

Analyzing Java programs with Octagons

September 20

Francesco Logozzo

ENS

Goals

We want a static analysis for **sequential** Java that

- ❑ Is sound
 - ❑ \neq ESC/Java, ESC/Java 2
- ❑ Is automatic
 - ❑ \neq Jive
- ❑ Infers **real** invariants
 - ❑ \neq Daikon
- ❑ Analyze classes in isolation
 - ❑ \neq Excelsion Flawdetector
- ❑ Does not rely on user annotations
 - ❑ \neq JML-based approaches as LOOP, Krakatoa

Outline

- ❑ Concrete Semantics
- ❑ Abstract Semantics
- ❑ Comparison with other tools
- ❑ Conclusions & Future work

Syntax

- ❑ A **class** is a triplet $\langle \text{init}, F, M \rangle$ where
 - ❑ `init` is the class constructor
 - ❑ `F` is a set of variables
 - ❑ `M` is a set of function definitions
- ❑ For simplicity assume
 - ❑ Just one constructor
 - ❑ All the fields are protected
 - Otherwise use `get_f/set_f`

Semantic domains

- ❑ The set of **environments** is $\text{Env} = [\text{Var} \rightarrow \text{Addr}]$ where
 - ❑ Var is a set of variables
 - ❑ Addr is a set of addresses
- ❑ The the set of **stores** is $\text{Store} = [\text{Addr} \rightarrow \text{Val}]$ where
 - ❑ Val is a set of values such that $\text{Env} \subseteq \text{Val}$
 - ❑ $\phi \in \text{Val}$ is the void value
- ❑ The environment of an object is stored in a certain memory location
- ❑ The address of such a location is the object identity
- ❑ A program state is $\sigma \in \text{Env} \times \text{Store}$

Constructor and method semantics

- The semantics of the **constructor** is a function

$$i[[init]] \in [Val \times Store \rightarrow \mathcal{P}(Env \times Store)]$$

- The semantics of a **method** m is a function

$$m[[m]] \in [Val \times Env \times Store \rightarrow \mathcal{P}(Val \times Env \times Store)]$$

- If $Addr \subseteq Val$ then the method may expose a part of object state to the context

Object semantics

- ❑ The semantics of an object is given by a set of traces
 - ❑ Each trace represents a possible evolution history of the object
- ❑ The first state represents the object right after its creation
- ❑ Each further state is the result of an interaction between the object and a context
 - ❑ The context invokes a method of the object; or
 - ❑ It modifies a memory location that is reachable from the object environment
- ❑ The semantics of a class considers all its instances

Initial states

- The set of **interaction** states is

$$\Sigma = \text{Env} \times \text{Store} \times \text{Val} \times \mathcal{P}(\text{Addr})$$

- The initial states of an **object** o is

$$S_0\langle v, s \rangle = \{ \langle e', s', \phi, \emptyset \rangle \mid \langle e', s' \rangle \in i[\text{init}](v, s) \}$$

- The initial states of a **class** are

$$\mathbb{I}[\text{init}] = \cup \{ S_0\langle v, s \rangle \mid \langle v, s \rangle \in \text{Val} \times \text{Store} \}$$

Direct interaction

- Let $\text{reachable} \in [(\text{Val} \times \text{Store}) \rightarrow \mathcal{P}(\text{Addr})]$
 - It determines all the memory addresses reachable from an address in Val
- The **collecting semantics** of a method m is

$$\begin{aligned} \mathbb{M}[[m]](S) = \{ \langle e', s', v', \text{Esc}' \rangle \mid & \langle e, s, v, \text{Esc} \rangle \in S, v_{\text{in}} \in \text{Val}, \\ & \langle v', e', s' \rangle \in m[[m]](v_{\text{in}}, e, s), \\ & \text{Esc}' = \text{Esc} \cup \text{reachable}(v', s') \} \end{aligned}$$

- The transition function for the **direct interaction** is

$$\text{next}_{\text{dir}}(\sigma) = \cup \{ \mathbb{M}[[m]](\{\sigma\}) \mid m \in M \}$$

Indirect interaction

- Let $\text{update} \in [\text{Addr} \times \text{Store} \rightarrow \mathcal{P}(\text{Store})]$ be defined as

$$\text{update}(a, s) = \{s' \mid \exists v \in \text{Val}. s' = s[a \mapsto v]\}$$

- The **context interactions** are

$$\begin{aligned} \text{Context}(S) = \{ \langle e, s', \phi, \text{Esc} \rangle \mid \langle e, s, v, \text{Esc} \rangle \in S, \\ \exists a \in \text{Esc}. s' \in \text{update}(a, s) \} \end{aligned}$$

- The transition function for the **indirect interaction** is

$$\text{next}_{\text{ind}}(\sigma) = \text{Context}(\{\sigma\})$$

Class semantics, $\mathbb{C}[\cdot]$

- **Th.** $\mathbb{C}[A]$ can be expressed in fixpoint form as

$$\begin{aligned}\mathbb{C}[A] = \text{lfp } \lambda T. \mathbb{I}[\text{init}] \cup \{ & \sigma_0 \rightarrow \dots \rightarrow \sigma_n \rightarrow \sigma' \mid \\ & \sigma_0 \rightarrow \dots \rightarrow \sigma_n \in T, \\ & \sigma' \in \text{next}_{\text{dir}}(\sigma_n) \cup \text{next}_{\text{ind}}(\sigma_n) \}\end{aligned}$$

- $\mathbb{C}[A]$ can be shown **sound** and **complete** w.r.t. a whole program object-oriented semantics
- The class reachable states of a class A are

$$\mathbb{C}[A] = \text{lfp } \lambda S. \mathbb{I}[\text{init}] \cup \bigcup_{m \in M} \mathbb{M}[m](S) \cup \text{Context}(S)$$

Outline

- ❑ Concrete Semantics
- ❑ Abstract Semantics
- ❑ Comparison with other tools
- ❑ Conclusions & Future work

Abstraction of the reachable states

- Assume an abstract domain \bar{D} such that

$$\langle \mathcal{P}(\Sigma), \subseteq, \emptyset, \Sigma, \cup, \cap \rangle \xrightleftharpoons[\alpha]{\gamma} \langle \bar{D}, \bar{\subseteq}, \bar{\perp}, \bar{\top}, \bar{\cup}, \bar{\cap} \rangle$$

- A sound approximation of the initial states

$$\mathbb{I}[\text{init}] \subseteq \gamma(\bar{\mathbb{I}}[\text{init}])$$

- A sound approximation of the method semantics

$$\forall S \in \mathcal{P}(\Sigma). \mathbb{M}[\mathbf{m}](S) \subseteq \gamma(\bar{\mathbb{M}}[\mathbf{m}](\alpha(S)))$$

- A sound approximation of the context

$$\forall S \in \mathcal{P}(\Sigma). \text{Context}(S) \subseteq \gamma(\overline{\text{Context}}(\alpha(S)))$$

Abstract class invariant

□ **Th.** (SAS'03) Let A be a class. Then

$$\bar{C}[[A]] = \text{lfp } \lambda X. \bar{I}[[\text{init}]] \sqcap \overline{\bigsqcup_{m \in M} \bar{M}[[m]](X)} \sqcap \overline{\text{Context}}(X)$$

is such that

$$c[[A]] \subseteq \gamma_{\Sigma} \circ \gamma(\bar{C}[[A]])$$

□ In practice:

□ We have to define $\bar{I}[[\cdot]]$ and $\bar{M}[[\cdot]]$

□ The abstract domain \bar{D} must track object aliasing

⇒ In particular the exposed objects

⇒ $\overline{\text{Context}}$ sets them to \bar{T}

Abstract domain: Abstract environment

- ❑ The abstract domain is $\overline{D} = \overline{Env} \times \overline{Store}$
- ❑ An abstract environment overapproximates the set of addresses a variable may be tied to
- ❑ The abstract environments are

$$\overline{Env} = [\text{Var} \rightarrow \mathcal{P}(\overline{Addr})]$$

- ❑ \overline{Addr} is the set of abstract addresses
- ❑ The meaning of a $\overline{\rho} \in \overline{Env}$ is

$$\gamma(\overline{\rho}) = \{\rho \in Env \mid \forall \mathbf{x}. \rho(\mathbf{x}) \in \gamma_{addr}(\overline{\rho}(\mathbf{x}))\}$$

Abstract domain: Abstract store

- Assume the only primitive type being `int`
- The abstract stores are

$$\overline{\text{Store}} = \underbrace{([\overline{\text{Addr}} \rightarrow \mathbb{N}] \times \text{Oct})}_{\text{Approximation of ints}} \times \underbrace{[\overline{\text{Addr}} \rightarrow \overline{\text{Env}}]}_{\text{Approximation of objects}} \times \underbrace{\mathcal{P}(\overline{\text{Addr}})}_{\text{Summary locations}}$$

- Let $\langle \langle g, o \rangle, f, S \rangle \in \overline{\text{Store}}$ then
 - $\text{dom}(f) \cap \text{dom}(g) = \emptyset$
 - \implies An address denotes objects **or** integers
 - \implies Typed addresses

- S is the set of summary locations
 - if $\bar{a} \in S$ then \bar{a} stands for (infinitely?) **many** concrete addresses
 - Otherwise \bar{a} stands for exactly **one** concrete address
- $f \in [\overline{\text{Addr}} \rightarrow \overline{\text{Env}}]$ associates an abstract environment to each abstract object
- $g \in [\overline{\text{Addr}} \rightarrow \mathbb{N}]$ maps `ints` to octagon's dimensions
 - A dimension is associated to each address
 - Useful in implementation
 - Max dimension not known statically

Example

```
class Box { int i = 0; }  
// ...  
  
//  $\langle \rho, \langle \langle g, o \rangle, f, S \rangle \rangle$ ,  $o$  with  $n$  dimensions  
Box b = null;  
if( ... )  
    b = new Box();  
else skip;  
//  $\langle \rho[b \mapsto \{a\}], \langle \langle g[a \mapsto (n + 1)], o[x_{n+1} = 0] \rangle, f, S \rangle \rangle$ 
```

Note: At the join point we add the constraint $x_{n+1} = 0$

\implies In the false branch: $x_{n+1} = \perp$ is tacitly assumed

Abstract operations and transfer functions

- The abstract operations are defined as one expects
 - e.g., \sqcup , \sqcap , ∇ are pointwisely defined
- Most of the transfer functions are as usual
 - Sequence, `if`, `while`, ...
- Next we describe the most interesting ones
 - Set up the initial state for the methods' analysis
 - Creation of a new object (`new`)
 - Assignment (`:=`)
 - Projection of exposed addresses

Handling parameters

- ❑ Let $I^k = \langle \rho, \langle \langle g, o \rangle, f, S \rangle \rangle$ be the state at the k -th iteration for the class invariant
- ❑ The initial state I_0 for a method m is given by
 - ❑ The class invariant computed so far (I^k)
 - ❑ The parameters of the method
- ❑ Parameters can be `ints` or objects
 - ❑ If it is an `int` we simply add it to the state, e.g.:

```
public void m(int v) { ... }
```


then $I_0 = \langle \rho[v \mapsto \bar{a}], \langle \langle g[\bar{a} \mapsto n + 1], o' \rangle, f, S \rangle \rangle$
where o' is the octagon o extended to the $n + 1$ dimension

- If they are objects, we have to pay attention to

- **aliasing**, eg:

```
class A { B bRef; }
class B { int i; }
class AnalyzeMe {
// ...
public void m(A a, B b) { ... } }
```

⇒ a.bRef and b may alias

- **cycles**, eg:

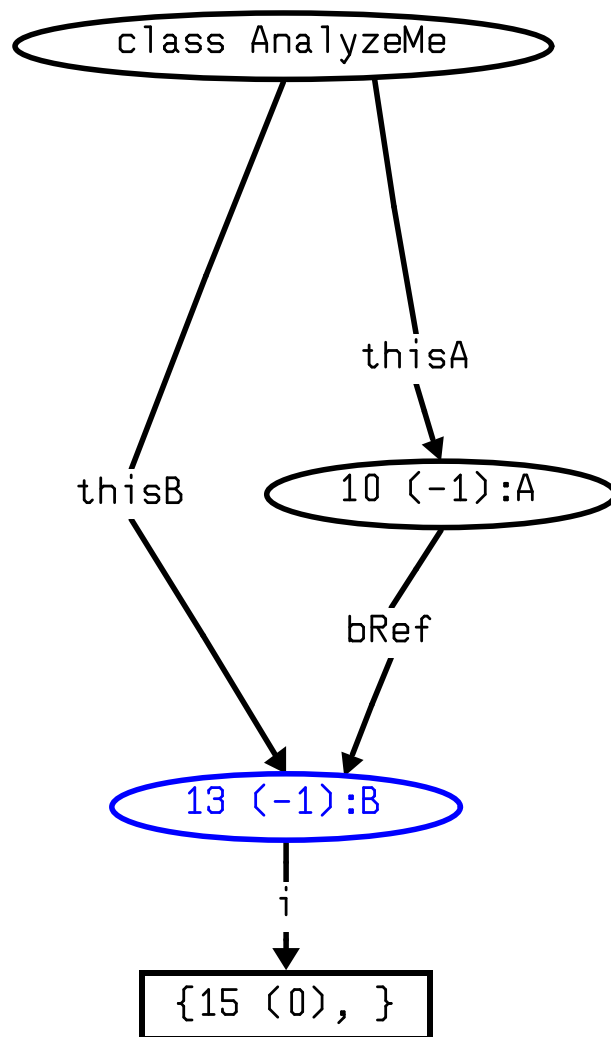
```
class Cycle_1 { Cycle_2 c1to2; // ... }
class Cycle_2 { Cycle_1 c2to1; // ... }
class AnalyzeMeWithCycles {
//
public void n(Cycle_1 c) { ... } }
```

⇒ c and c.c1to2.c2to1 may alias

- ❑ We use summary locations (we assume worst case) for objects of a type that appear **at least twice** in method's parameters
- ❑ Ex: AnalyzeMe: $I_0 = \langle \rho', \langle \langle g', o' \rangle, f', S' \rangle \rangle$ where
 - ❑ $\rho' = \rho[a \mapsto \{13\}, b \mapsto \{10\}]$
 - ❑ $g' = g[15 \mapsto (n + 1)]$
 - ❑ $o' = o$ extended with an $n + 1$ th dimension
 - ❑ $f' = f[13 \mapsto \langle i \mapsto \{15\} \rangle, 10 \mapsto \langle \text{bRef} \mapsto \{13\} \rangle]$
 - ❑ $S' = S \cup \{13, 15\}$
- ❑ (A more understandable picture will come)

Output of the Analyzer

No exceptions



```
# of dimensions: 8
# of allocated cells: 288
# of constraints: 0
# of addresses: 1 : [15 (0)]
```

new

- ❑ Objects are created through the `new` statement

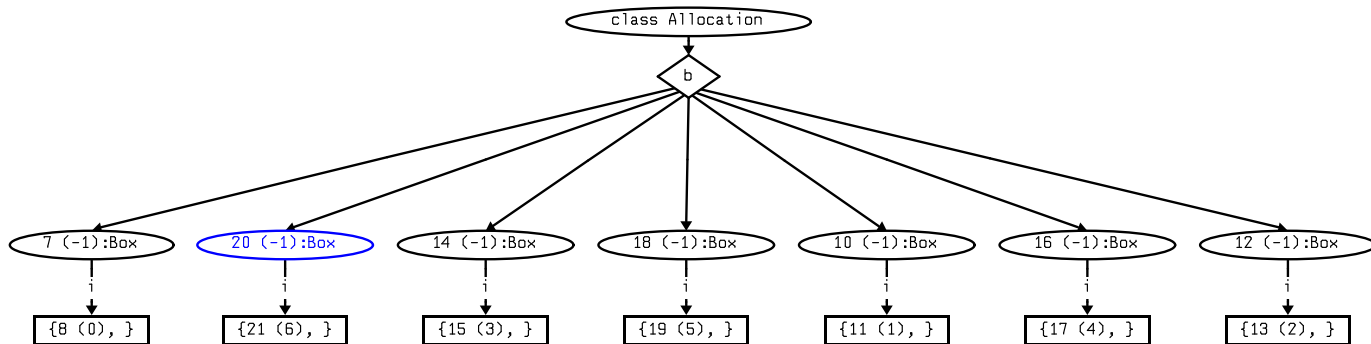
`A a = new A();`

- ❑ We must not to create **infinitely many** objects
 - ❑ Loops and successive invocations of the same method
- ❑ We handle the first k invocations of a `new` exactly and then we create a summary location
 - ❑ k is a parameter of the analysis

Example

```
class Box { int i; }
```

```
class Allocation {  
    Box b;  
    Allocation() { b = new Box(); }  
    public void m() { b = new Box(); }  
}
```



No exceptions

Assignment

- ❑ Let $e1 = e2$ be an assignment
- ❑ **Two cases**: whether $e1$ is an object or an `int`
 - ❑ Java is a typed language...
- ❑ If it is an object, then create an alias for $e2$
⇒ Update the abstract environment
- ❑ If it is an `int`, then we have to update the octagon

- e1 and e2 may evaluate to several addresses, eg:

$$a.x = b.c.y + 2$$

with, eg, $I_0 = \langle \rho, \langle \langle g, o \rangle, f, S \rangle \rangle$ and

- $\rho = \langle a \mapsto \{a_1, a_2\}, b \mapsto \{a_3\} \rangle$

- $g = \langle a_4 \mapsto d_0, a_5 \mapsto d_1, a_7 \mapsto d_2 \rangle$

- $f = \langle a_1 \mapsto \langle x \mapsto \{a_4\} \rangle, a_2 \mapsto \langle x \mapsto \{a_5\} \rangle, a_3 \mapsto \langle c \mapsto \{a_6\} \rangle, a_6 \mapsto \langle c \mapsto \{a_7\} \rangle \rangle$

The possible assignments are

$$a_4 = a_7 + 2 \text{ and } a_5 = a_7 + 2$$

that become the octagon constraints

$$d_0 = d_2 + 2 \text{ and } d_1 = d_2 + 2$$

finally, the octagon after the assignment is the union of the two possible octagons:

$$o' = o.\text{assign}(d_0 = d_2 + 2) \sqcup o.\text{assign}(d_1 = d_2 + 2)$$

so that the state after the assignment is

$$I_1 = \langle \rho, \langle \langle g, o' \rangle, f, S \rangle \rangle$$

Assignment of `ints`: Summary

- ❑ Let $e1 = e2 + c$ and $I_0 = \langle \rho, \langle \langle g, o \rangle, f, S \rangle \rangle$
- ❑ Let $A[[e1]](I_0)$ and $A[[e2]](I_0)$ the addresses $e1$ (resp. $e2$) may evaluate to in I_0
- ❑ Let

$$\text{Ass} = \{g(a_1) = g(a_2) + c \mid a_1 \in A[[e1]](I_0), a_2 \in A[[e2]](I_0)\}$$

- ❑ Then the octagon after the assignment is

$$o' = \bigsqcup_{d_1 = d_2 + c \in \text{Ass}} o.\text{assign}(d_1 = d_2 + c)$$

Method exit point

- ❑ At method's exit point we collect the addresses that **escape** from the class
 - ❑ Somehow similar to the concept of ownership types
- ❑ They are reachable
 - ❑ by the method's / constructor parameters
 - ❑ by the methods return value
- ❑ We set the corresponding location to \top
 - ❑ The context can do everything on an exposed location...

Outline

- ❑ Concrete Semantics
- ❑ Abstract Semantics
- ❑ Comparison with other tools
- ❑ Conclusions & Future work

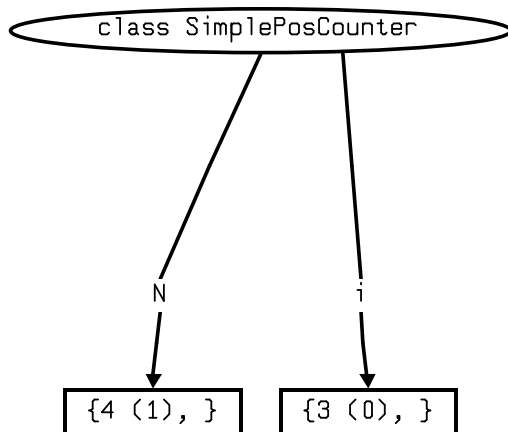
A Running example

```
class SimplePosCounter {
    private int i, N;

    SimplePosCounter(int N0) {
        assert -N0 <= -1;
        this.N = N0;
        this.i = 0; }
    public void add() {
        if(this.i - this.N < 0)
            i = i + 1; }
    public void sub() {
        if(-this.i <= this.N < 0)
            i = i - 1; }
    public void check() {
        assert N - i <= -1; }
}
```


My analyzer

- ❑ **Discover** that the assertion in `check()` is **always** violated
 - ❑ in 266ms on my laptop (Centrino 1.2Ghz, 768Mbyte, WinXp, Java 1.5_04)



No exceptions

```
# of dimensions: 8
# of allocated cells: 288
# of constraints: 4
# of addresses: 2 : [4 (1), 3 (0)]

-a3 <= 0
-a4 <= -1
a3 - a4 <= 0
-a3 - a4 <= -1
```

ESC/Java 2

- ❑ Let us consider ESC/Java 2
 - ❑ Unsound, but not for the example...
- ❑ It is made to discover errors...
- ❑ However, the wrong program passes without not even a warning...

```
< omissis >
```

```
SimplePosCounter: check() ...
```

```
[0.016 s 6625592 bytes] passed
```

- ❑ Is is even slower than ours: 297ms

Daikon

- ❑ It executes the program, and it infers properties from the traces
- ❑ Very popular...
- ❑ It needs that user supplies test cases
 - ❑ Unlike what suggested by the tool's home page
 - ⇒ Cannot analyze the example directly
 - ⇒ Write the test case by myself

```
class TesterSimplePosCounter {
    public static void main(String[] unused) {

// ... Repeat 10 times
        s = new SimplePosCounter( ... );

// ... repeat Math.random() times
        double tmp = Math.random();

        if(tmp <= 0.475)
            s.add();
        else if(tmp <= 0.999)
            s.sub();
        else
            s.check();
    }
}
```

- ❑ If the Daikon executes `s.check()` then it stops
⇒ Execute it many times before success
- ❑ Finds wrong “invariants” ! eg

<omissis>

```
SimplePosCounter:::OBJECT  
this.i >= 0  
this.i < this.N
```

Excelsior Flawdetector

- ❑ A commercial product (free trial) to statically verify:
 - ❑ `ArrayIndexOutOfBoundsException`
 - ❑ `NullPointerException`
 - ❑ Simple assertions
 - ❑ etc.
- ❑ Now: “*Excelsior FlawDetector is temporarily unavailble due to technical issues.*”
- ❑ However, I had a trial
 - ❑ not compatible with Java 5.0
- ❑ It requires a class with `main`
 - ⇒ Cannot analyze classes in isolation
 - ⇒ Just issues a warning for the `check()`

JML-based/style tools

- ❑ Krakatoa, LOOP, Jack and Spec#
 - ❑ Spec# is as JML but for C#
- ❑ They need annotations
- ❑ In particular the class invariant:

`invariant i <= N`

to prove the class

- ❑ In practice these tools require heavy annotations

Conclusions

- ❑ We presented the theory and some of the abstract operations of our analyzer
- ❑ There is a lot of work to be done:
 - ❑ Implementation details
 - Some optimizations, some syntax transformations, etc.
 - ❑ Modular handling of inheritance (VMCAI'04 paper)
 - ❑ Better handling of mutually recursive classes (AMAST'04 paper)
 - ❑ Does it scales up?
 - eg. Object caching, used in the implementation, does scales up?

Thank you!