (p_i^1, p_i^2) \\
\mathbb{T}[\text{newobj } \mathbb{T}(a_1, \dots, a_i), \\ \bar{s} :: t_1 :: \dots :: t_i, \bar{l}, \bar{a}, \bar{w}] &= x_i \text{store } idx_i; \dots; x_1 \text{store } idx_1; \text{new } \mathbb{T}; \text{dup}; \\ &\quad x_1 \text{load } idx_1; \dots; x_i \text{load } idx_i; \text{invokevirtual } \langle \text{init} \rangle (\text{arg}_1, \dots, \text{arg}_i) \\ &\quad \text{where } \forall j \in [1..i] : x_j = \text{JVMprefix}(a_j) \wedge idx_j = \text{freshIdx}(\text{newobj } \mathbb{T}(a_1, \dots, a_i), j) \\
\mathbb{T}[\text{ldfld } f, \bar{s} :: t_0, \bar{l}, \bar{a}, \bar{w}] &= \text{getfield } f \\
\mathbb{T}[\text{stfld } f, \bar{s} :: t_0 :: t_v, \bar{l}, \bar{a}, \bar{w}] &= \text{putfield } f \\
\mathbb{T}[\text{bgt } k, \bar{s} :: t_1 :: t_2, \bar{l}, \bar{a}, \bar{w}] &= \text{if\_icmpgt } k' \text{ where } k' = \text{statementIdx}(\text{getBody}(\text{bgt } k)(k)) \text{ if } t_1 = t_2 = \text{Int} \\
\mathbb{T}[\text{ldloca } i, \bar{s}, \bar{l}, \bar{a}, \bar{w}] &= \mathbb{T}[\text{newobj } \text{WrapRef}(); \text{dup2}; \text{stloc } j; \text{ldloc } i; \text{stfld value}, \bar{s}, \bar{l}, \bar{a}, \bar{w}] \\ &\quad \text{where } j = \text{freshIdx}(\text{ldloca } i, 0) \\
\mathbb{T}[\text{stind}, \bar{s}, \bar{l}, \bar{a}, \bar{w}] &= \mathbb{T}[\text{stfld value}, \bar{s}, \bar{l}, \bar{a}, \bar{w}] \\
\mathbb{T}[\text{ldind}, \bar{s}, \bar{l}, \bar{a}, \bar{w}] &= \mathbb{T}[\text{ldfld value}, \bar{s}, \bar{l}, \bar{a}, \bar{w}]
\end{aligned}$$

Figure 7: Translation of CIL statements into JB.

$\sigma_1 = (s_1, l_1^1 :: \dots :: l_1^n, h_1)$  and  $\sigma_2 = (s_2, l_2^1 :: \dots :: l_2^n, :: l_2^{n+1} :: \dots :: l_2^{n+k}, h_2)$ , then  $\sigma_1 =_{\bullet} \sigma_2$  iff  $s_1 = s_2$ , and  $\forall i \leq n : l_1^i = l_2^i$ , and  $h_1 = h_2$ . Note that instrumentation variables are present only in the JB state, hence in the right hand-side of the equality.<sup>6</sup>

**Running example:** Sec. 3.2 showed that, starting from  $\sigma_{\text{CIL}} = (\emptyset, \emptyset, [0 \mapsto 1], \emptyset)$ , the concrete semantics on the program in Fig. 2b ends up in  $\sigma'_{\text{CIL}} = ([0 \mapsto \#1], [0 \mapsto \#1, 1 \mapsto 1], [0 \mapsto 1], [\#1 \mapsto < \#2 >, \#2 \mapsto \{f \mapsto 0\}])$ . Then Sec. 3.3 showed that, starting from the corresponding  $\mathbb{T}_{\sigma}[\sigma_{\text{CIL}}]$  state, the JB concrete semantics leads to  $\sigma'_{\text{JB}} = ([0 \mapsto \#1], [0 \mapsto 1, 1 \mapsto \#1, 2 \mapsto 1, 3 \mapsto 0], [\#1 \mapsto < \#2 >, \#2 \mapsto \{f \mapsto 0\}])$ . So, by definition of  $\mathbb{T}_{\sigma}[\ ]$ , we get  $\mathbb{T}_{\sigma}[\sigma'_{\text{CIL}}] = \bullet \sigma'_{\text{JB}}$  since the two stacks and the two heaps are equal, the values of three local variables of the JB state correspond to the values of the argument and the two local variables of the CIL state, respectively, and  $=_{\bullet}$  projects out the fourth variable of the JB state  $\sigma'$ .

#### 4.4 Other Instructions

In this section, we informally discuss how our approach deals with CIL instructions that are slightly different from other instructions in

JB. We decided to handle these instructions informally since their translation is mostly straightforward. It is intended for readers that are expert of JB, CIL and more advanced C# features, such as generic type erasure in JB, or delegates in C#.

**Numerical and Reference Comparison** CIL compares numerical or reference values in two ways: through conditional branches (e.g., `beq` branches when the topmost two values on the stack are equals) and comparisons (e.g., `ceq` pushes 1 iff the topmost two values on the stack are equals, and 0 otherwise). As usual, these instructions are type independent and apply to numerical (int, float, long, ...) as well as reference values. JB uses a different approach, since its instructions are type dependent. If the topmost two values on the stack are integers, it uses a conditional branch instruction (`if_icmpeq`) similar to that of CIL (`beq`). However, JB has no comparison instruction on integers and we need to simulate it through a sequence of JB instructions relying on constants and branch. For instance, a `ceq` statement on integers is simulated as in Fig. 8a. Instead, if the topmost two values on the stack are long, JB uses a comparison statement `lcmp` that pushes to the stack 1, 0, or -1 iff the first value is less than, equal to, or greater than the

<sup>6</sup>Appendix A reports the detailed formal proof.



<pre> 1  if_icmpeq 4 2  iconst_0 3  goto 5 4  iconst_1 5  nop </pre> <p style="text-align: center;">(a) On integers</p>	<pre> 1  lcmp 2  iconst_0 3  if_icmpeq 6 4  iconst_0 5  goto 7 6  iconst_1 7  nop </pre> <p style="text-align: center;">(b) On longs</p>
---	--

**Figure 8: Translation of `ceq`.**

<pre> 1  ldc.i4.0 2  ldc.i4.0 3  callvirt !0 List&lt;A&gt;::get(int32) 4  stloc.1 </pre> <p style="text-align: center;">(a) In CIL</p>	<pre> 1  aload.0 2  iconst.0 3  invoke List.get:(I)Object; 4  checkcast A 5  astore.1 </pre> <p style="text-align: center;">(b) In JB</p>
--	---

**Figure 9: Getting an element from a list.**

second, respectively. Hence, we simulate CIL conditional branch and comparison instructions through `lcmp`, integer constants and a conditional branch on integers. For instance, `beq` is translated into the sequence `lcmp; ifne #i;` where `i` is the target JB instruction of `beq`. The comparisons over `long` is similar to `int`. Namely, `ceq` is translated into the code in Fig. 8b. Conditional branch and comparison work also on references. Equality and inequality statements are treated as for integers, since JB defines an `if_acmpeq` statement. Other CIL operators (e.g., `bgt`) can be applied to arbitrary references, as long as one of them is `null`.

**Generic Types** CIL keeps information about generic types, while JB erases it into `Object`. For instance, imagine that we have a local variable `list` of type `List<A>`. At source code level, a method call like `A a = list.get(0)` in Java or `A a = list[0]` in C# is legal since the elements of the list have type `A` in both languages. At bytecode level, getting an element from the list effectively returns an object of type `A` in CIL (see Fig. 9a), while it returns an object whose static type is `Object` in JB and casts it dynamically to `A` through a `checkcast` (Fig. 9b). Hence, our translation of a CIL method call with generic return type `T` adds a `checkcast` instruction to `T` after the call. Primitive types (e.g., `int` and `long`) can be passed as generic types in CIL but not in JB. Hence, when using a primitive type for the generic parameter or return value of a CIL method call, we box and unbox the primitive value into a Java wrapper class such as `java.lang.Integer`.

**Delegates** Lambda expressions have only been introduced in Java 8, while C# has been using *delegates* since its very beginning. C# implements delegates through CIL instructions that load a pointer to a method (`ldftn`) and execute it, sometime by using inner classes. Namely, C# accesses a pointer to the method through `ldftn` and calls the `Invoke` method of the delegate class. Consider for instance Fig. 10. The C# code in Fig. 10a uses a delegate to call a method. In Fig. 10b, this is compiled into a `ldftn` statement at line 4 followed by a call to `Invoke` at line 7. We translate this by using reflection and string constants. Namely, the signature of the method pointed by `ldftn` is represented by a string, passed to an instrumentation library call in class `Reflection`, that calls this method by reflection (Fig. 10c). However, many static analyzers (including Julia) are unsound for reflection. Hence, our translation marks all signatures accessed in this way as entry points (that is,

Library	# met.	# fail	% fail	Tr. t.	Mem.
mscorlib	28,344	870	3.07%	23"	158
Sys.Core	6,988	47	0.68%	4"	96
Sys.Design	13,509	4	0.03%	20"	180
Sys	17,851	242	1.36%	21"	142
Sys.Runtime.Serial	5,624	74	1.32%	5"	86
Sys.ServiceModel	34,603	80	0.23%	34"	156
Sys.Web	28,249	38	0.13%	37"	216
Sys.Web.Extensions	4,245	0	0.00%	4"	109
Sys.Windows.Forms	28,319	53	0.19%	42"	189
Sys.XML	12,727	171	1.34%	23"	146
Total	180,460	1,579	0.87%	3'33"	

**Table 2: Experimental results on libraries.**

methods that might be directly called from outside the application and therefore are analyzed under the most generic assumptions). This might cause a loss of precision, since contextual information on delegates is lost, but preserves soundness.

**Async and Await** In C#, an `async` method returns a `Task` object that allows the caller to execute the code of the method asynchronously. On the other hand, statement `await` waits until the execution of the asynchronous method ends and extracts the results of the computation. This pattern is compiled into method pointers and reflection at CIL bytecode level, in the same way delegates are treated. Therefore, we apply the same solution for delegates we described in Sec. 4.4.

## 5 EXPERIMENTAL RESULTS

We implemented our translation from CIL to JB through (i) a C# program that translates a CIL program to an intermediate XML representation (representing Java bytecode), and (ii) a Java program that produces a `jar` file from an XML representation. We had to split the implementation in this way since the library to read CIL bytecode (`Mono.Cecil`) is written in .NET, while the library writing `jar` bytecode (`BCEL`) is written in Java. The first part of the translation (CIL to JB) runs in parallel on different classes through the `System.Threading.Tasks` library (part of the standard .Net framework). We used an Intel Core i5-6600 CPU at 3.30GHz machine with 16 GB of RAM, 64-bit Windows 7 Professional, and Java SE Runtime Environment v.1.8.0\_111-b14.

As a first experiment to assess the efficiency and precision of our approach, we translated and analyzed the five most popular GitHub repositories (as on February 27th, 2017) written in C# and tagged as C# repositories<sup>7</sup>. Tab. 1 reports (i) the number of C# LOC of each projects (Column **LOC**)<sup>8</sup> (our benchmarks range between 17 and 120 KLOC, hence they are real world applications); (ii) the number of stars of the GitHub repository as on February 27th, 2017 (**GH \***); (iii) the total number of methods (**# meth.**), and for how many of them the translation failed because of unsafe code (**# fail**); (iv) the time (**Tr. t.**) and memory (**Mem.**, in MB) consumed by the translation from CIL to JB; (v) the time of Julia's analyses (**Analysis t.**), and (vi) the number of alarms (**Al.**), of false alarms because of loss of information introduced by the translation (**False al.**), and the precision (ratio of false alarms *w.r.t.* the total number of alarms, column **Precision**) of Julia's analysis.

<sup>7</sup>We consider the number of watchers as measure of popularity of a repository. We discarded some projects *tagged* as C# that actually mostly contain native code (corefx, coreclr, mono), that did not compile in Visual Studio (roslyn, powershell), that have been dismissed (shadowsocks), or that are particularly small (wavefunction, below 1KLOC).

<sup>8</sup>LOC are computed with `LocMetrics` <http://www.locmetrics.com/>

<pre> 1 delegate void Del(string message); 2 void DelegateMethod(string message) {...} 3 void go() { 4     Del handler = DelegateMethod; 5     handler("Hello World"); 6 } </pre> <p style="text-align: center;">(a) C# code</p>	<pre> 1 void go () { 2     ldarg.0 3     ldftn A::DelegateMethod(string) 4     newobj A/Del::ctor(object, int) 5     ldstr "Hello World" 6     call void A/Del::Invoke(string) 7 } </pre> <p style="text-align: center;">(b) CIL</p>	<pre> 1 void go() { 2     ldc "DelegateMethod(LString;)V" 3     invokestatic 4     Reflection.GetMethod:(LString;)LMethod; 5     ldc "Hello World" 6     invokevirtual A/Del.Invoke:(LString;) </pre> <p style="text-align: center;">(c) JB</p>
--	--	---

**Figure 10: An example of CIL delegate.**

Program	LOC	GH *	# met.	# fail	Tr. t.	Mem.	Analysis t.	Al.	False al.	Precision
CodeHub	32,510	7,718	4,887	0	0'07"	115	0'43"	9	1	89%
SignalR	71,207	6,285	6,610	3	0'07"	131	0'50"	8	1	88%
Dapper	22,513	5,815	1,058	0	0'07"	77	0'29"	13	3	77%
ShareX	171,580	5,208	11,568	14	0'58"	193	2'08"	57	0	100%
Nancy	109,139	4,969	8,817	0	0'07"	136	1'25"	18	1	94%
Total	406,949		32,940	17	1'26"		4'35"	105	6	94%

**Table 1: Experimental results on the 5 most starred Github C# projects.**

In order to assess the efficiency and library coverage of our approach, we also analyzed the 10 largest (based on the size of the .dll files) system libraries of the Microsoft .Net framework version 4.0.30319. They contain unsafe code (such as cryptographic code in mscorlib.dll) and might not be compiled from C#, but possibly from VB.Net. Tab. 2 reports the number of methods of the library (# **met.**), the number and percentage of methods where the translation fails because of unsafe code (# **fail** and % **fail**), and time (**Tr. t.**) and memory (**Mem.**, in MB) for the translation.

#### Research Question 1: Efficiency.

In 4 out of 5 top Github projects, our translation took 7" (Tab. 1); it took almost 1 minute for ShareX. These times are much shorter (overall, less than a third) than the analysis time. The memory consumed by the translation is small (below 200MB). The results for .Net framework libraries (Tab. 2) show a similar trend: we translated about 180K methods in about 3'30" (that is, a bit more than 1 msec per method) consuming at most 189MB of memory. So, our system deals with industrial-size software with a translation time comparable to the analysis time.

#### Research Question 2: Precision

We manually checked only the 105 high severity alarms issued by Julia on the top 5 Github projects, over a total of several thousands. Tab. 1 reports their number (column **Al.**) and the number of false alarms (**False al.**) due to our translation. The static analysis might generate false alarms as well, for instance because of disjunctive constraints not tracked by Julia; we do not count these as false alarms, since we want to evaluate the imprecision due to the translation, and not that inherent to Julia. In particular, 6 alarms out of 105 (about 6%) are false because of imprecision introduced by the translation. This shows that our approach satisfies Research Question 2. The origins of this imprecision are (i) `async` and `await` statements (in particular in `Dapper`), and (ii) `try-catch-finally` blocks (e.g., in `SignalR`). These would require to modify Julia to recognize these features more precisely (through automatic annotations produced by the translation).

#### Research Question 3: Libraries

We manually checked that all methods of the 5 top Github C# projects where our translation fails are actually unsafe. Column #

**fail** in Tab. 1 shows that there is no failure for CodeHub, Dapper and Nancy. Instead, there are 3 failures for SignalR, due to unsafe methods in class `Infrastructure.SipHashBasedStringEqualityComparer`, and 14 failures for ShareX, due to unsafe methods in two classes: (i) `GreenshotPlugin.Core` has methods setting or getting colors in fast implementations of bitmaps (`UnsafeBitmap` and subclasses); (ii) `ShareX.ImageHelpers` uses unsafe classes (such as `UnsafeBitmap`). This shows that our approach fails only for unsafe code with unsafe pointer manipulation (storing pointers in fields, returning them from methods, performing pointer arithmetic). Tab. 2 shows that the translation succeeds for 99.13% of the methods, with a worst case of 96.93%. Hence, our approach satisfies Research Question 3.

## 6 CONCLUSION

This article introduced, formalized and proved correct a translation from CIL to JB, for static analysis. To assess its feasibility and interest, it has been implemented and connected to the Julia analyzer. Experiments show positive results for efficiency, precision, and libraries' coverage. As future work, we plan to (i) improve Julia precision in the corner cases highlighted by our experiments, (ii) investigate new .Net properties of interest, and (iii) rely on `invokedynamic` when translating delegates.

## REFERENCES

- [1] Checkstyle. <http://checkstyle.sourceforge.net>.
- [2] Clr to jvm. <http://xmlvm.org/clr2jvm/>.
- [3] CodeSonar – Static Analysis SAST Software. <https://www.grammatech.com/products/codesonar>.
- [4] Findbugs<sup>TM</sup> – Find Bugs in Java Programs. <http://findbugs.sourceforge.net/>.
- [5] Ndepend. <http://www.ndepend.com>.
- [6] .Net Reflector Add-Ins. <https://www.microsoft.com/en-us/research/project/spec>.
- [7] PMD. <https://pmd.github.io/>.
- [8] WALA. [http://wala.sourceforge.net/wiki/index.php/Main\\_Page](http://wala.sourceforge.net/wiki/index.php/Main_Page).
- [9] M. Ameri and C. A. Furia. Why Just Boogie? Translating between Intermediate Verification Languages. In *Proceedings of IFM '16*, LNCS. Springer, 2016.
- [10] R. Atkey and D. Sannella. ThreadSafe: Static Analysis for Java Concurrency. *Electronic Comm. of the European Association of Software Science and Technology*, 72, 2015.
- [11] M. Barnett, K. Leino, and W. Schulte. The Spec# Programming System: An Overview. In *Proceedings of CASSIS '04*, 2004.
- [12] M. Bebenita, F. Brandner, M. Fähndrich, F. Logozzo, W. Schulte, N. Tillmann, and H. Venter. SPUR: a Trace-based JIT Compiler for CIL. In W. R. Cook, S. Clarke, and M. C. Rinard, editors, *Proceedings of OOPSLA '10*. ACM, 2010.

- [13] P. Cousot and R. Cousot. Abstract interpretation: a unified lattice model for static analysis of programs by construction or approximation of fixpoints. In *Proceedings of POPL '77*. ACM Press, 1977.
- [14] P. Cousot and R. Cousot. Abstract Interpretation: A Unified Lattice Model for Static Analysis of Programs by Construction or Approximation of Fixpoints. In *Proc. of Principles of Programming Languages (POPL'77)*, pages 238–252, 1977.
- [15] Coverity. Coverity Prevent<sup>TM</sup>. <http://www.coverity.com/library/pdf/coverity-prevent.pdf>.
- [16] ECMA. *Standard ECMA-335: Common Language Infrastructure (CLI)*. 2012.
- [17] T. A. S. Foundation. Apache Commons BCEL. <https://commons.apache.org/proper/commons-bcel>, checked on June 24, 2016.
- [18] JetBrains. Resharper. <https://www.jetbrains.com/resharper>.
- [19] T. Lindholm, F. Yellin, G. Bracha, and A. Buckley. *The Java Virtual Machine Specification, Java SE 7 Edition*. Addison-Wesley Professional, 1st edition, 2013.
- [20] F. Logozzo. Cibai: An Abstract Interpretation-Based Static Analyzer for Modular Analysis and Verification of Java Classes. In *Proceedings of VMCAI '07*, LNCS. Springer, 2007.
- [21] F. Logozzo and M. Fähndrich. Static Contract Checking with Abstract Interpretation. In *Proceedings of FoVeOOS '10*, LNCS. Springer, 2010.
- [22] Microsoft. FxCop. [https://msdn.microsoft.com/en-us/library/bb429476\(v=vs.80\).aspx](https://msdn.microsoft.com/en-us/library/bb429476(v=vs.80).aspx).
- [23] F. Spoto. The Julia Static Analyzer for Java. In *Proceedings of SAS '16*, LNCS. Springer, 2016.
- [24] Wikipedia. List of Tools for Static Code Analysis. [https://en.wikipedia.org/wiki/List\\_of\\_tools\\_for\\_static\\_code\\_analysis#Java](https://en.wikipedia.org/wiki/List_of_tools_for_static_code_analysis#Java).

## A PROOF OF CORRECTNESS

We present here the main steps of the correctness proof of our approach.

LEMMA A.1. *The function  $\text{convertLocals}$  is an identity embedding.*

PROOF. It is sufficient to observe that, by construction, the function concatenates the two arrays by shifting indexes when a 64 bits value occurs, hence preserving the values and the ordering of the elements' indexes.  $\square$

LEMMA A.2. *The translation of  $\text{ldloc}$  is correct.*

PROOF. Let us consider the case of integer arguments (the other case can be treated analogously). Let us prove that the translation of  $\text{ldloc}$  for integer values is correct, i.e., that  $\forall \sigma_{\text{CIL}} \in \Sigma_{\text{CIL}}$ , if  $\langle \text{ldloc } i, \sigma_{\text{CIL}} \rangle \rightarrow_{\text{CIL}} \sigma'_{\text{CIL}}$  and  $\langle \mathbb{T}[\text{ldloc } i, (\bar{s}, \bar{l}, \bar{a}, \bar{w})], \mathbb{T}[\sigma_{\text{CIL}}] \rangle \rightarrow_{\text{JB}} \sigma'_{\text{JB}}$  then  $\mathbb{T}[\sigma'_{\text{CIL}}] = \sigma'_{\text{JB}}$ . Let  $\sigma_{\text{CIL}} = (s, l, a, h)$  be arbitrary. By the  $\text{ldloc-CIL}$  rule we have  $\sigma'_{\text{CIL}} = (s :: 1[i], l, a, h)$ . By rule (3) of Fig. 7,  $\mathbb{T}[\text{ldloc } i, \bar{s}, \bar{l}, \bar{a}, \bar{w}] = \text{i load } j$  where  $j = |\bar{a}| + 64 \frac{|\bar{a}|}{a} + i + 64 \frac{i}{a}$ . By definition of  $\mathbb{T}[\sigma]$ :

$$\begin{aligned} \mathbb{T}[\sigma_{\text{CIL}}] &= \mathbb{T}[(s, l, a, h)] \\ &= (s', \text{convertLocals}(l, a), h'). \end{aligned}$$

By the  $\text{i load-JB}$  rule:

$$\begin{aligned} \langle \text{i load } j, (s, \text{convertLocals}(l, a), h) \rangle &\rightarrow_{\text{JB}} \\ (s' :: \text{convertLocals}(l, a)[j], \text{convertLocals}(l, a), h') &= \sigma'_{\text{JB}}. \end{aligned}$$

By definition of  $\mathbb{T}[\sigma]$  we have

$$\begin{aligned} \mathbb{T}[\sigma'_{\text{CIL}}] &= \mathbb{T}[(s :: 1[i], l, a, h)] \\ &= (s :: 1[i], \text{convertLocals}(l, a), h). \end{aligned}$$

By definition of  $\text{convertLocals}$  we have

$$\text{convertLocals}(l, a)[j] = 1[i],$$

which implies  $\mathbb{T}[\sigma'_{\text{CIL}}] = \sigma'_{\text{JB}}$ , and thus  $\mathbb{T}[\sigma'_{\text{CIL}}] = \bullet \sigma'_{\text{JB}}$ .  $\square$

LEMMA A.3. *The translation of  $\text{call}$  is correct.*

PROOF. Let us consider the case of static call (the other case can be treated analogously). Moreover, let us denote by  $\ell_{[i,j]}$  the sequence  $\ell_i, \dots, \ell_j$ , and by  $s_{(i,j)}$  the sequence of types  $[k_1, \dots, k_i, t_1, \dots, t_j]$ .

We show that  $\forall \sigma_{\text{CIL}} = ([u_{[1,n]}, v_{[1,i]}], l, a, h)$ , if  $\langle \text{call } m(\text{arg}_{[1,i]}), \sigma_{\text{CIL}} \rangle \rightarrow_{\text{CIL}} \sigma'_{\text{CIL}} = ([u_{[1,n]}], l, a, h')$  and  $\langle \mathbb{T}[\text{call } m(\text{arg}_{[1,i]}), (s_{(n,i)}, \bar{l}, \bar{a}, \bar{w}_{[1,i]})], \mathbb{T}[\sigma_{\text{CIL}}] \rangle \rightarrow_{\text{JB}} \sigma'_{\text{JB}}$ , then  $\mathbb{T}[\sigma'_{\text{CIL}}] = \bullet \sigma'_{\text{JB}}$ .

Rule  $\langle \text{call } m(\text{arg}_{[1,i]}), \sigma_{\text{CIL}} \rangle \rightarrow_{\text{CIL}} \sigma'_{\text{CIL}}$  requires, to be applied, that the condition

$$\begin{aligned} \langle \text{body}(m(\text{arg}_{[0,i]}), (v_{[1,i]})), ([], \emptyset, [j-1 \mapsto v_j : j \in [1..i]], h) \rangle \\ \downarrow_{\text{CIL}} \\ (s', l', a', h') \end{aligned}$$

is satisfied.

Observe that for  $a = [j-1 \mapsto v_j : j \in [1..i]]$ , we get  $\mathbb{T}[\sigma'_{\text{CIL}}] = ([], \hat{l}, h)$ , where for each  $j \in [1..i]$ ,  $\hat{l}[j + 64 \frac{j}{a}] = v_j$ . Moreover,

$\mathbb{T}[\text{body}(m(\text{arg}_{[0,i]})), (t_1, \dots, t_i), \emptyset, a, h] = \text{body}(m(\text{arg}_{[0,i]}))$ . Therefore, by inductive hypothesis, we get that  $\langle \text{body}(m(\text{arg}_{[0,i]}))(v_{[1,i]}) \rangle$ ,  $([], [j-1 \mapsto v_j : j \in [1..i]], h) \rangle \rightarrow_{\text{JB}} (\hat{s}, \hat{l}, \hat{h})$  is such that  $\mathbb{T}[\sigma'_{\text{CIL}}] = \bullet (s', l', a', h')$  (where  $(\hat{s}, \hat{l}, \hat{h})$ , and in particular  $h' = \hat{h}$ ). It is sufficient now to recall that by Fig. 7 (static call)  $\mathbb{T}[\text{call } m(\text{arg}_{[1,i]}), s_{(n,i)}, \bar{l}, \bar{a}, \bar{w}_{[n,i]}]$  is obtained by applying  $\text{invokestatic } m(\text{arg}_{[1,i]})$  followed by the update of all the local variables passed by reference to the called method, and that by the  $\text{invokestatic}$  rule of Fig. 6,  $\sigma'_{\text{JB}} = ([u_{[1,n]}], \hat{l}, h)$ . Finally, by the definition of  $\mathbb{T}[\sigma]$ , we get  $\mathbb{T}[\sigma'_{\text{CIL}}] = \bullet \sigma'_{\text{JB}}$ .  $\square$

LEMMA A.4. *The translation of  $\text{bgt}$  is correct.*

PROOF. Consider  $\mathbb{T}[\text{bgt } k, \bar{s} :: t_1 :: t_2, \bar{l}, \bar{a}, \bar{w}]$  in the case  $t_1 = t_2 = \text{int}$ , and assume  $\sigma_{\text{CIL}} = ([s_{[1,n]} :: v_1, v_2], l, a, h)$  with  $v_2 > v_1$  (the other case is similar), yielding to  $\langle \text{bgt } l, \sigma_{\text{CIL}} \rangle \rightarrow_{\text{CIL}} \sigma'_{\text{CIL}} = \langle l, \sigma_{\text{CIL}} \rangle$ .

We show that if  $\langle \mathbb{T}[\text{bgt } k, \bar{s} :: t_1 :: t_2, \bar{l}, \bar{a}, \bar{w}], \mathbb{T}[\sigma_{\text{CIL}}] \rangle \rightarrow_{\text{JB}} \sigma'_{\text{JB}}$ , then  $\mathbb{T}[\sigma'_{\text{CIL}}] = \bullet \sigma'_{\text{JB}}$ . By the corresponding rule in Fig. 7,  $\mathbb{T}[\text{bgt } k, \bar{s} :: t_1 :: t_2, \bar{l}, \bar{a}, \bar{w}] = \text{if\_icmpgt } k'$  where  $k' = \text{statementIdx}(\text{getBody}(\text{bgt } k)(k))$ .

By the semantics of  $\text{if\_icmpgt}$  we have that  $\langle \text{if\_icmpgt } k', ([s_{[1,n]} :: v_1 :: v_2], \text{convertLocals}(l, a), h) \rangle \rightarrow_{\text{JB}} \langle k', (s_{[1,n]}, \text{convertLocals}(l, a), h) \rangle$ . As  $\mathbb{T}[\sigma_{\text{CIL}}] = ([s_{[1,n]} :: v_1 :: v_2], \text{convertLocals}(l, a), h)$ , and by definition of  $\text{statementIdx}()$ , we get that  $\mathbb{T}[\langle l, \sigma_{\text{CIL}} \rangle] = \langle k', \mathbb{T}[\sigma_{\text{CIL}}] \rangle = \sigma'_{\text{JB}}$ , and thus  $\mathbb{T}[\langle l, \sigma_{\text{CIL}} \rangle] = \bullet \sigma'_{\text{JB}}$ .  $\square$

LEMMA A.5. *The translation of  $\text{newobj}$  is correct.*

PROOF. Assume  $\sigma_{\text{CIL}} = ([s_{[1,n]} :: v_{[1,i]}], l, a, h)$ . By definition,  $\langle \text{newobj } \tau(\text{arg}_{[1,i]}), \sigma_{\text{CIL}} \rangle \rightarrow_{\text{CIL}} \sigma'_{\text{CIL}} = (s :: r, (l, a, h'))$ , where  $r$  and  $h'$  satisfy the constraints of the corresponding rule of Fig. 5. In particular,  $\text{fresh } \tau h = (r, h')$  allocates the memory for an object of type  $\tau$  on heap  $h$  and returns (i) the reference  $r$  of the freshly allocated object, and (ii) the heap  $h'$  resulting from the allocation of memory on  $h$ .

Let  $\sigma'_{\text{JB}} = \langle \mathbb{T}[\text{newobj } \tau(\text{arg}_{[1,i]}), \bar{s} :: t_{[1,i]}, \bar{l}, \bar{a}, \bar{w}], \mathbb{T}[\sigma_{\text{CIL}}] \rangle$ , and compare  $\mathbb{T}[\sigma'_{\text{CIL}}]$  and  $\sigma'_{\text{JB}}$  componentwise. We may observe that in both cases the store is equal to  $s :: r$ , as the new elements added to the store during the translation in order to implement the object initialization are finally removed when applying the  $\text{invokevirtual}$  call, whose correctness is granted by structural inductive hypothesis. Moreover, the single array in JB for both local variables and method arguments is updated properly by storing and loading the values of the constructor arguments in the expected ordering. Finally, the heap  $h'$  results in both cases from the allocation of corresponding memory on  $h$ .  $\square$

LEMMA A.6. *The translation of  $\text{stind}$  is correct.*

PROOF. By Fig. 7, we have that

$$\mathbb{T}[\text{stind } \bar{s}, \bar{l}, \bar{a}, \bar{w}] = \mathbb{T}[\text{stfld value } \bar{s}, \bar{l}, \bar{a}, \bar{w}] = \text{putfield value}$$

By Fig. 5, we have that  $\langle \text{stind } (s :: r_i :: v, l, a, h) \rangle \rightarrow_{\text{CIL}} (s, l, a, h[r_i \mapsto v])$ . By definition of the translation of concrete states (assuming that the only direct reference in the stack is  $r_i$ ), we have that  $\mathbb{T}[\sigma'_{\text{CIL}}] = (s :: r' :: v, \text{cnvrtLoc}(l, a), h[r' \mapsto \text{value} \mapsto v])$ .

$l(i)]])$  where  $r'$  is a freshly allocated references pointing to a wrapper. Then by Fig. 6, we have that  $\langle \text{putfield value}, (s :: r' :: v, \text{cnvrtLoc}(l, a), h[r' \mapsto [\text{value} \mapsto l(i)]]) \rangle \rightarrow_{\text{JB}} (s, \text{cnvrtLoc}(l, a), h[r' \mapsto [\text{value} \mapsto v]])$ . Finally, we obtain that  $\mathbb{T}_{\sigma} \llbracket (s, l, a, h[r_i \mapsto v]) \rrbracket =_{\bullet} (s, \text{cnvrtLoc}(l, a), h[r' \mapsto [\text{value} \mapsto v]])$  proving the soundness of the translation of `stind`.  $\square$

LEMMA A.7. *The translation of `ldind` is correct.*

PROOF. By Figure 7, we have that

$$\mathbb{T} \llbracket \text{ldind}, \bar{s}, \bar{l}, \bar{a}, \bar{w} \rrbracket = \mathbb{T} \llbracket \text{ldfld value}, \bar{s}, \bar{l}, \bar{a}, \bar{w} \rrbracket = \text{getfield value}$$

By Figure 5, we have that  $\langle \text{ldind}, (s :: r_i, l, a, h) \rangle \rightarrow_{\text{CIL}} (s :: h(r_i), l, a, h)$ . By definition of the translation of concrete states (assuming that the only direct reference in the stack is  $r_i$ ), we have that  $\mathbb{T}_{\sigma} \llbracket (s :: r_i, l, a, h) \rrbracket = (s :: r', \text{cnvrtLoc}(l, a), h[r' \mapsto [\text{value} \mapsto l(i)]])$  where  $r'$  is a freshly allocated references pointing to a wrapper. Then by Fig. 6, we have that  $\langle \text{getfield value}, (s :: r', \text{cnvrtLoc}(l, a), h[r' \mapsto [\text{value} \mapsto l(i)]]) \rangle \rightarrow_{\text{JB}} (s :: l(i)), \text{cnvrtLoc}(l, a), h[r' \mapsto [\text{value} \mapsto v]])$ . Finally, we obtain that  $\mathbb{T}_{\sigma} \llbracket (s :: l(i), l, a, h) \rrbracket =_{\bullet} (s, \text{cnvrtLoc}(l, a), h[r' \mapsto [\text{value} \mapsto l(i)]])$  proving the soundness of the translation of `ldind`.  $\square$

LEMMA A.8. *The translation of `ldloca` is correct.*

PROOF. By Fig. 5, we have that  $\langle \text{ldloca } i, (s, l, a, h) \rangle \rightarrow_{\text{CIL}} (s :: r_i, l, a, h)$  where  $r_i$  is the direct reference pointing to the  $i$ -th local variable. By Fig. 7, we have that `ldloca i` is translated into a sequence of statements that (i) creates a wrapper object containing the value of the  $i$ -th local variable, (ii) stores its reference into an instrumentation variable, and (iii) leaves its reference on the operand stack as well. Then, by definition of the concrete semantics of JB, we obtain a final state  $\sigma'_{\text{JB}}$  appending to the initial stack a reference to the wrapper object whose value is the one of the  $i$ -th local variable. Therefore,  $\sigma'_{\text{JB}} =_{\bullet} \mathbb{T}_{\sigma} \llbracket (s :: r_i, l, a, h) \rrbracket$  since  $=_{\bullet}$  ignores the instrumentation variables. This proves the soundness of the translation of `ldloca`.  $\square$

THEOREM A.9. *The translation of a CIL program into JB code is correct.*

PROOF. We prove that the translation of each statement from CIL to JB depicted in Fig. 7 satisfies the correctness property introduced in Sec. 4.3. In fact, by Lemma A.2, the translation of `ldloc` is correct, and a similar proof can be provided for the `add`, `stloc`, `ldarg`, `ldfld`, and `stfld` statements. The correctness of a new object creation is proved by Lemma A.5. The correctness proof of static and dynamic calls translation has been given in Lemma A.3, and that of `stind`, `ldind`, `ldloca` was proved by Lemmas A.6, A.7 and A.8, respectively. Finally, the correctness of the comparison statements translation is shown in Lemma A.4.  $\square$