

Distinguished Paper Award  
FormalISE 2018

## CIL to Java-bytecode Translation for Static Analysis Leveraging

Pietro Ferrara  
JuliaSoft SRL

Agostino Cortesi  
University of Venice

Fausto Spoto  
University of Verona

# Static Program Analysis

- Automatic analysis of the behavior of software
- Target some specific properties
- Static: without executing the program
  - Need a model of the semantic of the programming language
- Sound: detect all bugs w.r.t. the property of interest
  - False alarms: the analysis fails to prove correct a correct sw

## Properties

- SQL Injection
- NullPointerExceptions
- ...

## Software

```
it(""); } $("#unique").click()  
from_string($("#film")  
= use_unique(array_from  
if (c < 2 * b - 1) {  
    this.trigger("click");  
    if (a[b] && " != a[b] || for  
    if (c.length > 0) {  
        val(); c = array_from  
        for (b = 0; b < c
```

SPA

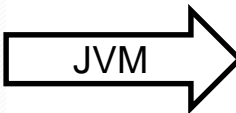
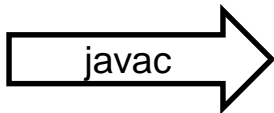


[Injection] possible SQL-injection

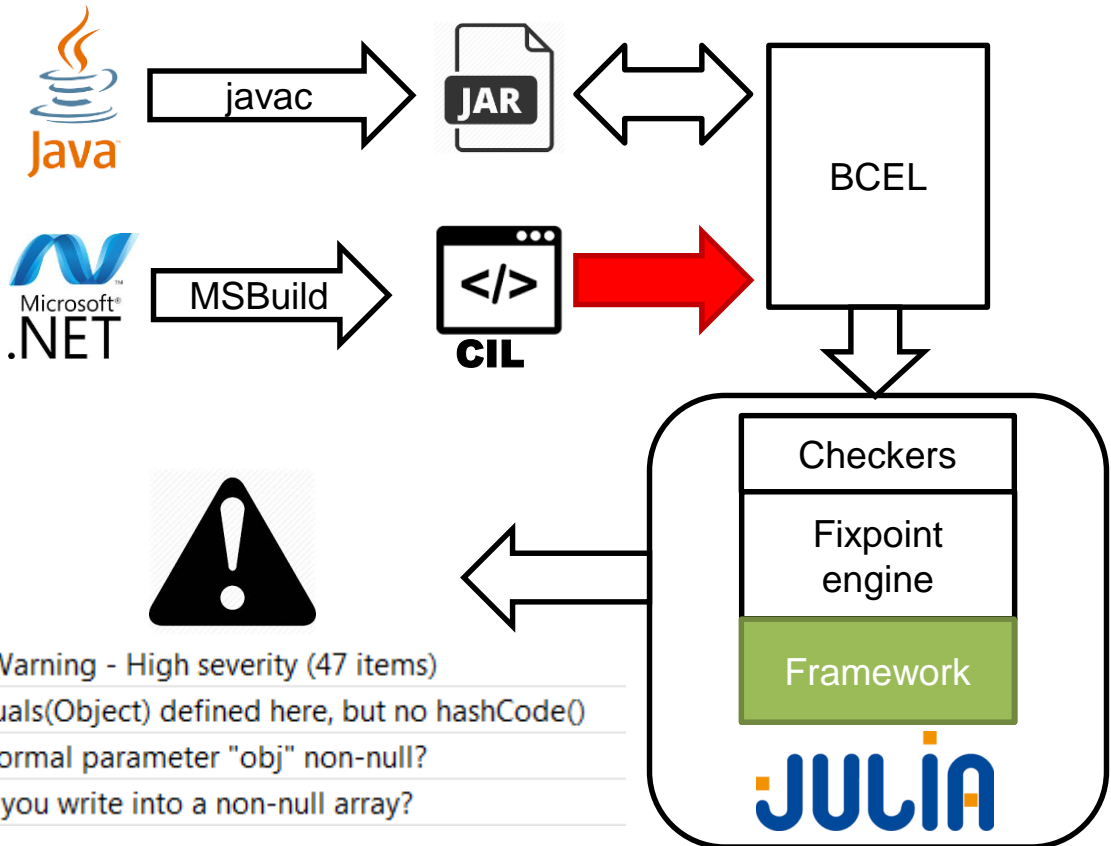
[BasicNullness] is the receiver of  
"getProperty" non-null?

# Bytecode languages

- Machine-independent low-level languages
- Interpreted or Just-In-Time compiled
- Based on local variables, stack of values, heap
- Object-oriented (for the purposes of this work)



# The Julia static analyzer



▲ ⚠ Julia Warning - High severity (47 items)

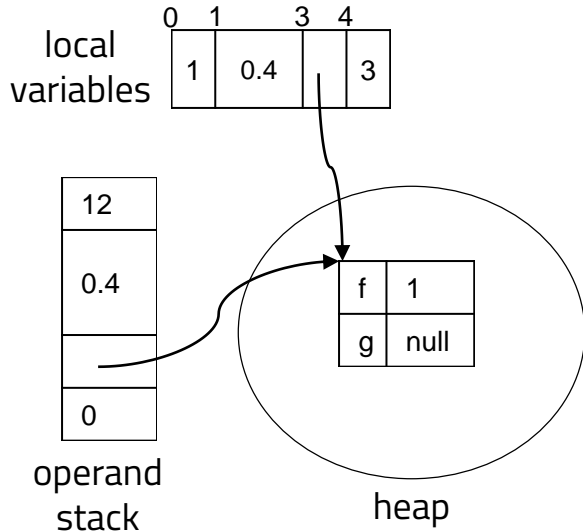
⚠ `equals(Object)` defined here, but no `hashCode()`

⚠ is formal parameter `obj` non-null?

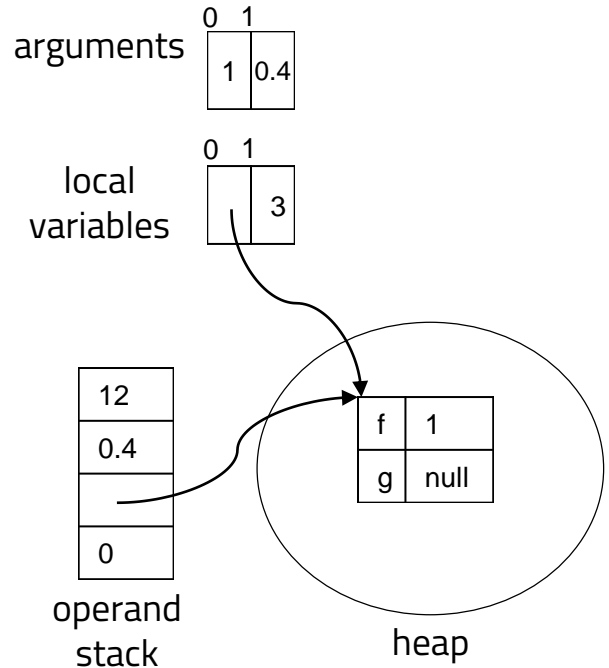
⚠ do you write into a non-null array?

# Concrete states

## Java Bytecode



## MS CIL

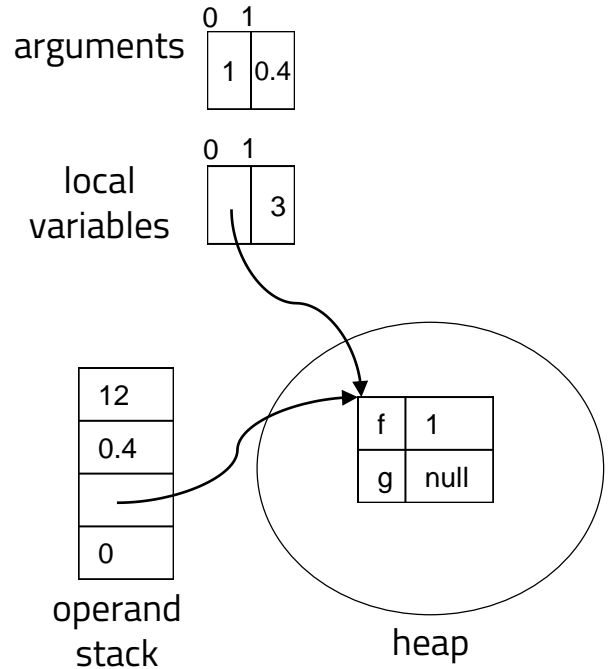
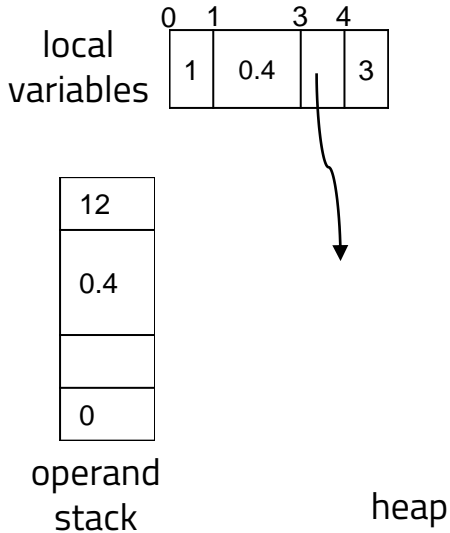


# States translation $T_\sigma [ ] : \Sigma_{CIL} \rightarrow \Sigma_{JB}$

$\Sigma_{JB}$

$T_\sigma [ ]$

$\Sigma_{CIL}$



# Java Bytecode (JB) vs. CIL

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)
invokeresult invoketarget	call	(call)
new getfield putfield if_icmpeq		
dup dup2		(stack)
	ldloca i stind ldind	(pointers)

- CIL more expressive than JB
- Some uses of pointers do not have an equivalent in JB
- Focus on CIL derived from safe C# code

JB: typed  
CIL: untyped

JB: local vars  
CIL: local vars + args

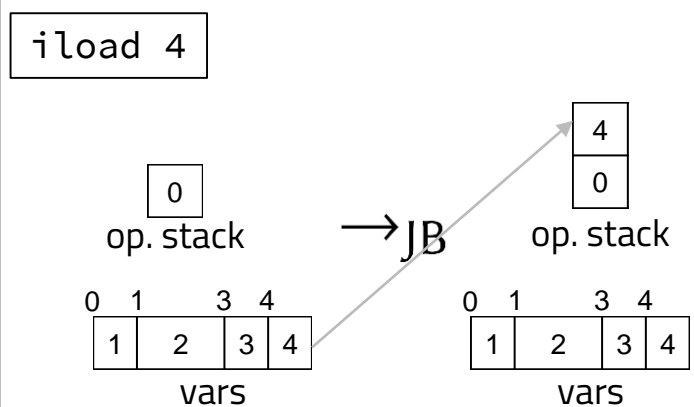
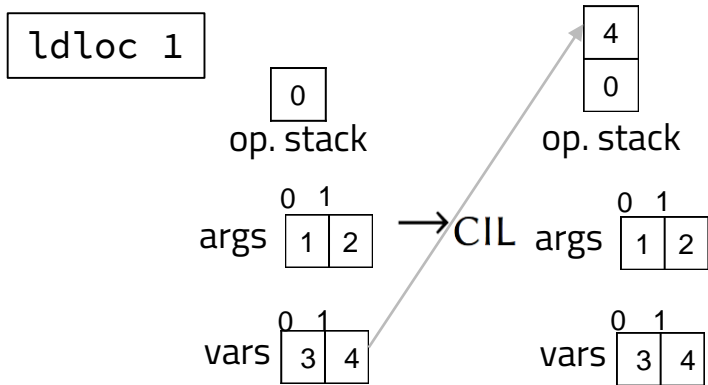
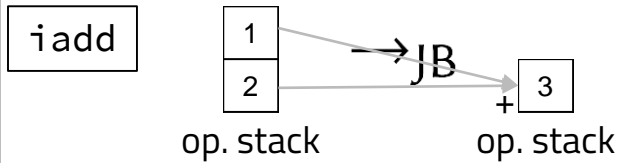
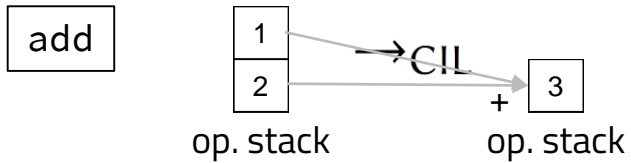
JB: static and dynamic  
CIL: generic

JB: typed  
CIL: untyped

CIL: unique values  
JB: 32- and 64-bit values

CIL: direct pointers

# Concrete semantics





# Concrete semantics $\rightarrow_{\text{CIL}}$

$$\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)} \text{ (add)} \qquad \frac{}{\langle \text{ldloc } i, (s, l, a, h) \rangle \rightarrow_{\text{CIL}} (s :: l(i), l, a, h)} \text{ (ldloc)} \\
 \frac{}{\langle \text{stloc } i, (s :: v, l, a, h) \rangle \rightarrow_{\text{CIL}} (s, l[i \mapsto v], a, h)} \text{ (stloc)} \qquad \frac{}{\langle \text{ldarg } i, (s, l, a, h) \rangle \rightarrow_{\text{CIL}} (s :: a(i), l, a, h)} \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)), ([], \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')} \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)), ([], \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')} \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)), ([], \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')} \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)} \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'])} \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 1, (s, l, a, h) \rangle} \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)} \text{ (bgt false)} \\
 \frac{}{\langle \text{ldloca } i, (s, l, a, h) \rangle \rightarrow_{\text{CIL}} (s :: r_i, l, a, h)} \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])} \text{ (stind)} \\
 \frac{}{\langle \text{dup}, (s :: v, l, a, h) \rangle \rightarrow_{\text{CIL}} (s :: v :: v, l, a, h)} \text{ (dup)} \qquad \frac{}{\langle \text{ldind}, (s :: r_i, l, a, h) \rangle \rightarrow_{\text{CIL}} (s :: h(r_i), l, a, h)} \text{ (ldind)}
 \end{array}$$

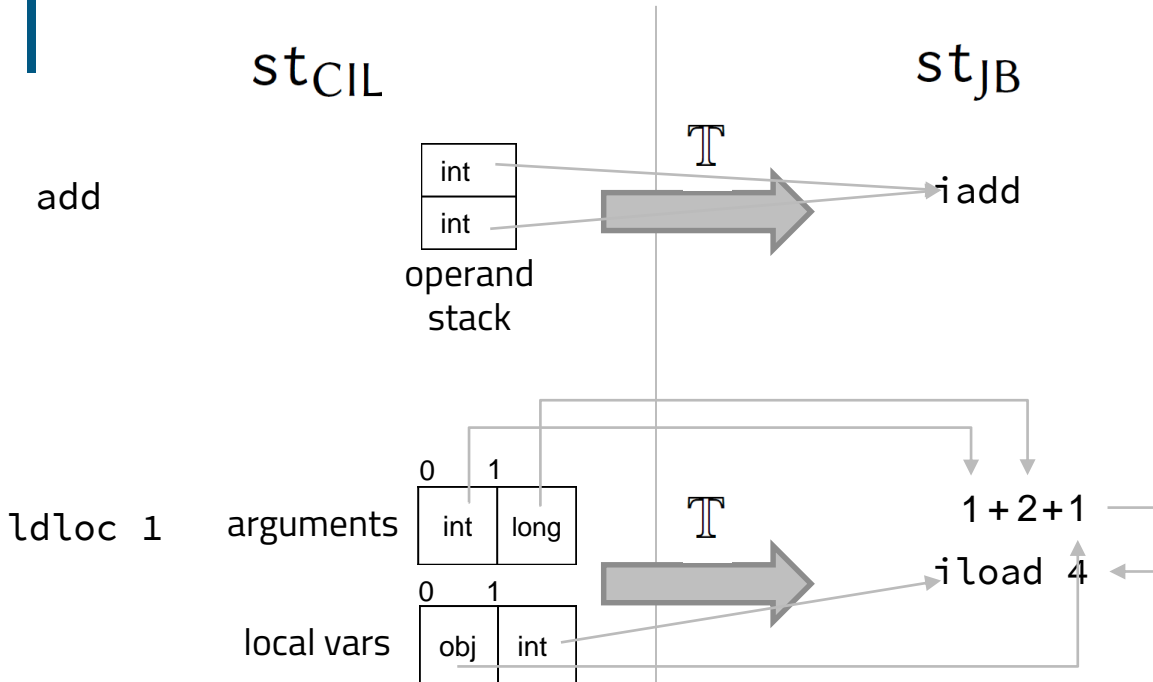
Figure 5: Concrete CIL semantics.

# Concrete semantics $\rightarrow_{\text{JB}}$

$$\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)} \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)} \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)} \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)} \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)} \text{ (ladd)} \\
 \\
 \frac{x = \mathcal{JVM}\text{prefix}(\text{typeOf}(l(i)))}{\langle \text{xload } i, (s, l, h) \rangle \rightarrow_{\text{JB}} (s :: l(i), l, h)} \text{ (xload)} \qquad \frac{x = \mathcal{JVM}\text{prefix}(\text{typeOf}(v))}{\langle \text{xstore } i, (s :: v, l, h) \rangle \rightarrow_{\text{JB}} (s, [i \mapsto v], h)} \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')} \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')} \text{ (invokestatic)} \\
 \\
 \frac{\text{fresh}(T, h) = (r, h')}{\langle \text{new } T, (s, l, h) \rangle \rightarrow_{\text{JB}} (s :: r, l, h')} \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)} \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'])} \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} \left( \begin{array}{c} \text{if\_icmpgt} \\ \text{true} \end{array} \right) \\
 \\
 \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)} \left( \begin{array}{c} \text{if\_icmpgt} \\ \text{false} \end{array} \right)
 \end{array}$$

Figure 6: Concrete JB semantics.

# Statement translation $T[st_{CIL}, K] = st_{JB}$



# Statement translation $\mathbb{T}[\text{st}_{\text{CIL}}, \mathbb{K}] = \text{st}_{\text{JB}}$

$$\mathbb{T}[\text{dup}, \bar{s} :: t, \bar{l}, \bar{a}, \bar{w}] = \begin{cases} \text{dup} & \text{if } t \neq \text{Long} \\ \text{dup2} & \text{if } t = \text{Long} \end{cases}$$

$$\mathbb{T}[\text{add}, \bar{s} :: t_1 :: t_2, \bar{l}, \bar{a}, \bar{w}] = \begin{cases} \text{iadd} & \text{if } t_1 = t_2 = \text{Int} \\ \text{ladd} & \text{if } t_1 = t_2 = \text{Long} \end{cases}$$

$$\mathbb{T}[\text{ldloc } i, \bar{s}, \bar{l}, \bar{a}, \bar{w}] = x\text{load } j \text{ where } j = |\bar{a}| + 64 \frac{|\bar{a}|}{a} + i + 64 \frac{i}{a} \wedge x = \text{JVMprefix}(\text{typeOf}(\bar{l}(i)))$$

$$\mathbb{T}[\text{stloc } i, \bar{s} :: t, \bar{l}, \bar{a}, \bar{w}] = x\text{store } j \text{ where } j = |\bar{a}| + 64 \frac{|\bar{a}|}{a} + i + 64 \frac{i}{a} \wedge x = \text{JVMprefix}(\text{typeOf}(\bar{l}(i)))$$

$$\mathbb{T}[\text{ldarg } i, \bar{s}, \bar{l}, \bar{a}, \bar{w}] = x\text{load } j \text{ where } j = i + 64 \frac{i}{a} \wedge x = \text{JVMprefix}(\text{typeOf}(\bar{a}(i)))$$

$$\begin{aligned} \mathbb{T}[\text{call } m(\text{arg}_1, \dots, \text{arg}_i), \\ \bar{s} :: t_1 :: \dots :: t_i, \bar{l}, \bar{a}, \bar{w} :: p_1 :: \dots :: p_i] &= \text{invoke} ; \text{aload } p_{idx_1}^1 ; \text{getfield value} ; x_{idx_1} \text{store } p_{idx_1}^2 ; \dots \\ &\quad \dots \text{aload } p_{idx_j}^1 ; \text{getfield value} ; x_{idx_j} \text{store } p_{idx_j}^2 ; \\ &\quad \text{where invoke} = \begin{cases} \text{invokestatic } m(\text{arg}_1, \dots, \text{arg}_i) & \text{if } \text{isStatic}(m(\text{arg}_1, \dots, \text{arg}_i)) \\ \text{invokevirtual } m(\text{arg}_1, \dots, \text{arg}_i) & \text{otherwise} \end{cases} \\ &\quad \{idx_1, \dots, idx_j\} = \{k : \text{arg}_k \in \text{Ref}_{\text{Loc}}\} \\ &\quad \forall k \in [1..j] : x_{idx_k} = \text{JVMprefix}(\text{typeOf}(\bar{l}(p_{idx_k}^2))) \wedge \forall r \in [1..i] : p_i = (p_i^1, p_i^2) \end{aligned}$$

$$\begin{aligned} \mathbb{T}[\text{newobj } T(a_1, \dots, a_i), \\ \bar{s} :: t_1 :: \dots :: t_i, \bar{l}, \bar{a}, \bar{w}] &= x_i \text{store } idx_i ; \dots ; x_1 \text{store } idx_1 ; \text{new } T ; \text{dup} ; \\ &\quad x_1 \text{load } idx_1 ; \dots ; x_i \text{load } idx_i ; \text{invokevirtual } < \text{init} > (\text{arg}_1, \dots, \text{arg}_i) \\ &\quad \text{where } \forall j \in [1..i] : x_j = \text{JVMprefix}(a_j) \wedge idx_j = \text{freshIdx}(\text{newobj } T(a_1, \dots, a_i), j) \end{aligned}$$

$$\mathbb{T}[\text{ldfld } f, \bar{s} :: t_0, \bar{l}, \bar{a}, \bar{w}] = \text{getfield } f$$

$$\mathbb{T}[\text{stfld } f, \bar{s} :: t_0 :: t_v, \bar{l}, \bar{a}, \bar{w}] = \text{putfield } f$$

$$\mathbb{T}[\text{bgt } k, \bar{s} :: t_1 :: t_2, \bar{l}, \bar{a}, \bar{w}] = \text{if\_icmpgt } k' \text{ where } k' = \text{statementIdx}(\text{getBody}(\text{bgt } k)(k)) \text{ if } t_1 = t_2 = \text{Int}$$

$$\begin{aligned} \mathbb{T}[\text{ldloca } i, \bar{s}, \bar{l}, \bar{a}, \bar{w}] &= \mathbb{T}[\text{newobj } \text{WrapRef}() ; \text{dup2} ; \text{stloc } j ; \text{ldloc } i ; \text{stfld value}, \bar{s}, \bar{l}, \bar{a}, \bar{w}] \\ &\quad \text{where } j = \text{freshIdx}(\text{ldloca } i, 0) \end{aligned}$$

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

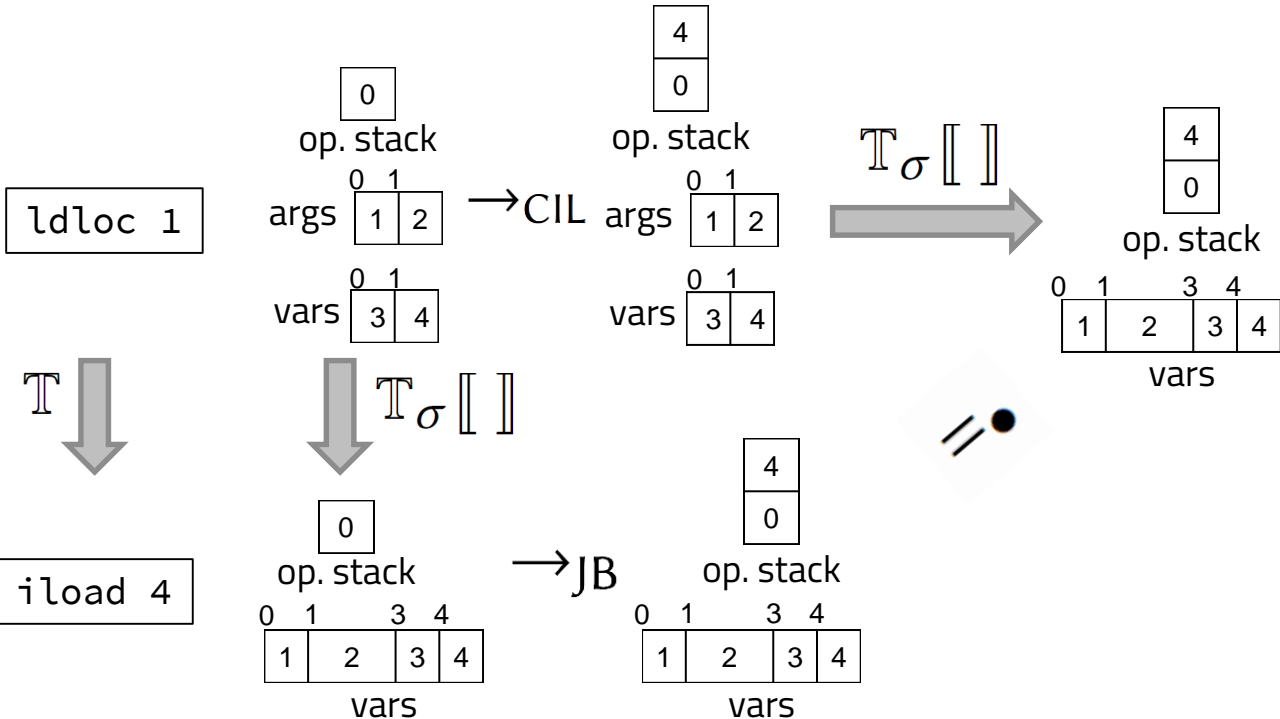
$$\mathbb{T}[\text{ldind}, \bar{s}, \bar{l}, \bar{a}, \bar{w}] = \mathbb{T}[\text{ldfld value}, \bar{s}, \bar{l}, \bar{a}, \bar{w}]$$

# Correctness

$$\forall st \in St_{CIL}, \sigma_{CIL} \in \Sigma_{CIL} : \langle st, \sigma_{CIL} \rangle \rightarrow_{CIL} \sigma'_{CIL} \text{ and } \langle T[st, K], T_\sigma[\sigma_{CIL}] \rangle \rightarrow_{JB} \sigma'_{JB}$$

$$\Downarrow$$

$$T_\sigma[\sigma'_{CIL}] = \bullet \sigma'_{JB}$$



# Experimental results

## Efficiency and precision

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

Program	LOC	met.	fail	Tr. t.	An.t.	Al	F	Prec.
CodeHub	32,510	4,887	0	0'07"	0'43"	9	1	89%
SignalR	71,207	6,610	3	0'07"	0'50"	8	1	88%
Dapper	22,513	1,058	0	0'07"	0'29"	13	3	77%
ShareX	171,580	11,568	14	0'58"	2'08"	57	0	100%
Nancy	109,139	8,817	0	0'07"	1'25"	18	1	94%
Total	406,949	32,940	17	1'26"	4'35"	105	6	94%

- 5 most popular GitHub repositories written in C#
- Efficiency:
  - Analyze industrial-size software in a few minutes
  - Translation time comparable to the analysis time.
- Precision:
  - 6/105 alarms (6%) are false because the translation

# Experimental results

## Libraries

Table 2: Experimental results on libraries.

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

- 10 largest system libraries of .NET framework
- Study how much code we can translate
  - 99.13% of the methods, with a worst case of 96.93%

# Limitations and future work

- Native [and unsafe] code
  - Code written in languages other than bytecode
  - Linked to and executed by the virtual machine/runtime env.
  - Static analysis requires a (manually written) model
  - Cannot be translated by our approach
    - We cannot execute the translation of a CIL program with the JVM
- “Naming conventions” introduced by the compiler
  - Need some polishing to link warnings back to the source code
- Support .NET frameworks
  - ASP.NET, etc..



The contents of this document are property and copyright © of JuliaSoft S.r.l. The document is intended solely for informative purposes. All trademarks are property of their owners. The content of this document cannot be totally or partially copied, reproduced, transferred, uploaded, published or distributed in any way without the previous written consent of JuliaSoft S.r.l.



**JULIASOFT**  
code analysis reinvented



Corvallis Group  
Management and Coordination by Corvallis Holding SpA

Lungadige Galtarossa 21  
37133 Verona, Italy  
Tel. + 39 045.2081901  
info@juliasoft.com  
www.juliasoft.com

Administrative Offices  
Via Savelli, 56 - 35129 Padova, Italy  
T +39 049.8434511 | F +39 049.8434555



**Thank you.**  
**Questions?**

**C O N T A C T S :**

*pietro.ferrara@juliasoft.com*