

# Analizator statyczny bajtkodu Javy

Michał Ziemba

20 października 2010

# Plan prezentacji

- 1 Wprowadzenie
- 2 Bajtkod
- 3 Analiza statyczna
- 4 Podsumowanie

## O czym jest praca?

- Celem jest napisanie narzędzia, umożliwiającego prostą analizę statyczną kodu wykonywalnego Javy
- Narzędzie będzie zintegrowane ze środowiskiem Umbra
- Analiza będzie obejmować wykrywanie dereferencji wskaźników, które mogą być *NULL*, ...

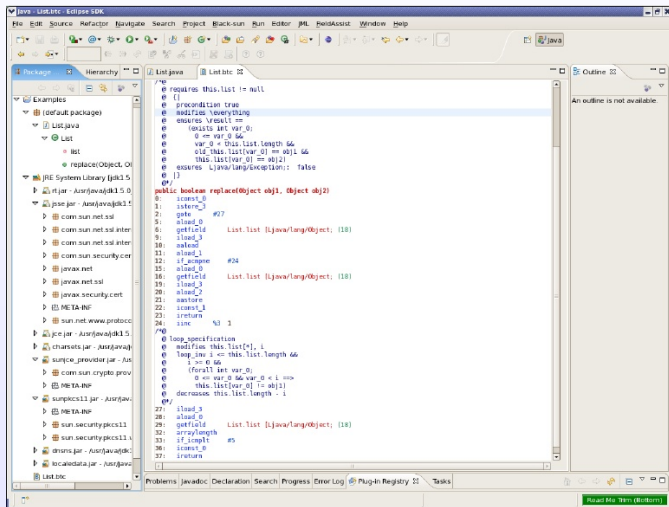
# Motywacja

- Bajtkod składa się z 200 instrukcji
- Możliwa jest faktoryzacja tych instrukcji do 12 grup instrukcji, co znacznie ułatwia analizę
- Nie ma potrzeby zajmowania się kodem źródłowym

# Umbra

- Plugin do Eclipse, umożliwiający generowanie bajtkodu z pliku *.class*
- Jest w fazie rozwojowej
- Umożliwia pisanie specyfikacji w BML

# Umбра



# Bajtkod

```
outer:  
for (int i = 2; i < 1000; i++) {  
  for (int j = 2; j < i; j++) {  
    if (i % j == 0)  
      continue outer;  
  }  
  System.out.println (i);  
}
```

```
0:   iconst 2  
1:   istore 1  
2:   iload 1  
3:   sipush 1000  
6:   if icmpge      44  
9:   iconst 2  
10:  istore 2  
11:  iload 2  
12:  iload 1  
13:  if icmpge      31  
16:  iload 1  
17:  iload_2  
18:  irem  
19:  ifne 25  
22:  goto 38  
25:  iinc 2, 1  
28:  goto 11  
31:  getstatic      #84; //Field java/lang/System.out:Ljava/io/PrintStream;  
34:  iload 1  
35:  invokevirtual #85; //Method java/io/PrintStream.println:(I)V  
38:  iinc 1, 1  
41:  goto 2  
44:  return
```

## Dziedziny semantyczne

Podstawowa forma pojedynczego kroku:

$$P \vdash h, ts \rightarrow h', ts'$$

Struktury semantyczne:

$$\text{Heap} = [\text{Loc} \times \text{Thrld} \rightarrow_{fin} (\text{Cnames} \times \text{Monitor} \times [\text{Fnames} \rightarrow_{fin} \text{Val}])]$$

$$\text{Thread} = \text{Thrld} \times \text{ThreadStatus} \times \text{EvalState} \times \text{FrameStack}$$

$$\text{FrameStack} = \text{MethodFrame}^*$$

$$\text{MethodFrame} = \text{Cnames} \times \text{Mnames} \times \text{LVals} \times \text{OpStack} \times \text{PC}$$



## $load(k, n)$

Instrukcja, która odczytuje lokalne wartości i wkłada je na stos.  
 Typy dzielimy na 32 – bitowe ( $Cat1^\bullet$ ) i 64 – bitowe ( $Cat2$ ).

$$\frac{lv_{ts} = (k, m) \quad ostck' = (k, m) \cdot ostck_{ts} \quad pc' = next(cmthd, pc_{ts}) \\ P@pc_{ts}.mnm_{ts}.cnm_{ts} = load(k, n) \quad k \in Cat1^\bullet \quad est_{ts} = null}{P \vdash h, ts \rightarrow h, ts [ostck \leftarrow ostck'] [pc \leftarrow pc']}$$

$$\frac{lv_{ts}(n) = (k, m_1) \quad lv_{ts}(n+1) = (k, m_2) \\ ostck' = (k, k(m_1, m_2)) \cdot ostck_{ts} \quad pc' = next(cmthd, pc_{ts}) \\ P@pc_{ts}.mnm_{ts}.cnm_{ts} = load(k, n) \quad k \in Cat2 \quad est_{ts} = null}{P \vdash h, ts \rightarrow h, ts [ostck \leftarrow ostck'] [pc \leftarrow pc']}$$

## *store(k, n)*

Instrukcja, która zdejmuje wartość ze stosu i umieszcza ją w tablicy zmiennych lokalnych

$$\frac{
 \begin{array}{l}
 lv' = lv_{ts} [n \leftarrow (k, m)] \\
 ostck_{ts} = (k, m) \cdot ostck' \quad pc' = next(cmthd, pc_{ts}) \\
 P @ pc_{ts}.mnm_{ts}.cnm_{ts} = store(k, n) \quad k \in Cat1 \bullet \quad est_{ts} = null
 \end{array}
 }{
 P \vdash h, ts \rightarrow h, ts [ostck \leftarrow ostck'] [pc \leftarrow pc'] [lv \leftarrow lv']
 }$$

$$\frac{
 \begin{array}{l}
 lv' = lv_{ts} [n \leftarrow (k, m_1)] [n + 1 \leftarrow (half, m_2)] \\
 ostck' = (k, k(m_1, m_2)) \cdot ostck_{ts} \quad pc' = next(cmthd, pc_{ts}) \\
 P @ pc_{ts}.mnm_{ts}.cnm_{ts} = store(k, n) \quad k \in Cat2 \quad est_{ts} = null
 \end{array}
 }{
 P \vdash h, ts \rightarrow h, ts [ostck \leftarrow ostck'] [pc \leftarrow pc']
 }$$

## *stackop*(*op*)

Instrukcja, która używa wyłącznie stosu.  $kinds_{stackop}(op)$  jest zbiorem trójek: lista typów wejściowych, funkcja  $f$  i lista typów wyjściowych. Na przykład

$$kinds_{stackop}(iadd) = \{([int, int], f_{iadd}, [int])\}$$

$$\frac{\begin{array}{l} (l, f, l') \in kinds_{stackop}(op) \quad ostck_{ts} = s \cdot r \\ check(s, l) \quad ostck' = f(s) \cdot r \quad pc' = next(cmthd, pc_{ts}) \\ P @ pc_{ts}.mnm_{ts}.cnm_{ts} = stackop(op) \quad est_{ts} = null \end{array}}{P \vdash h, ts \rightarrow h, ts [ostack \leftarrow ostack'] [pc \leftarrow pc']}$$

## *cond*(*op*, *d*)

Instrukcja, która może wpływać na przebieg sterowania w programie, ale nie modyfikuje stosu ramek (m.in. skoki warunkowe i bezwarunkowe)

Funkcja  $f$  zwraca zmodyfikowany stos oraz nową wartość PC

$$\begin{array}{c}
 (l, f, l') \in \text{kinds}_{\text{cond}}(\text{op}) \quad \text{ostck}_{ts} = s \cdot r \\
 \text{check}(s, l) \quad (s', pc') = f(d, s, \text{next}(\text{cmthd}, pc_{ts})) \quad \text{ostck}' = s' \cdot r \\
 P @ pc_{ts}.mnm_{ts}.cnm_{ts} = \text{cond}(\text{op}, d) \quad \text{est}_{ts} = \text{null} \\
 \hline
 P \vdash h, ts \rightarrow h, ts [\text{ostack} \leftarrow \text{ostack}'] [pc \leftarrow pc']
 \end{array}$$

## *iinc*(*n*, *c*)

Tylko jedna instrukcja(*iinc*), używająca wyłącznie tablicy zmiennych lokalnych

$$\begin{array}{l}
 lv_{ts}(n) = (int, m) \quad lv' = lv_{ts} [n \leftarrow (int, m +_{int} c)] \\
 pc' = next(cmthd, ps_{ts}) \\
 P @ pc_{ts}.mnm_{ts}.cnm_{ts} = iinc(n, c) \quad est_{ts} = null \\
 \hline
 P \vdash h, ts \rightarrow h, ts [lv \leftarrow lv'] [pc \leftarrow pc']
 \end{array}$$

# $get(op, d)$

Instrukcja odczytująca wartość z *Heap* i wstawiająca ją na *OpStack*

$$\begin{array}{l}
 (l, f, l') = kinds_{get}(op, d) \quad ostck_{ts} = s \cdot r \quad check(s, l) \\
 s' = f(d, s, h) \quad s' \in OpStack \quad ostck' = s' \cdot r \quad est_{ts} = null \\
 pc' = next(cmthd, ps_{ts}) \quad P@pc_{ts}.mnm_{ts}.cnm_{ts} = get(op, d) \\
 \hline
 P \vdash h, ts \rightarrow h, ts [lv \leftarrow lv'] [pc \leftarrow pc']
 \end{array}$$

$$\begin{array}{l}
 (l, f, l') = kinds_{get}(op, d) \quad ostck_{ts} = s \cdot r \quad check(s, l) \\
 e = f(d, s, g) \quad e \in Loc^\bullet \quad P@pc_{ts}.mnm_{ts}.cnm_{ts} = iinc(n, c) \\
 est_{ts} = null \\
 \hline
 P \vdash h, ts \rightarrow h, ts [est \leftarrow e]
 \end{array}$$

# *put*(*op*, *d*)

Instrukcja odczytująca wartość z *OpStack* i wstawiająca ją na *Heap*

$$\begin{array}{l}
 (l, f, l') = \text{kinds}_{\text{put}}(\text{op}, d) \quad \text{ostck}_{ts} = s \cdot r \quad \text{check}(s, l) \quad \text{ostck}' = r \\
 h' = f(d, s, h) \quad h' \in \text{Heap} \quad \text{pc}' = \text{next}(\text{cmthd}, \text{ps}_{ts}) \\
 P @ \text{pc}_{ts} . \text{mnm}_{ts} . \text{cnm}_{ts} = \text{put}(\text{op}, d) \quad \text{est}_{ts} = \text{null} \\
 \hline
 P \vdash h, ts \rightarrow h', ts [\text{ostck} \leftarrow \text{ostck}'] [\text{pc} \leftarrow \text{pc}']
 \end{array}$$

$$\begin{array}{l}
 (l, f, l') = \text{kinds}_{\text{put}}(\text{op}, d) \quad \text{ostck}_{ts} = s \cdot r \quad \text{check}(s, l) \quad \text{ostck}' = r \\
 e = f(d, s, h) \quad e \in \text{Loc}^\bullet \\
 P \text{ pc}_{ts} . \text{mnm}_{ts} . \text{cnm}_{ts} = \text{put}(\text{op}, d) \quad \text{est}_{ts} = \text{null} \\
 \hline
 P \vdash h, ts \rightarrow h, ts [\text{est} \leftarrow e]
 \end{array}$$

## $new(op, d)$

Instrukcja przydzielająca pamięć tworzonej strukturze. W odróżnieniu od  $put$ ,  $new$  tworzy nowe lokacje na stosie, zamiast modyfikować już istniejące

$$\begin{array}{l}
 (l, f, l') \in kinds_{new}(op, d) \quad ostck_{ts} = s \cdot r \quad check(s, l) \\
 pc' = next(cmthd, ps_{ts}) \quad ostck' = s' \cdot r \\
 (s', h') = f(d, s, h) \quad s' \in OpStack \quad h' \in Heap \\
 P@pc_{ts}.mnm_{ts}.cnm_{ts} = new(op, d) \quad est_{ts} = null \\
 \hline
 P \vdash h, ts \rightarrow h', ts [lv \leftarrow lv'] [pc \leftarrow pc']
 \end{array}$$

$$\begin{array}{l}
 (l, f, l') \in kinds_{new}(op, d) \quad ostck_{ts} = s \cdot r \quad check(s, l) \\
 e = f(d, s, h) \quad e \in Loc^\bullet \quad P@pc_{ts}.mnm_{ts}.cnm_{ts} = new(op, d) \\
 est_{ts} = null \\
 \hline
 P \vdash h, ts \rightarrow h, ts [est \leftarrow e]
 \end{array}$$



## *monitor*(*op*)

Instrukcja modyfikująca stan wątków poprzez zajęcie (*op* = *enter*)  
 lub zwolnienie (*op* = *exit*) monitora

$$\begin{array}{c}
 (l, f, l') \in \text{kinds}_{\text{new}}(\text{op}, d) \quad \text{ostck}_{ts} = s \cdot r \quad \text{check}(s, l) \\
 \text{ostck}' = s' \cdot r \\
 (s', h') = f(d, s, h) \quad s' \in \text{OpStack} \quad h' \in \text{Heap} \\
 pc' = \text{next}(\text{cmthd}, ps_{ts}) \quad P@pc_{ts}.\text{mnm}_{ts}.\text{cnm}_{ts} = \text{new}(\text{op}, d) \\
 \text{est}_{ts} = \text{null} \\
 \hline
 P \vdash h, ts \rightarrow h', ts [lv \leftarrow lv'] [pc \leftarrow pc']
 \end{array}$$

## *monitor(op) cd.*

$$\frac{\begin{array}{l} (l, f, l') = \text{kinds}_{\text{new}}(\text{op}, d) \quad \text{ostck}_{ts} = s \cdot r \quad \text{check}(s, l) \\ e = f(d, s, h) \quad e \in \text{Loc}^\bullet \quad P @ pc_{ts}.mnm_{ts}.cnm_{ts} = \text{new}(\text{op}, d) \\ \text{est}_{ts} = \text{null} \end{array}}{P \vdash h, ts \rightarrow h, ts [\text{est} \leftarrow e]}$$

## *invoke(mode, cnm, mnm)*

- Instrukcja modyfikująca stos i stos ramek
- *mode* może być jednym z *interface*, *special*, *static* lub *virtual*
- Funkcja *dispatch* między innymi sprawdza czy flagi metody nie są sprzeczne z *mode* oraz czy prawa dostępu są zachowane (*private*, *protected*)
- Funkcja *initlv* umieszcza wartości ze stosu w tablicy zmiennych lokalnych

## *invoke(mode, cnm, mnm) cd.*

$$\begin{array}{l}
 (l, l') = \text{kindsof}_{\text{invoke}}(\text{mode}, \text{cnm}, \text{mnm}) \quad \text{ostck}_{ts} = s \cdot r \quad \text{check}(s, l) \\
 \text{cmn}' = \text{dispatch}(\text{mode}, \text{cnm}, \text{mnm}, s, h) \quad \text{tfs}' = \text{tfs}_{ts} [\text{ostck} \leftarrow r] \\
 \text{lv}' = \text{initlv}(\text{lvlength}(P@\text{mnm}.\text{cnm}), s) \\
 \text{tfs}'' = \langle \text{cmn}', \text{mnm}, \text{lv}', [ ], 0 \rangle \cdot \text{tfs}' \\
 P@\text{pc}_{ts}.\text{mnm}_{ts}.\text{cnm}_{ts} = \text{invoke}(\text{mode}, \text{cnm}, \text{mnm}) \quad \text{est}_{ts} = \text{null} \\
 \hline
 P \vdash h, ts \rightarrow h, ts [\text{tfs} \leftarrow \text{tfs}'']
 \end{array}$$

$$\begin{array}{l}
 (l, l') = \text{kindsof}_{\text{invoke}}(\text{mode}, \text{cnm}, \text{mnm}) \quad \text{ostck}_{ts} = s \cdot r \quad \text{check}(s, l) \\
 e = \text{dispatch}(\text{mode}, \text{cnm}, \text{mnm}, s, h) \quad e \in \text{Loc} \\
 P@\text{pc}_{ts}.\text{mnm}_{ts}.\text{cnm}_{ts} = \text{invoke}(\text{mode}, \text{cnm}, \text{mnm}) \quad \text{est}_{ts} = \text{null} \\
 \hline
 P \vdash h, ts \rightarrow h, ts [\text{est} \leftarrow e]
 \end{array}$$

# *return(l)*

Instrukcja usuwająca aktualną ramkę ze stosu oraz modyfikująca poprzednią ramkę

$$\begin{array}{c}
 ostck_{ts} = s \cdot r \quad check(s, l) \\
 tfs_{ts} = f_1 \cdot \langle cnm', mnm', lv', ostck', pc' \rangle \cdot tfs^{tail} \quad f_1 \in MethodFrame \\
 tfs' = \langle cnm', mnm', lv', vsc(s) \cdot ostck', next(P@mm'.cnm', pc') \rangle \cdot tfs^{tail} \\
 P@pc_{ts}.mnm_{ts}.cnm_{ts} = return(l) \quad est_{ts} = null \\
 \hline
 P \vdash h, ts \rightarrow h, ts [tfs \leftarrow tfs']
 \end{array}$$

## *return(l) cd.*

$$\frac{\begin{array}{l} \text{ostck}_{ts} = s \cdot r \quad \text{check}(s, l) \quad \text{tfs}_{ts} = [f] \quad f \in \text{MethodFrame} \\ \text{tfs}' = [] \quad \text{tstatus}' = \text{TERMINATED} \\ P@pc_{ts}.mnm_{ts}.cnm_{ts} = \text{return}(l) \quad \text{est}_{ts} = \text{null} \end{array}}{P \vdash h, ts \rightarrow h, ts [\text{tfs} \leftarrow \text{tfs}'] [\text{tstatus} \leftarrow \text{tstatus}']}$$

## *throw*

Instrukcja nie mająca argumentów, odczytująca i usuwająca lokację wyjątku ze stosu i zmieniająca stan wykonania aktualnego wątku  
 Do kolejnego stanu można przejść albo wykonując instrukcję `throw`, albo rzucając wyjątek

$$\frac{ostck_{ts} = e \cdot r \quad e \in Loc^\bullet \quad P@pc_{ts}.mnm_{ts}.cnm_{ts} = throw \quad est_{ts} = null}{P \vdash h, ts \rightarrow h, ts [est \leftarrow e]}$$

$$\frac{ostck' = [e] \quad (pc_{ts}, h(e)@cnm) \in dom(P@etable.nmn_{ts}.cnm_{ts}) \quad pc' = P@etable.mnm_{ts}.cnmts(pc_{ts}, h(e)@cnm) \quad est_{ts} = e \in Loc^\bullet}{P \vdash h, ts \rightarrow h, ts [ostck \leftarrow ostck'] [pc \leftarrow pc'] [est \leftarrow null]}$$

*throw cd.*

$$\frac{(pc_{ts}, h(e)@cnm) \notin \text{dom}(P@etable.nmn_{ts}.cnm_{ts}) \quad tfs_{ts} = f_1 \cdot f_2 \cdot tfs_{ts}^{tail} \quad f_1, f_2 \in \text{MethodFrame} \quad est_{ts} = e \in \text{Loc}^\bullet}{P \vdash h, ts \rightarrow h, ts \left[ tfs \leftarrow f_2 \cdot tfs_{ts}^{tail} \right]}$$

$$\frac{(pc_{ts}, h(e)@cnm) \notin \text{dom}(P@etable.nmn_{ts}.cnm_{ts}) \quad tfs_{ts} = [f] \quad f \in \text{MethodFrame} \quad est_{ts} = e \in \text{Loc}^\bullet}{P \vdash h, ts \rightarrow h, ts \left[ tfs \leftarrow [ ] \right] \left[ tstatus \leftarrow \text{TERMINATED} \right]}$$



# NonNull

- Niektóre języki obiektowe i funkcyjne nie pozwalają odwoływać się do obiektów w czasie ich tworzenia, jednak w Javie nie ma tego ograniczenia
- Potrzebujemy anotacji opisujących obiekty, które mogą być NULL, a także obiekty które mogą być częściowo zainicjalizowane
- Dzięki takim anotacjom w kodzie można przeprowadzić podstawową analizę

```
class A {
    /*@ non_null */ string name;

    public A (/*@ non_null */ string s) {
        this.name = s;
        this.m(55);
    }

    virtual void m(int x) {...}
}

class B extends A {
    /*@ non_null */ string path;

    public B (/*@ non_null */ string p; /*@ non_null */ string s) {
        super();
        this.path = p;
    }

    override void m(int x) {
        ...
        this.path
        ...
    }
}
```

- Analizator statyczny bajtkodu Javy może być przydatnym narzędziem, gdy chcemy wykonać proste sprawdzenie kodu bez oglądania i parsowania kodu źródłowego
- Dzięki uproszczeniu bajtkodu do 12 grup instrukcji analiza będzie łatwiejsza
- Początkowo przewidziana jest tylko analiza w zakresie dpreferencji NULLi, ale będzie możliwość rozszerzenia narzędzia o inne sprawdzenia

## Bibliografia

- 1 J. Chrzęszcz, P. Czarnik, A. Schubert, *A dozen instructions make Java bytecode*
- 2 M. Fähndrich, K. R. M. Leino, *Declaring and Checking Non-null Types in an Object-Oriented Language*
- 3 <http://www.mimuw.edu.pl/~alx/umbra/bml-bytecode.pdf>
- 4 [http://java.sun.com/docs/books/jvms/second\\_edition/html/Overview.doc.html](http://java.sun.com/docs/books/jvms/second_edition/html/Overview.doc.html)