# Extracting a Data Flow Analyser in Constructive Logic

David Cachera, Thomas Jensen, David Pichardie and Vlad Rusu

APPSEM'04, Tallinn

IRISA
institut de recherche en informatique
et systèmes aléatoires

# Static program analysis

The goals of static program analysis

- ▶ To prove properties about the run-time behaviour of a program
- ▶ In a fully automatic way
- ▶ Without actually executing this program
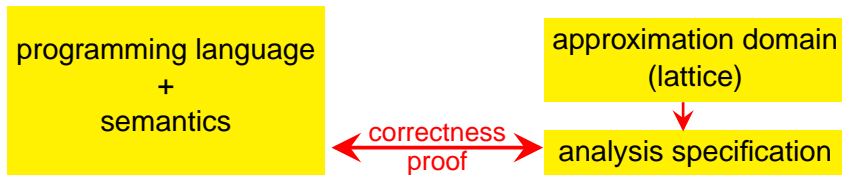
# Static program analysis

The goals of static program analysis

- ▶ To prove properties about the run-time behaviour of a program
- ▶ In a fully automatic way
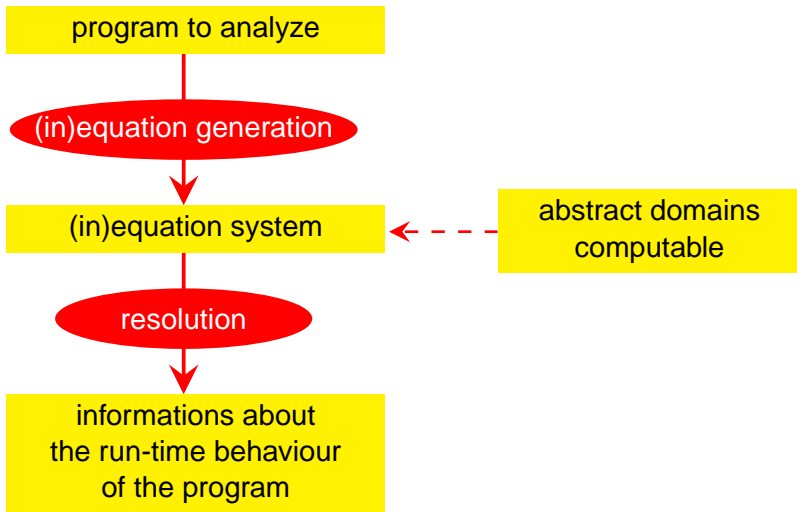- ▶ Without actually executing this program

Solid foundations for designing an analyser

- ▶ Formalization and correctness proof by abstract interpretation
- ▶ Resolution of constraints on lattices by iteration and symbolic computation

# Formalization

# Resolution

So what's the problem ?

# Formalization part

$\ddot{\alpha}\llbracket P \rrbracket(\text{Post}\llbracket \textbf{if } B \textbf{ then } S_t \textbf{ else } S_f \textbf{ fi}\rrbracket)$

$=$    $\langle$def. (110) of $\ddot{\alpha}\llbracket P \rrbracket\rangle$

     $\ddot{\alpha}\llbracket P \rrbracket \circ \text{Post}\llbracket \textbf{if } B \textbf{ then } S_t \textbf{ else } S_f \textbf{ fi}\rrbracket \circ \ddot{\gamma}\llbracket P \rrbracket$

$=$    $\langle$def. (103) of Post$\rangle$

     $\ddot{\alpha}\llbracket P \rrbracket \circ \text{post}[\tau^\bullet \llbracket \textbf{if } B \textbf{ then } S_t \textbf{ else } S_f \textbf{ fi}\rrbracket] \circ \ddot{\gamma}\llbracket P \rrbracket$

$=$    $\langle$big step operational semantics (93)$\rangle$

     $\ddot{\alpha}\llbracket P \rrbracket \circ \text{post}\,[(1_{\Sigma\llbracket P \rrbracket} \cup \tau^{\mathcal{B}}) \circ \tau^\bullet \llbracket S_t \rrbracket \cup (1_{\Sigma\llbracket P \rrbracket} \cup \tau^{\mathcal{B}}) \circ \tau^\bullet \llbracket S_f \rrbracket \circ (1_{\Sigma\llbracket P \rrbracket} \cup \tau^{f})] \circ \ddot{\gamma}\llbracket P \rrbracket$

$=$    $\langle$Galois connection (98) so that post preserves joins$\rangle$

     $\ddot{\alpha}\llbracket P \rrbracket \quad \circ \quad (\text{post}\,[(1_{\Sigma\llbracket P \rrbracket} \cup \tau^{\mathcal{B}}) \circ \tau^\bullet \llbracket S_t \rrbracket \quad \dot{\cup} \quad (1_{\Sigma\llbracket P \rrbracket} \cup \tau^{t})] \quad \dot{\cup}$   $\text{post}[(1_{\Sigma\llbracket P \rrbracket} \cup \tau^{\mathcal{B}}) \circ \tau^\bullet \llbracket S_f \rrbracket \circ (1_{\Sigma\llbracket P \rrbracket} \cup \tau^{f})]) \circ \ddot{\gamma}\llbracket P \rrbracket$

$=$    $\langle$Galois connection (106) so that $\ddot{\alpha}\llbracket P \rrbracket$ preserves joins$\rangle$

     $(\ddot{\alpha}\llbracket P \rrbracket \quad \circ \quad \text{post}[(1_{\Sigma\llbracket P \rrbracket} \cup \tau^{\mathcal{B}}) \circ \tau^\bullet \llbracket S_t \rrbracket \circ (1_{\Sigma\llbracket P \rrbracket} \cup \tau^{t}) \quad \circ \quad \ddot{\gamma}\llbracket P \rrbracket) \quad \dot{\dot{\cup}} \quad (\ddot{\alpha}\llbracket P \rrbracket \quad \circ$   $\text{post}[(1_{\Sigma\llbracket P \rrbracket} \cup \tau^{\mathcal{B}}) \circ \tau^\bullet \llbracket S_f \rrbracket \circ (1_{\Sigma\llbracket P \rrbracket} \cup \tau^{f})] \circ \ddot{\gamma}\llbracket P \rrbracket)$

$\stackrel{\cdot}{=}$    $\langle$lemma (5.3) and similar one for the **else** branch$\rangle$

     $\lambda J \cdot \text{let } J^{t'} = \lambda l \in \text{in}_P\llbracket P \rrbracket \cdot (l = \text{at}_P\llbracket S_t \rrbracket \; ? \; J_{\text{at}_P\llbracket S_t \rrbracket} \; \dot{\cup} \; \text{Abexp}\llbracket B \rrbracket (J_\ell) \; \natural \; J_l) \text{ in}$      (120)

        $\text{let } J^{t''} = \text{APost}\llbracket S_t \rrbracket (J^{t'}) \text{ in}$

           $\lambda l \in \text{in}_P\llbracket P \rrbracket \cdot (l = \ell' \; ? \; J^{t''}_{\ell'} \; \dot{\cup} \; J^{t''}_{\text{after}_P\llbracket S_t \rrbracket} \; \natural \; J_l^{t''})$

     $\ddot{\cup}$

        $\text{let } J^{f'} = \lambda l \in \text{in}_P\llbracket P \rrbracket \cdot (l = \text{at}_P\llbracket S_f \rrbracket \; ? \; J_{\text{at}_P\llbracket S_f \rrbracket} \; \dot{\cup} \; \text{Abexp}\llbracket T(\neg B) \rrbracket (J_\ell) \; \natural \; J_l) \text{ in}$

        $\text{let } J^{f''} = \text{APost}\llbracket S_f \rrbracket (J^{f'}) \text{ in}$

           $\lambda l \in \text{in}_P\llbracket P \rrbracket \cdot (l = \ell' \; ? \; J^{f''}_{\ell'} \; \dot{\cup} \; J^{f''}_{\text{after}_P\llbracket S_f \rrbracket} \; \natural \; J_l^{f''})$

$=$    $\langle$by grouping similar terms$\rangle$

     $\lambda J \cdot \text{let } J^{t'} = \lambda l \in \text{in}_P\llbracket P \rrbracket \cdot (l = \text{at}_P\llbracket S_t \rrbracket \; ? \; J_{\text{at}_P\llbracket S_t \rrbracket} \; \dot{\cup} \; \text{Abexp}\llbracket B \rrbracket (J_\ell) \; \natural \; J_l)$

     $\text{and } J^{f'} = \lambda l \in \text{in}_P\llbracket P \rrbracket \cdot (l = \text{at}_P\llbracket S_f \rrbracket \; ? \; J_{\text{at}_P\llbracket S_f \rrbracket} \; \dot{\cup} \; \text{Abexp}\llbracket T(\neg B) \rrbracket (J_\ell) \; \natural \; J_l) \text{ in}$

        $\text{let } J^{t''} = \text{APost}\llbracket S_t \rrbracket (J^{t'})$

     $\text{and } J^{f''} = \text{APost}\llbracket S_f \rrbracket (J^{f'}) \text{ in}$

        $\lambda l \in \text{in}_P\llbracket P \rrbracket \cdot (l = \ell' \; ? \; J^{t''}_{\ell'} \; \dot{\cup} \; J^{t''}_{\text{after}_P\llbracket S_t \rrbracket} \; \dot{\cup} \; J^{f''}_{\ell'} \; \dot{\cup} \; J^{f''}_{\text{after}_P\llbracket S_f \rrbracket} \; \natural \; J_l^{t''} \; \dot{\cup} \; J_l^{f''})$

$=$    $\langle$by locality (113) and labelling scheme (59) so that in particular $J^{t''}_{\ell'} = J^{t''}_{\ell'} = J^{f}_{\ell'} = J^{f}_{\ell'}$

     $= J^{t''}_{\ell'} = J^{f''}_{\ell'}$ and $\text{APost}\llbracket S_t \rrbracket$ and $\text{APost}\llbracket S_f \rrbracket$ do not interfere$\rangle$

# Formalization part

$$\ddot{a}[\![P]\!](\mathrm{Post}[\![\mathbf{if}\ B\ \mathbf{then}\ S_t\ \mathbf{else}\ S_f\ \mathbf{fi}]\!])$$
$$=\quad \big\{\text{def. (110) of }\ddot{a}[\![P]\!]\big\}$$
$$\ddot{a}[\![P]\!]\circ\mathrm{Post}[\![\mathbf{if}\ B\ \mathbf{then}\ S_t\ \mathbf{else}\ S_f\ \mathbf{fi}]\!]\circ\ddot{\gamma}[\![P]\!]$$
$$=\quad \big\{\text{def. (103) of Post}\big\}$$
$$\ddot{a}[\![P]\!]\circ\mathrm{post}[\![\tau^\star[\![\mathbf{if}\ B\ \mathbf{then}\ S_t\ \mathbf{else}\ S_f\ \mathbf{fi}]\!]]\!]\circ\ddot{\gamma}[\![P]\!]$$
$$=\quad \big\{\text{big step operational semantics (93)}\big\}$$
$$\ddot{a}[\![P]\!]\circ\mathrm{post}[\![(1_{\Sigma[\![P]\!]}\cup\tau^t)\circ(1_{\Sigma[\![P]\!]}\cup\tau^{\neg B})\circ\tau^\star[\![S_f]\!]\circ(1_{\Sigma[\![P]\!]}\cup\tau^t)]\!]\circ\ddot{\gamma}[\![P]\!]$$
$$=\quad \big\{\text{Galois connection (98) so that post preserves joins}\big\}$$
$$\ddot{a}[\![P]\!]\quad\circ\quad(\mathrm{post}[\![(1_{\Sigma[\![P]\!]}\cup\tau^{\neg B})\circ\tau^\star[\![S_t]\!]]\!]\quad\circ\quad(1_{\Sigma[\![P]\!]}\cup\tau^t))\quad\dot{\cup}$$
$$\mathrm{post}[\![(1_{\Sigma[\![P]\!]}\cup\tau^{\neg B})\circ\tau^\star[\![S_f]\!]\circ(1_{\Sigma[\![P]\!]}\cup\tau^t)]\!])\circ\ddot{\gamma}[\![P]\!]$$
$$=\quad \big\{\text{Galois connection (106) so that }\ddot{a}[\![P]\!]\text{ preserves joins}\big\}$$
$$(\ddot{a}[\![P]\!]\quad\circ\quad\mathrm{post}[\![(1_{\Sigma[\![P]\!]}\cup\tau^{\neg B})\circ\tau^\star[\![S_t]\!]\circ(1_{\Sigma[\![P]\!]}\cup\tau^t)]\!]\quad\circ\quad\ddot{\gamma}[\![P]\!])\quad\dot{\dot{\cup}}\quad(\ddot{a}[\![P]\!]\quad\circ$$
$$\mathrm{post}[\![(1_{\Sigma[\![P]\!]}\cup\tau^{\neg B})\circ\tau^\star[\![S_f]\!]\circ(1_{\Sigma[\![P]\!]}\cup\tau^t)]\!]\circ\ddot{\gamma}[\![P]\!])$$
$$\stackrel{\cdot}{\sqsubseteq}\quad \big\{\text{lemma (5.3) and similar one for the \textbf{else} branch}\big\}$$
$$\lambda\,J\cdot\mathrm{let}\ J^{t'}=\lambda l\in\mathrm{in}_P[\![P]\!]\cdot(l=\mathrm{at}_P[\![S_t]\!]\ ?\ J_{\mathrm{at}_P[\![S_t]\!]}\ \dot{\cup}\ \mathrm{Abexp}[\![B]\!](J_\ell)\ \dot{\iota}\ J_l)\ \mathrm{in}\qquad\qquad(120)$$
$$\mathrm{let}\ J^{t''}=\mathrm{APost}[\![S_t]\!](J^{t'})\ \mathrm{in}$$
$$\lambda l\in\mathrm{in}_P[\![P]\!]\cdot(l=\ell'\ ?\ J_\ell^{t''}\ \dot{\cup}\ J_{\mathrm{after}^t[\![S_t]\!]}^{t''}\ \dot{\iota}\ J_l^{t''})$$
$$\ddot{\cup}$$
$$\mathrm{let}\ J^{f'}=\lambda l\in\mathrm{in}_P[\![P]\!]\cdot(l=\mathrm{at}_P[\![S_f]\!]\ ?\ J_{\mathrm{at}_P[\![S_f]\!]}\ \dot{\cup}\ \mathrm{Abexp}[\![T(\neg B)]\!](J_\ell)\ \dot{\iota}\ J_l)\ \mathrm{in}$$
$$\mathrm{let}\ J^{f''}=\mathrm{APost}[\![S_f]\!](J^{f'})\ \mathrm{in}$$
$$\lambda l\in\mathrm{in}_P[\![P]\!]\cdot(l=\ell'\ ?\ J_\ell^{f''}\ \dot{\cup}\ J_{\mathrm{after}^f[\![S_f]\!]}^{f''}\ \dot{\iota}\ J_l^{f''})$$
$$=\quad \big\{\text{by grouping similar terms}\big\}$$
$$\lambda\,J\cdot\mathrm{let}\ J^{t'}=\lambda l\in\mathrm{in}_P[\![P]\!]\cdot(l=\mathrm{at}_P[\![S_t]\!]\ ?\ J_{\mathrm{at}_P[\![S_t]\!]}\ \dot{\cup}\ \mathrm{Abexp}[\![B]\!](J_\ell)\ \dot{\iota}\ J_l)$$
$$\mathrm{and}\ J^{f'}=\lambda l\in\mathrm{in}_P[\![P]\!]\cdot(l=\mathrm{at}_P[\![S_f]\!]\ ?\ J_{\mathrm{at}_P[\![S_f]\!]}\ \dot{\cup}\ \mathrm{Abexp}[\![T(\neg B)]\!](J_\ell)\ \dot{\iota}\ J_l)\ \mathrm{in}$$
$$\mathrm{let}\ J^{t''}=\mathrm{APost}[\![S_t]\!](J^{t'})$$
$$\mathrm{and}\ J^{f''}=\mathrm{APost}[\![S_f]\!](J^{f'})\ \mathrm{in}$$
$$\lambda l\in\mathrm{in}_P[\![P]\!]\cdot(l=\ell'\ ?\ J_\ell^{t''}\ \dot{\cup}\ J_{\mathrm{after}^t[\![S_t]\!]}^{t''}\ \dot{\cup}\ J_\ell^{f''}\ \dot{\cup}\ J_{\mathrm{after}^f[\![S_f]\!]}^{f''}\ \dot{\iota}\ J_l^{t''}\ \dot{\cup}\ J_l^{f''})$$
$$=\quad \big\{\text{by locality (113) and labelling scheme (59) so that in particular }J_{\ell'}^{t''}=J_\ell^{t''}=J_\ell^{f''}=J_\ell^f$$
$$=J_{\ell'}^{f''}=J_\ell^{f''}\text{ and }\mathrm{APost}[\![S_t]\!]\text{ and }\mathrm{APost}[\![S_f]\!]\text{ do not interfere}\big\}$$

# Implementation part

```c
int main(int argc, char **argv)
{
    int i, j, t, parent_a, parent_b;
    int **swap, **newpop, **oldpop;
    double *fit, *normfit;

    get_options(argc, argv, options, help_string);
    srandom(seed);
    read_specs(specs);
    size += (size / 2 * 2 != size);
    newpop = xmalloc(sizeof(int *) * size);
    oldpop = xmalloc(sizeof(int *) * size);
    fit = xmalloc(sizeof(double) * size);
    normfit = xmalloc(sizeof(double) * size);
    for(i = 0; i < size; i++) {
        newpop[i] = xmalloc(sizeof(int) * len);
        oldpop[i] = xmalloc(sizeof(int) * len);
        for(j = 0; j < len * 2; j++)
            random_solution(oldpop[i]);
    }
    for(t = 0; t < gens; t++) {
        compute_fitness(oldpop, fit, normfit);
        dump_stats(t, oldpop, fit);
        for(i = 0; i < size; i += 2) {
            parent_a = select_one(normfit);
            parent_b = select_one(normfit);
            reproduce(oldpop, newpop, parent_a, parent_b, i);
        }
        swap = newpop; newpop = oldpop; oldpop = swap;
    }
    exit(0);
}
```

$\ddot{\alpha}\llbracket F \rrbracket(\mathrm{Post}\llbracket \mathbf{if}\ B\ \mathbf{then}\ S_t\ \mathbf{else}\ S_f\ \mathbf{fi}\rrbracket)$

$=$　$\{$def. (110) of $\ddot{\alpha}\llbracket F \rrbracket\}$

　$\ddot{\alpha}\llbracket F \rrbracket \circ \mathrm{Post}\llbracket \mathbf{if}\ B\ \mathbf{then}\ S_t\ \mathbf{else}\ S_f\ \mathbf{fi}\rrbracket \circ \ddot{\gamma}\llbracket F \rrbracket$

$=$　$\{$def. (103) of Post$\}$

　$\ddot{\alpha}\llbracket F \rrbracket \circ \mathrm{post}\lceil \tau^\star \llbracket \mathbf{if}\ B\ \mathbf{then}\ S_t\ \mathbf{else}\ S_f\ \mathbf{fi}\rrbracket \circ \ddot{\gamma}\llbracket F \rrbracket$

$=$　$\{$big step operational semantics (93)$\}$

　$\ddot{\alpha}\llbracket F \rrbracket \circ \mathrm{post}\lceil (1_{\Sigma\llbracket F\rrbracket} \cup \tau^B) \circ \tau^\star \llbracket S_t\rrbracket \circ (1_{\Sigma\llbracket F\rrbracket} \cup \tau^B) \circ \tau^\star \llbracket S_f\rrbracket \circ (1_{\Sigma\llbracket F\rrbracket} \cup \tau^f) \rceil \circ \ddot{\gamma}\llbracket F \rrbracket$

$=$　$\{$Galois connection (98) so that post preserves joins$\}$

　$\ddot{\alpha}\llbracket F \rrbracket \quad \circ \quad (\mathrm{post}\lceil (1_{\Sigma\llbracket F\rrbracket} \cup \tau^B) \quad \circ \quad \tau^\star \llbracket S_t\rrbracket \quad (1_{\Sigma\llbracket F\rrbracket} \cup \tau^f)\rceil \quad \dot{\cup}$
　$\mathrm{post}\lceil (1 \to \tau^B) \circ \tau^\star \llbracket S_f\rrbracket \circ (1 \cup \tau^f)\rceil \circ \ddot{\gamma}\llbracket F \rrbracket$



$\lambda J \cdot \text{let } J^{t'} = \lambda J \in \mathrm{in}_F\llbracket F \rrbracket \cdot (l = \mathrm{at}_F\llbracket S_t\rrbracket\ ?\ J_{\mathrm{at}_F\llbracket S_t\rrbracket} \dot{\cup} \mathrm{Abexp}\llbracket B \rrbracket (J_\ell) \dot{\cup} J_l) \text{ in} \qquad (120)$
　　let $J^{t''} = \mathrm{APost}\llbracket S_t\rrbracket (J^{t'})$ in
　　$\lambda J \in \mathrm{in}_F\llbracket F \rrbracket \cdot (l = \ell'\ ?\ J^{t''}_{\ell'} \dot{\cup} J^{t''}_{\mathrm{after}_F\llbracket S_t\rrbracket} \dot{\cup} J^{t''}_l)$
　　$\dot{\cup}$
　　let $J^{f'} = \lambda J = \lambda J \in \mathrm{in}_F\llbracket F \rrbracket \cdot (l = \mathrm{at}_F\llbracket S_f\rrbracket\ ?\ J_{\mathrm{at}_F\llbracket S_f\rrbracket} \dot{\cup} \mathrm{Abexp}\llbracket T(\neg B)\rrbracket (J_\ell) \dot{\cup} J_l) \text{ in}$
　　　let $J^{f''} = \mathrm{APost}\llbracket S_f\rrbracket (J^{f'})$ in
　　　$\lambda J \in \mathrm{in}_F\llbracket F \rrbracket \cdot (l = \ell'\ ?\ J^{f''}_{\ell'} \dot{\cup} J^{f''}_{\mathrm{after}_F\llbracket S_f\rrbracket} \dot{\cup} J^{f''}_l)$

$=$　$\{$by grouping similar terms$\}$

$\lambda J \cdot \text{let } J^{t'} = \lambda J \in \mathrm{in}_F\llbracket F \rrbracket \cdot (l = \mathrm{at}_F\llbracket S_t\rrbracket\ ?\ J_{\mathrm{at}_F\llbracket S_t\rrbracket} \dot{\cup} \mathrm{Abexp}\llbracket B \rrbracket (J_\ell) \dot{\cup} J_l)$
　　and $J^{t''} = \lambda J \in \mathrm{in}_F\llbracket F \rrbracket \cdot (l = \mathrm{at}_F\llbracket S_f\rrbracket\ ?\ J_{\mathrm{at}_F\llbracket S_f\rrbracket} \dot{\cup} \mathrm{Abexp}\llbracket T(\neg B)\rrbracket (J_\ell) \dot{\cup} J_l) \text{ in}$
　　let $J^{t''} = \mathrm{APost}\llbracket S_t\rrbracket (J^{t'})$
　　and $J^{f''} = \mathrm{APost}\llbracket S_f\rrbracket (J^{f'})$ in
　　$\lambda J \in \mathrm{in}_F\llbracket F \rrbracket \cdot (l = \ell'\ ?\ J^{t''}_{\ell'} \dot{\cup} J^{t''}_{\mathrm{after}_F\llbracket S_t\rrbracket} \dot{\cup} J^{f''}_{\ell'} \dot{\cup} J^{f''}_{\mathrm{after}_F\llbracket S_f\rrbracket} \dot{\cup} J^{t''}_l \dot{\cup} J^{f''}_l)$

$=$　$\{$by locality (113) and labelling scheme (59) so that in particular $J^{t''}_{\ell'} = J^{t''}_{\ell'}, J^{t''}_{\ell'} = J^f_{\ell'}$
　　$= J^{f''}_{\ell'} = J^{f''}_{\ell'}$ and $\mathrm{APost}\llbracket S_t\rrbracket$ and $\mathrm{APost}\llbracket S_f\rrbracket$ do not interfere$\}$

©P. Cousot

```c
int main(int argc, char **argv)
{
  int i, j, t, parent_a, parent_b;
  int **swap, **newpop, **oldpop;
  double *fit, *normfit;

  get_options(argc, argv, options, help_string);
  srandom(seed);
  read_specs(specs);
  size += (size / 2 * 2 != size);
  newpop = xmalloc(sizeof(int *) * size);
```

## Do both parts talk about the same ?

```c
    oldpop[i] = xmalloc(sizeof(int *) * len);
    for(j = 0; j < len * 2; j++)
      random_solution(oldpop[i]);
  }
  for(t = 0; t < gens; t++) {
    compute_fitness(oldpop, fit, normfit);
    dump_stats(t, oldpop, fit);
    for(i = 0; i < size; i += 2) {
      parent_a = select_one(normfit);
      parent_b = select_one(normfit);
      reproduce(oldpop, newpop, parent_a, parent_b, i);
    }
    swap = newpop; newpop = oldpop; oldpop = swap;
  }
  exit(0);
}
```

# Static Analysis for real-life languages

Example of real-life language : bytecode JavaCard

- ▶ 180 instructions
- ▶ Real need of static analysis to verify properties about security, memory management, ...

For this kind of languages,

- ▶ Abstract domains can be complex
- ▶ Correctness proofs become long and tiresome
- ▶ Implementation and maintenance of the analyser become a software engineering task

# In this work

We propose a technique based on the Coq proof assistant

- ▶ To specify a static analysis,
- ▶ To prove its correctness wrt. the semantics of the language,
- ▶ To extract a static analyser from the proof of existence of a correct program analysis result

Program-as-proofs paradigm:

Write a function $f$ which verifies a specification $P$ $\forall x,\ P(x, f(x))$ $\iff$ Make a constructive proof of $\forall x,\ \exists y,\ P(x, y)$

# Outline

- Motivation

- A Static Analysis for Carmel

- Building a certified static analyser

- Conclusion

# Outline

# Case study : a static analysis for Carmel

We follow the analysis proposed by René Rydhof Hansen[1]

- ▸ Carmel : an intermediate representation of Java Card byte code
- ▸ Construction of a certified data flow analyser for Carmel

[1]René Rydhof Hansen. Flow Logic for Carmel. SECSAFE-IMM-001, 2002

## Syntax of Carmel

Instruction ::=

```
nop
push c
pop                    ⎫
numop op               ⎬ stack manipulation
                       ⎭
load x                 ⎫
store x                ⎬ local variables manipulation
                       ⎭
if pc                  ⎫
goto pc                ⎬ jump
                       ⎭
new cl                 ⎫
putfield f             ⎬ heap manipulation
getfield f             ⎭
invokevirtual m_id     ⎫
return                 ⎬ method call and return
                       ⎭
```

## Semantic domains

$$
\begin{aligned}
\text{Val} \quad ::= \quad & \text{num } n & n \in \mathbb{N} \\
& \text{ref } r & r \in \text{Reference} \\
& \text{null} \\
\text{Stack} \quad = \quad & \text{Val}^* \\
\text{LocalVar} \quad = \quad & \text{Var} \rightarrow \text{Val} \\
\text{Frame} \quad = \quad & \text{PointProg} \times \text{NameMethod} \\
& \qquad \times \text{LocalVar} \times \text{Stack} \\
\text{CallStack} \quad = \quad & \text{Frame}^* \\
\text{Object} \quad = \quad & \text{FieldName} \rightarrow \text{Val} \\
\text{Heap} \quad = \quad & \text{Reference} \rightarrow \text{Object}_\perp \\
\text{State} \quad = \quad & \text{Heap} \times \text{CallStack}
\end{aligned}
$$

Example :

$$
(H, \langle m, pc, L, v :: S \rangle :: SF)
$$

## Dynamic semantics

Operational semantics with rules like

$$\frac{\text{instructionAt}_P(m, pc) = \texttt{push } c}{(H, \langle m, pc, L, S \rangle :: SF) \Rightarrow (H, \langle m, pc + 1, L, c :: S \rangle :: SF)}$$

$$\frac{\begin{array}{rcl} \text{instructionAt}_P(m, pc) &=& \texttt{invokevirtual } m_{id} \\ m' &=& \text{methodLookup}(m_{id}, h(loc)) \\ f' &=& \langle m', 1, V, \varepsilon \rangle \\ f'' &=& \langle m, pc, I, s \rangle \end{array}}{(h, \langle m, pc, I, loc :: V :: s \rangle :: sf) \Rightarrow (h, f' :: f'' :: sf)}$$

# A Static Analysis for Carmel

We want to calculate an approximation $\left( \hat{H}, \hat{L}, \hat{S} \right)$ on the domain

$$\widehat{\text{State}} = \widehat{\text{Heap}} \ \times \ \left( \text{NameMethod} \times \text{PointProg} \rightarrow \widehat{\text{LocalVar}} \right)$$
$$\times \ \left( \text{NameMethod} \times \text{PointProg} \rightarrow \widehat{\text{Stack}} \right)$$

- ▶ An approximation for all reachable heaps
- ▶ For each program points, an approximation of the operand stack and the local variables
- ▶ An object is abstracted to its class
- ▶ Numeric values are abstracted using Killdall's Constant Propagation domain

Each instruction impose constraints on $\left(\hat{H}, \hat{L}, \hat{S}\right)$.
**Example**

```
0 :  push 1
1 :  push 2
2 :  store 0
3 :  load 0
4 :  numop mult
5 :  goto 1
```

Each instruction impose constraints on $\left(\hat{H}, \hat{L}, \hat{S}\right)$.

**Example**

$$\widehat{\text{nil}} \sqsubseteq \hat{S}(m, 0) \qquad \top \sqsubseteq \hat{L}(m, 0)$$

0 : push 1

1 : push 2

2 : store 0

3 : load 0

4 : numop mult

5 : goto 1

## Analysis specification

Each instruction impose constraints on $(\hat{H}, \hat{L}, \hat{S})$.

**Example**

0 : push 1

$\widehat{\text{nil}} \sqsubseteq \hat{S}(m, 0)$          $\top \sqsubseteq \hat{L}(m, 0)$

$\widehat{\text{push}}(\hat{1}, \hat{S}(m, 0)) \sqsubseteq \hat{S}(m, 1)$      $\hat{L}(m, 0) \sqsubseteq \hat{L}(m, 1)$

1 : push 2

2 : store 0

3 : load 0

4 : numop mult

5 : goto 1

## Analysis specification

Each instruction impose constraints on $(\hat{H}, \hat{L}, \hat{S})$.

**Example**

0 : push 1

1 : push 2

2 : store 0

3 : load 0

4 : numop mult

5 : goto 1

$$\widehat{\text{nil}} \sqsubseteq \hat{S}(m, 0) \qquad\qquad \top \sqsubseteq \hat{L}(m, 0)$$

$$\widehat{\text{push}}(\hat{1}, \hat{S}(m, 0)) \sqsubseteq \hat{S}(m, 1) \qquad \hat{L}(m, 0) \sqsubseteq \hat{L}(m, 1)$$

$$\widehat{\text{push}}(\hat{2}, \hat{S}(m, 1)) \sqsubseteq \hat{S}(m, 2) \qquad \hat{L}(m, 1) \sqsubseteq \hat{L}(m, 2)$$

$$\widehat{\text{pop}}(\hat{S}(m, 2)) \sqsubseteq \hat{S}(m, 3) \qquad \hat{L}(m, 2)[0 \mapsto \widehat{\text{top}}(\hat{S}(m, 2))] \sqsubseteq \hat{L}(m, 3)$$

$$\widehat{\text{push}}(\hat{L}(m, 3)[0], \hat{S}(m, 3)) \sqsubseteq \hat{S}(m, 4) \qquad \hat{L}(m, 3) \sqsubseteq \hat{L}(m, 4)$$

$$\cdots \qquad\qquad \hat{L}(m, 4) \sqsubseteq \hat{L}(m, 5)$$

$$\hat{S}(m, 5) \sqsubseteq \hat{S}(m, 1) \qquad \hat{L}(m, 5) \sqsubseteq \hat{L}(m, 1)$$

The smallest value which verifies all constraints.
**Example**

| | | |
|---|---|---|
| 0 : push 1 | $\widehat{nil}$ | $[0 \mapsto \top; \ 1 \mapsto \top]$ |
| 1 : push 2 | $< \hat{2} >$ | $[0 \mapsto \hat{1}; \ 1 \mapsto \top]$ |
| 2 : store 0 | $< \hat{1} :: \hat{2} >$ | $[0 \mapsto \hat{1}; \ 1 \mapsto \top]$ |
| 3 : load 0 | $< \hat{2} >$ | $[0 \mapsto \hat{1}; \ 1 \mapsto \top]$ |
| 4 : numop mult | $< \hat{1} :: \hat{2} >$ | $[0 \mapsto \hat{1}; \ 1 \mapsto \top]$ |
| 5 : goto 1 | $< \hat{2} >$ | $[0 \mapsto \hat{1}; \ 1 \mapsto \top]$ |
| | $\bot$ | $\bot$ |

# Outline

# Building a certified static analyser



- A puzzle with 8 pieces,
- Each piece interacts with its neighbors

# Building a certified static analyser



- ▶ Each semantic domain is modeled with a type
- ▶ Following exactly the definitions already seen in a previous slide

# Building a certified static analyser



- ► Each semantic domain is in relation with an abstract domain
- ► an abstract domain is a lattice (formalization of lattices in Coq to follow...)

# Building a certified static analyser



- ▸ A relation $\sim$ between State and $\widehat{\text{State}}$
- ▸ $s \sim \widehat{\Sigma}$ interprets as "$\widehat{\Sigma}$ is a correct approximation of $s$"
- ▸ $\sim$ must be monotone :
  $\forall s \in \text{State}, \forall \widehat{\Sigma}_1, \widehat{\Sigma}_2 \in \widehat{\text{State}},$
  if $s \sim \widehat{\Sigma}_1$ and $\widehat{\Sigma}_1 \sqsubseteq \widehat{\Sigma}_2$ then $s \sim \widehat{\Sigma}_2$

# Building a certified static analyser



- ► The transition relation $\cdot \Rightarrow \cdot$ is defined using Coq inductive types
- ► Collecting semantics :
  $$[\![P]\!] = \{s \mid \exists s_0 \text{ an initial state, with } s_0 \Rightarrow^* s\}$$

We want to compute a correct approximation of $[\![P]\!]$

# Building a certified static analyser



- we define a predicate $P \vdash \widehat{\Sigma}$ which imposes a set of constraints on an abstract state $\widehat{\Sigma}$

# Building a certified static analyser



$$\forall P : \text{Program}, \ \forall \widehat{\Sigma} : \widehat{\text{State}}, \quad P \vdash \widehat{\Sigma} \Rightarrow [\![P]\!] \sim \widehat{\Sigma}$$

- One case by instruction
- With a special treatment for the `invokevirtual/return` instructions

# Building a certified static analyser



| semantic domains | correctness relations | abstract domains | |
|---|---|---|---|
| semantic rules | correctness proofs | analysis specification | constraint generator |

- ▸ Collects all constraint of a given program

# Building a certified static analyser



| semantic domains | correctness relations | abstract domains | constraints solver |
| --- | --- | --- | --- |
| semantic rules | correctness proofs | analysis specification | constraint generator |

$$\forall P : \text{Program}, \ \exists \widehat{\Sigma} : \widehat{\text{State}}, \quad P \vdash \widehat{\Sigma}$$

► In fact, a stronger result : there exists a smallest solution

# Building a certified static analyser



| semantic domains | correctness relations | abstract domains | constraints solver |
|---|---|---|---|
| semantic rules | correctness proofs | analysis specification | constraint generator |

## Final result

$$\left.\begin{array}{l} \forall P,\ \forall \widehat{\Sigma},\ P \vdash \widehat{\Sigma} \Rightarrow [\![P]\!] \sim \widehat{\Sigma} \\ \forall P,\ \exists \widehat{\Sigma},\ P \vdash \widehat{\Sigma} \end{array}\right\} \quad \forall P,\ \exists \widehat{\Sigma},\ [\![P]\!] \sim \widehat{\Sigma}$$

The `lattice` type is a big structure :

```
Record Lattice [A:Set] :  Type := {
  eq :  A → A → Prop;

  order :  A → A → Prop;

  join :  A → A → A;

  eq_dec :  A → A → bool;

  bottom :  A;

  top :  A;


}.
```

The `lattice` type is a big structure :
```
Record Lattice [A:Set] :  Type := {
  eq :  A → A → Prop;
    eq_prop :  ...; // eq is an equivalence relation
  order :  A → A → Prop;
    order_prop :  ...; // order is an order relation
  join :  A → A → A;
    join_prop :  ...; // join is a correct binary least upper bound
  eq_dec :  A → A → bool;
    eq_dec_prop :  ...; // eq_dec is a correct equality test
  bottom :  A;
    bottom_prop :  ...; // bottom is the smallest element
  top :  A;
    top_prop :  ...; // top is the biggest element
    acc_prop :  ...; // ⊐ is well founded (ascending chain condition)
}.
```

## Abstract domains

The `lattice` type is a big structure :

```
Record Lattice [A:Set] :  Type := {
  eq :  A → A → Prop;
    eq_prop :  ...; // eq is an equivalence relation
  order :  A → A → Prop;
    order_prop :  ...; // order is an order relation
  join :  A → A → A;
    join_prop :  ...; // join is a correct binary least upper bound
  eq_dec :  A → A → bool;
    eq_dec_prop :  ...; // eq_dec is a correct equality test
  bottom :  A;
    bottom_prop :  ...; // bottom is the smallest element
  top :  A;
    top_prop :  ...; // top is the biggest element
    acc_prop :  ...; // ⊒ is well founded (ascending chain condition)
}.
```

- ▶ Two base lattices
  - ▶ Flat lattice of constants
  - ▶ Lattice of sets over a finite subset of integer
- ▶ Four functions to combine lattices
  - ▶ Product of lattice
  - ▶ Sum of lattice
  - ▶ Arrays whose elements live in a lattice and whose size is bounded (efficient functional structure)
  - ▶ List whose elements live in a lattice

This modular construction saves a considerable amount of time and effort.

For this analysis :
$$\begin{array}{l}(\text{array (array (list (finiteSet} + \text{constants}))))\\ \times\ (\text{array (array (array (finiteSet} + \text{constants}))))\\ \times\ (\text{array (array (finiteSet} + \text{constants})))\end{array}$$

1. A generic fixed point solver

   $$\forall L : (\text{lattice } A), \ \forall f : A \to A, \ f \text{ monotone},$$
   $$\exists x : A, \ x \text{ is the least fixed point of } f$$

   proof : the sequence $\bot, f(\bot), f^2(\bot), \ldots$ stabilizes on the least fixed point

2. We use it to solve the functional constraints, using the fact that

   $$\begin{array}{ccc} x \text{ is the least solution of} & \iff & x \text{ is the least fixed point of} \\ f_1(x) \sqsubseteq x, \ldots, f_n(x) \sqsubseteq x & & \hat{f}_1 \circ \cdots \circ \hat{f}_n \end{array}$$

   with $\hat{f}(x) = f(x) \sqcup x$

3. Combined with the constraint generator, we obtain

   $$\forall P, \ \exists \widehat{\Sigma}, \ P \vdash \widehat{\Sigma}$$

# Outline

# Where is the proof effort ?



- ▶ Most technical part in the lattice library
- ▶ Correctness part does not require specific competences in Coq
- ▶ A majority of proof a reusable to develop others analysis

# Where is the programming effort ?



| semantic domains | correctness relations | abstract domains | constraints solver |
| semantic rules | correctness proofs | analysis specification | constraint generator |

- ▶ The extraction mechanism only keeps the computational content of proofs
- ▶ The corresponding parts require a high attention to obtain an efficient analyser

# Conclusion on the work

- ▶ We proposed a technique based on the Coq proof assistant
  - ▶ To develop a certified static analyser
  - ▶ To extract a correct analyser in Ocaml

- ▶ We illustrated this technique with a data flow analysis for the Carmel language
  - ▶ 10000 lines of Coq converted in 2000 lines of OCaml
  - ▶ With a reasonable efficiency of the analyser :
    - ▶ About 1 minute to analyse 1000 lines of Carmel byte code

# Further works

- ▶ Construction of an efficient certified work-set based program instead of the actual naive resolution
  - ▶ This program must be independent of the abstract domains
- ▶ Lattice of infinite height like intervals
- ▶ Automatization of the correctness proof ?
  - ▶ So as to quickly extend the number of language instructions
- ▶ A more extensive use of the abstract interpretation formalism
  - ▶ We must find a compromise between the reusability possibilities and the technical efforts in Coq
- ▶ Application of this technique to others languages

constraints
solver

relational interpretation
(predicat on $\widehat{State}$)

analysis
specification

constraint
generator

► Ideal for proofs
► Extraction is
compromised

functional interpretation
$(F(\widehat{\Sigma}) \sqsubseteq \widehat{\Sigma})$

constraints

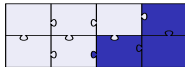- Difficult to use in proofs
- Can be extracted

solver

relational interpretation
(predicat on $\widehat{State}$)

- Ideal for proofs
- Extraction is compromised

analysis

specification

constraint

generator

functional interpretation
$(F(\widehat{\Sigma}) \sqsubseteq \widehat{\Sigma})$

- ▶ Difficult to use in proofs
- ▶ Can be extracted

relational interpretation
(predicat on $\widehat{State}$)

- ▶ Ideal for proofs
- ▶ Extraction is compromised

$\mathcal{F}[\cdot]$

intermediate interpretation

$\mathcal{R}[\cdot]$

# Correctness proof

$$\forall P : \text{Program}, \ \forall \widehat{\Sigma} : \widehat{\text{State}}, \quad P \vdash \widehat{\Sigma} \Longrightarrow [\![P]\!] \sim \widehat{\Sigma}$$

We prove it by well-founded induction on the length of the program execution.

Induction step:

- for all instructions $I$, except `return`

$$\forall P, \ \forall \widehat{\Sigma} : \widehat{\text{State}}, \ P \vdash \widehat{\Sigma} \Longrightarrow$$
$$\forall s_1 : \text{State}, \ s_1 \sim \widehat{\Sigma} \Longrightarrow$$
$$\forall s_2 : \text{State}, \ [s_1 \Rightarrow_I s_2] \Longrightarrow s_2 \sim \widehat{\Sigma}$$

- for the return instruction

$$\forall P, \ \forall \widehat{\Sigma} : \widehat{\text{State}}, \ P \vdash \widehat{\Sigma} \Longrightarrow$$
$$\forall s_1 : \text{State}, \ \left( \forall s, \ s \Rightarrow^* s_1 \Longrightarrow s \sim \widehat{\Sigma} \right) \Longrightarrow$$
$$\forall s_2 : \text{State}, \ [s_1 \Rightarrow_{\text{return}} s_2] \Longrightarrow s_2 \sim \widehat{\Sigma}$$