# Abstract Code Search

Isabel Garcia-Contreras    Jose F. Morales    Manuel V. Hermenegildo

New York, ICLP 2016, October 18, 2016

IMDEA Software Institute, and

Technical University of Madrid

## Introduction

**Motivation:** Find code to reuse.

**Problem:** The large amount of code available makes it impossible for a programmer to find useful fragments manually.

**Current code finders** are mainly based on:

- Documentation keyword indexing (Maarek et al. 1991).
- Signature matching (Rollins and Wing 1991 – $\lambda$ Prolog).
- Combinations of both (Mitchell 2008 – Hoogle/Haskell).

## Introduction

**Motivation:** Find code to reuse.

**Problem:** The large amount of code available makes it impossible for a programmer to find useful fragments manually.

**Current code finders** are mainly based on:

- Documentation keyword indexing (Maarek et al. 1991).
- Signature matching (Rollins and Wing 1991 – $\lambda$ Prolog).
- Combinations of both (Mitchell 2008 – Hoogle/Haskell).

**Our approach:**

- Find code based on its *semantic characteristics*.
  (that can be inferred automatically).

## (Ciao) Assertions

Assertions are a fundamental component of our approach.

They are **linguistic constructions** for expressing abstractions of the meaning and behavior of programs.

---

pred **assertions (subset)**

Allow specifying certain conditions on the state (current substitution or constraint store) that must hold at certain points of program execution.

$$:\text{- pred } Head : Pre \Rightarrow Post.$$

- *Head*: normalized atom that denotes the predicate that the assertion applies to.
- *Pre* and *Post*: conjunctions of "prop" atoms.

---

## (Ciao) Assertions - example

```
1  :- pred check_length(L,N) : (list(L), int(N)).
2  check_length(L,N) :- length(L,N).
3
4  :- pred gen_list(L,N) : (var(L),  var(N)) => (list(L), int(N)).
5  gen_list(L,N) :- length(L,N).
6
7  % length implementation ...
```

Herein we will use *assertions* for:

- **Programming**: To specify calling modes for predicates.
- **[New!] Searching**: To express properties of the code to be found.

## Our approach

Main idea of our method:

1. A **set of modules** is specified within which code is to be found.

## Our approach

Main idea of our method:

1. A **set of modules** is specified within which code is to be found.

```
:- module(lengths, [check_length/2, gen_list/2], [assertions]).

:- pred check_length(L,N) : (list(L), int(N)).
check_length(L,N) :- length(L,N).

:- pred gen_list(L,N) : (var(L), var(N)) => (list(L), int(N)).
gen_list(L,N) :- length(L,N).

% length implementation ...
```

## Our approach

Main idea of our method:

1. A **set of modules** is specified within which code is to be found.
2. A **static pre-analysis** is made to:
   - Infer semantic properties in one or more abstract domains
     (e.g.: shapes/types, variable sharing, inst. modes, polyhedra, ...).

# Static analysis inference

### Shapes and Modes/Sharing analysis:

```
1  % :- pred check_length(L,N)     : (list(L), int(N)).
2  :- true pred check_length(L,N)  : (list(L), int(N))
3                                   => (list(L), int(N)).
4  :- true pred check_length(L,N)  : (mshare([[L],[L,N],[N]])
5                                   => (mshare([[L]]), ground([N])).
6  check_length(L,N) :- length(L,N).
7
8  %:- pred gen_list(L,N)      : ((var(L),  var(N) ).
9  :- true pred gen_list(L,N) : ((term(L), term(N))
10                              => ((list(L), int(N)).
11 :- true pred gen_list(L,N) : (mshare([[L],[L,N],[N]])), var(L), var(N))
12                              => (mshare([[L]]), ground([N])).
13 gen_list(L,N) :- length(L,N).
14
15 % length implementation ...
```

## Our approach

Main idea of our method:

1. A **set of modules** is specified within which code is to be found.
2. A **static pre-analysis** is made to:
   - Infer semantic properties in one or more abstract domains (e.g.: shapes/types, variable sharing, inst. modes, polyhedra, ...).
3. User specifies **semantic properties** in a query, using a new kind of assertions that we call **query assertions**:
   - Example: `:- pred P(X,Y) : list(X) => sorted(Y).`

   Based on (Stulova et al. 2014).
4. Look within the modules for predicates that meet those properties (by comparing to inferred information).

## Inferring semantic properties

**Abstract interpretation:**

- To simulate the execution using an *abstract domain* $D_\alpha$.
- It guarantees:
    - Analysis termination, provided that $D_\alpha$ meets some conditions.
    - Results are **safe approximations** of the concrete semantics.

| **PLAI algorithm (in Ciao)** | |
|---|---|
| Input | **P**: Program |
| | $D_\alpha^i$: Abstract Domain(s) |
| | $\mathcal{Q}_\alpha = L{:}\lambda$: Initial call pattern |
| Output | $Analysis(P, L{:}\lambda, D_\alpha) = \{\langle L_1, \lambda_1^c, \lambda_1^s\rangle, \ldots, \langle L_n, \lambda_n^c, \lambda_n^s\rangle\}$, where: |
| | - $L_i$ is an atom |
| | - $\lambda_i^c$ are abstract call state in $D_\alpha$ |
| | - $\lambda_i^s$ are abstract success state in $D_\alpha$. |

## Querying for predicates

### Predicate query

?- findp({ *As* }, M:Pred/A, Residue, Status).

| | |
|---|---|
| Input | **As**: Set of query assertions. |
| Output | **M:Pred/A**: Module, predicate descriptor and arity. |
| | **Residue**: Information of condition matching. |
| Input or Output | **Status**: Result of the proof for the whole set of conditions. |
| | - *checked* if all conditions are proved to be checked. |
| | - *false* if any condition is false. |
| | - *check* if neither *checked* nor *false* can be proved. |

```
?- findp({ :- pred P(L, S) => (list(L), num(S)). }, M:P/A, Res, St).

 P/A = check_length/2    St = checked
```

## Normalizing the query

### Assertion Conditions

Given a predicate represented by a normalized atom *Head*, and a corresponding set of assertions $\mathcal{A} = \{A_1 \ldots A_n\}$, with $A_i = \text{":- pred } Head : Pre_i \Rightarrow Post_i.\text{"}$. The set of *assertion conditions* for *Head* determined by $\mathcal{A}$ is $\{C_0, C_1, \ldots, C_n\}$, with:

$$C_i = \begin{cases} \text{calls}(Head, \bigvee_{j=1}^{n} Pre_j) & i = 0 \\ \text{success}(Head, Pre_i, Post_i) & i = 1..n \end{cases}$$

```
1  :- pred my_length(L,N) : ((list(L), var(N)) => ((list(L), int(N)).
2  :- pred my_length(L,N) : ((list(L), int(N)) => ((list(L), int(N)).
3  my_length(L,N) :- length(L,N).
```

Assertion conditions from `my_length/2`:

$$C_i = \begin{cases} \text{calls(} & length(L, N), & (list(L) \wedge var(N)) \;\vee\; (list(L) \wedge int(N))), & \\ \text{success(} & length(L, N), & (list(L) \wedge var(N)), & (list(L), int(N))), \\ \text{success(} & length(L, N), & (list(L) \wedge int(N)), & (list(L), int(N))) \end{cases}$$

## Matching *success* conditions

Match $C = \texttt{success}(X(V_1, \ldots, V_n), Pre, Post)$ against *analysis*$(P, Q_\alpha)$

**Checked matches**

If *Pre* holds at the time of calling the matching predicate and the execution succeeds then the *Post* conditions hold.

$C$ is abstractly 'checked' for predicate $p \in P$ w.r.t. $Q_\alpha$ in $D_\alpha$ iff
$\exists L = p(V'_1, \ldots, V'_n)$ s.t. $\forall \langle L, \lambda^c, \lambda^s \rangle \in$ *analysis*$(P, Q_\alpha)$ s.t. $\exists \sigma \in$ *ren*, $L = p(V'_1, \ldots, V'_n) = X(V_1, \ldots, V_n)\sigma$, $\lambda^c \sqsupseteq \lambda^+_{TS(Pre\ \sigma, P)} \to \lambda^s \sqsubseteq \lambda^-_{TS(Post\ \sigma, P)}$

**False matches**

If *Pre* holds at the time of calling the matching predicate and the execution succeeds then its conditions and *Post* are disjoint.
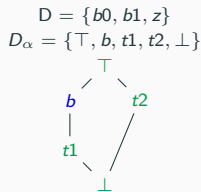
$C$ is abstractly false for $p \in P$ w.r.t. $Q_\alpha$ in $D_\alpha$ iff $\exists L = p(V'_1, \ldots, V'_n)$ s.t. $\forall \langle L, \lambda^c, \lambda^s \rangle \in$ *analysis*$(P, Q_\alpha)$ s.t. $\exists \sigma \in$ *ren*, $L = p(V'_1, \ldots, V'_n) = X(V_1, \ldots, V_n)\sigma$, $\lambda^c \sqsubseteq \lambda^-_{TS(Pre\ \sigma, P)} \wedge (\lambda^s \sqcap \lambda^+_{TS(Post\ \sigma, P)} = \bot)$

```
1
2   perfect(b1).                           mixed(b0).
3   perfect(b0).                           mixed(b1).
4                                          mixed(z).
5
6   reduced(b1).
7                                          hard(X) :- functor(b1(_), X, _).
8
9   outb(z).
10
11  :- regtype b/1.
12  b(b0).
13  b(b1).
```

# Matching *success* conditions - Example

```
1   :- true pred perfect(A) => b(A).          :- true pred mixed(X) => top(X).
2   perfect(b1).                               mixed(b0).
3   perfect(b0).                               mixed(b1).
4                                              mixed(z).
5   :- true pred reduced(A) => t1(A).
6   reduced(b1).                               :- true pred hard(X) => top(X).
7                                              hard(X) :- functor(b1(_), X, _).
8   :- true pred outb(A) => t2(A).
9   outb(z).
10
11  :- regtype b/1.        :- regtype t1/1    :- regtype t2/1.
12  b(b0).                 t1(b1).            t2(z).
13  b(b1).
```

$$D = \{b0, b1, z\}$$
$$D_\alpha = \{\top, b, t1, t2, \bot\}$$

```
1   :- true pred perfect(A) => b(A).           :- true pred mixed(X) => top(X).
2   perfect(b1).                               mixed(b0).
3   perfect(b0).                               mixed(b1).
4                                              mixed(z).
5   :- true pred reduced(A) => t1(A).
6   reduced(b1).                               :- true pred hard(X) => top(X).
7                                              hard(X) :- functor(b1(_), X, _).
8   :- true pred outb(A) => t2(A).
9   outb(z).
10
11  :- regtype b/1.         :- regtype t1/1    :- regtype t2/1.
12  b(b0).                  t1(b1).            t2(z).
13  b(b1).
```

```
?- findp({ :- pred P(V) => b(V). }, M:P/A, Res, St).
   P/A = perfect/1   St = checked
   P/A = reduced/1   St = checked
   P/A = outb/1      St = false
   P/A = mixed/1     St = check
   P/A = hard/1      St = check
```
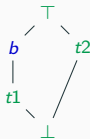
$D = \{b0, b1, z\}$
$D_\alpha = \{\top, b, t1, t2, \bot\}$

## Matching *calls* conditions - Example

```prolog
1  :- pred check_length(L,N) : ((list(L), int(N)).
2  check_length(L,N) :- length(L,N).
3
4  :- pred gen_list(L,N) : ((var(L), var(N)).
5  gen_list(L,N) :- length(L,N).
6
7  :- pred get_length(L,N) : var(N).
8  get_length(L,N) :- length(L,N).
```

```prolog
?- findp({ :- pred P(L, Size) : var(L), var(Size)). }, M:P/A, Res, St).

 P/A = check_length/2    St = false
 P/A = gen_list/2        St = checked
 P/A = get_length/2      St = check
```

## Combining abstract domains

```
1  :- true pred check_length(L,N) : (list(L), int(N)) => (list(L), int(N)).
2  :- true pred check_length(L,N) : (mshare([[L],[L,N],[N]])
3                                 => (mshare([[L]]), ground([N])).
4  check_length(L,N) :- length(L,N).
5
6  :- true pred gen_list(L,N) : ((term(L), term(N)) => ((list(L), int(N)).
7  :- true pred gen_list(L,N) : (mshare([[L],[L,N],[N]]), var(L), var(N))
8                                 => (mshare([[L]]), ground([N])).
9  gen_list(L,N) :- length(L,N).
10
```

?- findp({ :- pred P(L, Size) : list(L), num(Size)). }, M:P/A, Res, St).

| PredName/A | regtypes proof | shfr proof | combined proof |
|---|---|---|---|
| check_length/2 | checked | check | checked |
| gen_list/2 | check | false | false |

**Demo**

## Conclusions

**Finding code by its semantic characteristics:**

- Ensures that the found code behaves correctly.
- Reasons with relations between properties (implication, abstraction).
- Is independent from the documentation.
- Implemented in Ciao, in combination with other types of search (which it complements).

**Future work:**

- Use more domains.
- Extend to other programming languages and combinations.
- Other uses: finding duplicated code.

**Thanks!**

## Matching *calls* conditions

$$C = \texttt{calls}(X(V_1, \ldots, V_n), Pre)$$

### Checked matches

The admissible calls of the matching predicate are within the set of *Pre* conditions.

---

$C$ is abstractly 'checked' for a $p \in P$ w.r.t. $Q_\alpha$ in $D_\alpha$ iff
$\forall \langle L, \lambda^c, \lambda^s \rangle \in analysis(P, D_\alpha, Q_\alpha)$ s.t. $\exists \sigma \in ren,\ L = p(V_1, \ldots, V_n) = X(V_1, \ldots, V_n)\sigma, \lambda^c \sqsubseteq \lambda^-_{TS(Pre\ \sigma, P)}$.

### False matches

The admissible calls of the matching predicate and the set *Pre* conditions are disjoint.

---

$C$ is abstractly 'false' for a $p \in P$ w.r.t. $Q_\alpha$ in $D_\alpha$ iff
$\forall \langle L, \lambda^c, \lambda^s \rangle \in analysis(P, D_\alpha, Q_\alpha)$ s.t. $\exists \sigma \in ren,\ L = p(V_1, \ldots, V_n) = X(V_1, \ldots, V_n)\sigma, \lambda^c \sqcap \lambda^+_{TS(Pre\ \sigma, P)} = \bot$.

# Performance

Search times in $\mu s$.

| Ar\Cnds | 1 | 1 (AVG) | 2 | 2 (AVG) | 3 | 3 (AVG) | 4 | 4 (AVG) |
|---|---|---|---|---|---|---|---|---|
| 1 (85 pr) | 19,064 | 224 | 53,530 | 630 | 180,246 | 2,121 | 298,292 | 3,509 |
| 2 (74 pr) | 110,092 | 1,488 | 207,871 | 2,809 | 221,061 | 2,987 | 477,440 | 6,452 |
| 3 (47 pr) | 294,962 | 6,276 | 3,757,208 | 79,941 | 3,806,917 | 80,998 | 6,127,015 | 130,362 |
| 4 (12 pr) | 5,116 | 426 | 12,939 | 1,078 | 22,508 | 1,876 | 30,300 | 2,525 |