

Generación de código con certificado asociado

Ricardo Peña Marí

Catedrático de Lenguajes y Sistemas Informáticos
Departamento de Sistemas Informáticos y Computación
Universidad Complutense de Madrid

F. Informática S. Sebastián, enero 2008

Motivation (1)

- **Proof Carrying Code (PCC)** is an emerging research trend aimed at endowing programs with certificates guaranteeing that certain properties are satisfied. These certificates are in fact mathematical correctness proofs that can be automatically checked by appropriate tools.
- PCC is developing quickly, specially in the area of mobile code, where consumers of (normally untrusted) code wish to ensure that their safety policy is fulfilled, before running the code in their machines.
- Our project SELF (TIN2004-07943-C04-04) is developing **SAFE**, a functional language targeted at mobile applications with limited resources, and aimed at automatically generating certificates about the following properties:
 - ① Absence of dangling pointers
 - ② Termination
 - ③ Bounded heap and stack
- **SAFE** may use an explicit construction to liberate memory and does not need a garbage collector.

Motivation (2)

FUNCTIONAL LANGUAGES

Implicit Approach (Garbage Collection)

High-level
programming

vs.

Time delay
+
No reasoning about
heap needs

IMPERATIVE LANGUAGES

Explicit Approach (Pointers)

Control over
memory usage

vs.

Error-prone
programs
(dangling references,
undesired sharing)

SAFE

Semi-explicit Approach (Destructive Pattern Matching)

No garbage collection
Reasoning about the heap needed by the program
"Error-free" programs

Conventional version of mergesort

$$\begin{aligned}
 \textit{split } 0 \textit{ } xs &= ([], xs) \\
 \textit{split } n \textit{ } [] &= ([], []) \\
 \textit{split } n \textit{ } (x : xs) &= (x : xs_1, xs_2) \\
 &\quad \textbf{where } (xs_1, xs_2) = \textit{split } (n - 1) \textit{ } xs
 \end{aligned}$$

$$\begin{aligned}
 \textit{merge } [] \textit{ } ys &= ys \\
 \textit{merge } xs \textit{ } [] &= xs \\
 \textit{merge } (x : xs) \textit{ } (y : ys) & \\
 \quad | x \leq y &= x : \textit{merge } xs \textit{ } (y : ys) \\
 \quad | \textit{otherwise} &= y : \textit{merge } (x : xs) \textit{ } ys
 \end{aligned}$$

$$\begin{aligned}
 \textit{msort } xs & \\
 \quad | n \leq 1 &= xs \\
 \quad | \textit{otherwise} &= \textit{merge } (\textit{msort } xs_1) (\textit{msort } xs_2) \\
 \quad \textbf{where } (xs_1, xs_2) &= \textit{split } (n \textit{ 'div' } 2) \textit{ } xs \\
 \quad n &= \textit{length } xs
 \end{aligned}$$

Allocated memory in $\Theta(n \log n)$ and resident memory in $\Theta(n)$

SAFE version of mergesort

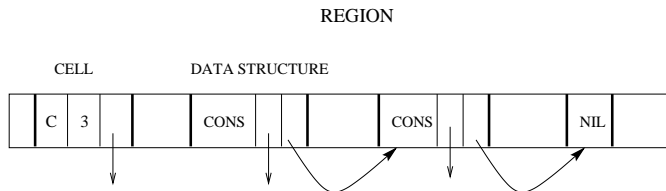
$$\begin{aligned}
 \text{splitD} &:: \forall a, \rho. \text{Int} \rightarrow [a]!@ \rho \rightarrow \rho \rightarrow ([a]@ \rho, [a]@ \rho)@ \rho \\
 \text{splitD } 0 \text{ } xs! &= ([], xs!) \\
 \text{splitD } n \text{ } []! &= ([], []) \\
 \text{splitD } n \text{ } (x : xs)! &= (x : xs_1, xs_2) \\
 &\quad \text{where } (xs_1, xs_2) = \text{splitD } (n - 1) \text{ } xs \\
 \text{mergeD} &:: \forall a, \rho. [a]!@ \rho \rightarrow [a]!@ \rho \rightarrow \rho \rightarrow [a]@ \rho \\
 \text{mergeD } []! \text{ } ys! &= ys! \\
 \text{mergeD } xs! \text{ } []! &= xs! \\
 \text{mergeD } (x : xs)! \text{ } (y : ys)! & \\
 &\quad | \ x \leq y \quad = x : \text{mergeD } xs \text{ } (y : ys)! \\
 &\quad | \ \text{otherwise} = y : \text{mergeD } (x : xs!) \text{ } ys \\
 \text{msortD} &:: \forall a, \rho. [a]!@ \rho \rightarrow \rho \rightarrow [a]@ \rho \\
 \text{msortD } xs & \\
 &\quad | \ n \leq 1 \quad = xs! \\
 &\quad | \ \text{otherwise} = \text{mergeD } (\text{msortD } xs_1) \text{ } (\text{msortD } xs_2) \\
 &\quad \text{where } (xs_1, xs_2) = \text{splitD } (n \text{ 'div' } 2) \text{ } xs \\
 &\quad \quad n \quad = \text{length } xs
 \end{aligned}$$

Resident memory in $\Theta(n)$. Additional heap memory in $\Theta(1)$

Index

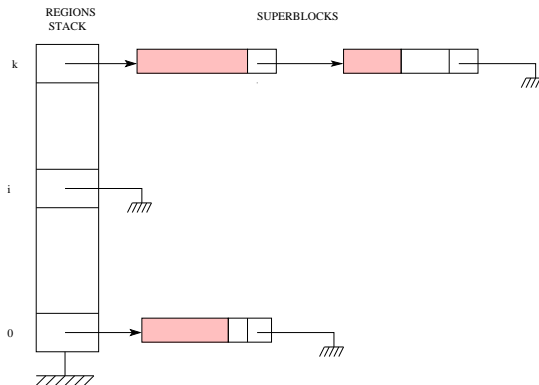
- ① Language Design
- ② Region inference
- ③ Sharing Analysis
- ④ Safe type system and type inference
- ⑤ Termination analysis
- ⑥ Memory consumption analysis
- ⑦ Certificate generation

Language Design: memory allocation scheme



- A data structure (DS) completely resides in one region.
- Two DSs can be one of them part of the other, or they can share a third DS.
- The basic values —integers, booleans, etc.— do not allocate cells in regions. They live inside the cells of DSs, or in the stack.

Memory Management



Action	Cost in time
Region creation	Constant
Region disposal	Constant
Cell allocation	Constant
Cell disposal	Constant

Core SAFE normalised syntax

$prog$	$\rightarrow dec_1; \dots; dec_n; e$	{program}
dec	$\rightarrow f \overline{x}_i^n @ \overline{r}_j^m = e$	{single-recursive, polymorphic function definition}
a	$\rightarrow c$	{atom}
	x	{variable}
	e	{expression}
e	$\rightarrow a$	{atom}
	$x@r$	{data structure copying}
	$x!$	{data structure reusing}
	$f \overline{a}_i^n @ \overline{r}_j^m$	{function application}
	let $x_1 = be$ in e	{non-recursive, monomorphic}
	case x of \overline{alt}_i^n	{read-only case}
	case! x of \overline{alt}_i^n	{destructive case}
alt	$\rightarrow C \overline{x}_i^n \rightarrow e$	{case alternative}
	e	{binding expression}
be	$\rightarrow C \overline{a}_i^n @ r$	{constructor application}
	e	{normal expression}

Example: List inversion

Non-destructive

$$\text{revF} :: \forall a, \rho_1, \rho_2. [a]@_{\rho_1} \rightarrow \rho_2 \rightarrow [a]@_{\rho_2}$$

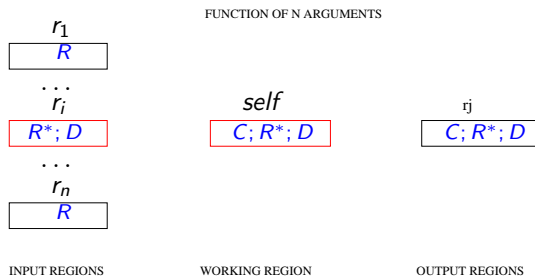
$$\begin{aligned} \text{revF } xs @r &= \mathbf{let} \ ys' = [] @r \mathbf{in} \ (\text{revauxF } xs \ ys') @r \\ \text{revauxF} &:: \forall a, \rho_1, \rho_2. [a]@_{\rho_1} \rightarrow [a]@_{\rho_2} \rightarrow \rho_2 \rightarrow [a]@_{\rho_2} \\ \text{revauxF } xs \ ys @r &= \mathbf{case} \ xs \ \mathbf{of} \\ &\quad [] \rightarrow ys \\ &\quad x : xx \rightarrow \mathbf{let} \ ys' = x : ys @r \mathbf{in} \\ &\quad \quad (\text{revauxF } xx \ ys') @r \end{aligned}$$

Destructive (in-place inversion)

$$\text{revD} :: \forall a, \rho_1, \rho_2. [a]!@_{\rho_1} \rightarrow \rho_2 \rightarrow [a]@_{\rho_2}$$

$$\begin{aligned} \text{revD } xs @r &= \mathbf{let} \ ys' = [] @r \mathbf{in} \ (\text{revauxD } xs \ ys') @r \\ \text{revauxD } xs \ ys @r &= \mathbf{case!} \ xs \ \mathbf{of} \\ &\quad [] \rightarrow ys \\ &\quad x : xx \rightarrow \mathbf{let} \ ys' = x : ys @r \mathbf{in} \\ &\quad \quad (\text{revauxF } xx \ ys') @r \end{aligned}$$

Region usage

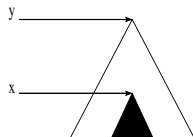
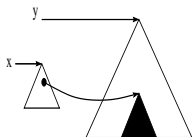
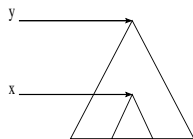
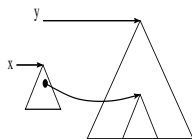


- A function may only read a DS which is a read-only parameter R .
- A function may read (before destroying it), and must destroy, a DS which is a condemned parameter $R^*; D$.
- A function may construct, read, or destroy DSs, in either its output regions or its (single) working region $C; R^*; D$.

Region inference

- Region inference takes place as a by-product of the Hindley-Milner type inference phase of the compiler.
- Program region variables r are mandatory at constructions $C \bar{a}_i^n @ r$ and at function application $f \bar{a}_i^n @ \bar{r}_j^m$ for known functions f . The compiler assigns to them **fresh** region types.
- The algorithm collects in three different sets,
 - ① the region types of the mandatory region variables,
 - ② the region types of the data structure produced by the function, and
 - ③ the region types of the input data structures.
- By solving a collection of restrictions, it takes the following decisions:
 - ① Whether zero or more output regions are needed by the function
 - ② For each program region variable r appearing in the text, it decides whether it must be assigned to an output region or to the **self** region.
- The algorithm never fails and it maximises the number of constructions done in the **self** region. This implies that as much garbage as possible will be collected when the function returns.

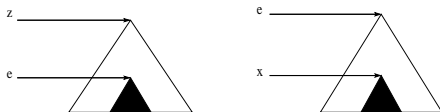
Sharing Relations

 $x \triangleleft \sim y$ (*SubR*)

 $x \triangleleft \sim y$ (*ShR*)

 $x \triangleleft y$ (*Sub*)

 $x \triangleleft y$ (*Sh*)


Sharing Analysis

- The interpretation traverses **top-down** the expression, accumulating the sharing relations (R):

$$S \llbracket e \rrbracket \text{SubR}' \text{ShR}' \text{Sub}' \text{Sh}' \rho = (\text{SubRP}, \text{ShRP}, \text{SubP}, \text{SubR}, \text{ShR}, \text{Sub}, \text{Sh})$$



- Function signatures** are kept in environment ρ :

$$(\{\text{Int}\}, \{\text{Int}\}, \{\text{Int}\}, \{\text{Int}\}, \{\text{Int}\}, \{\text{Int}\}, \{\text{Int}\})$$

- Sets $\text{sharerec}(x_1, e)$ and $\text{shareall}(x_1, e)$ are obtained from a global analysis of the whole program which annotates expression e with inherited information. They are respectively sets $\text{ShR}'(x_1)$ and $\text{Sh}'(x_1)$.

Examples

- **Reverse:** the inverse list only shares the elements with the original list:

$$\rho(\text{revaux}D) = (\{2\}, \{2\}, \{2\}, \{2\}, \{2\}, \{2\}, \{1, 2\})$$

$$\rho(\text{rev}D) = (\emptyset, \emptyset, \emptyset, \emptyset, \emptyset, \emptyset, \{1\})$$

- **Tree insert:** the new tree shares with the original tree both the elements and the recursive parts that are not on the path from the root to the inserted element:

$$\text{insert}D :: \forall a, \rho. a \rightarrow \text{Tree } a!@ \rho \rightarrow \rho \rightarrow \text{Tree } a@ \rho$$

$$\text{insert}D \times \text{Empty}! = \text{Node } \text{Empty} \times \text{Empty}$$

$$\text{insert}D \times (\text{Node } l \text{ y } r)!$$

$$| x < y = \text{Node } (\text{insert}D \times l) \text{ y } r!$$

$$| x == y = \text{Node } l! \text{ y } r!$$

$$| x > y = \text{Node } l! \text{ y } (\text{insert}D \times r)$$

$$\rho(\text{insert}D) = (\emptyset, \{1, 2\}, \emptyset, \emptyset, \{2\}, \{1\}, \{1, 2\})$$

Type Expressions

-- Polymorphic, with region type variables

$$\begin{aligned} \text{revF} &:: \forall a, \rho_1, \rho_2. [a]@_{\rho_1} \rightarrow \rho_2 \rightarrow [a]@_{\rho_2} \\ \text{insertF} &:: \forall a, \rho. a \rightarrow \text{Tree } a@_{\rho} \rightarrow \rho \rightarrow \text{Tree } a@_{\rho} \end{aligned}$$

-- Condemned parameters are reflected in the type

$$\begin{aligned} \text{revD} &:: \forall a, \rho_1, \rho_2. [a]!@_{\rho_1} \rightarrow \rho_2 \rightarrow [a]@_{\rho_2} \\ \text{insertD} &:: \forall a, \rho. a \rightarrow \text{Tree } a!@_{\rho} \rightarrow \rho \rightarrow \text{Tree } a@_{\rho} \end{aligned}$$

-- A data structure completely resides in one region

$$\begin{aligned} [] &:: \forall a, \rho. \rho \rightarrow [a]@_{\rho} \\ (:) &:: \forall a, \rho. a \rightarrow [a]@_{\rho} \rightarrow \rho \rightarrow [a]@_{\rho} \\ \text{Empty} &:: \forall a, \rho. \rho \rightarrow \text{Tree } a@_{\rho} \\ \text{Node} &:: \forall a, \rho. \text{Tree } a@_{\rho} \rightarrow a \rightarrow \text{Tree } a@_{\rho} \rightarrow \rho \rightarrow \text{Tree } a@_{\rho} \end{aligned}$$

Rules of the Type System (1)

$$\frac{
 \begin{array}{l}
 P = \{x_i \mid t_i = T_i \overline{s_{ij}^k} @ \rho_i, i = 1, \dots, n\}, \quad \text{fresh } \rho', \quad \rho' \notin \text{regions}(s) \\
 \Gamma^P + \overline{[x_i : t_i]^n} + [r : \rho] + [\text{self} : \rho'] + [f : \overline{t_i^n} \rightarrow \rho \rightarrow s] \vdash e : s
 \end{array}
 }{
 \{\Gamma^\emptyset\} \quad f \overline{x_i^n} r = e \quad \{\Gamma^\emptyset + [f : \text{gen}(\overline{t_i^n} \rightarrow \rho \rightarrow s, \Gamma)]\}
 } \quad \text{[FUNB]}$$

- DSs may only be constructed in region *self* and in the output region.
- The result of the function is contained in visible regions.
- P collects the read-only parameters.

$$\frac{
 \begin{array}{l}
 \Gamma^P(x) = d \quad R = \text{sharerec}(x, x!) - \{x\} \\
 \Gamma_R = \{y : \text{danger}(\text{type}(y)) \mid y \in R\}
 \end{array}
 }{
 \Gamma_R + \Gamma^P \vdash x! : \overline{d}
 } \quad \text{[REUSE]}$$

- Condemned DSs can be reused by changing its type, but they and their *sharerec* may not be accessed again.

Rules of the Type System (2)

$$\frac{\Gamma_1^P \vdash e_1 : s_1 \quad \Gamma_2^P + [x_1 : d_1] \vdash e : s \quad d_1 = \bar{s}_1 \quad P \cap R = \emptyset}{\Gamma_1^P \triangleright^{fv(e)} \Gamma_2^P \vdash \mathbf{let} \ x_1 = e_1 \ \mathbf{in} \ e : s} \text{ [LET2]}$$

Incorrect typings

let $x = \mathbf{case!} \ xs \ \mathbf{of} \ \dots$
in $\dots xs \dots$

A variable destructively used in e_1 may not be referenced in e

let $ys = (2 : xs)@self \ \mathbf{in}$
let $zs = (3 : xs)@self \ \mathbf{in}$
case! $zs \ \mathbf{of} \ \dots$

Variables sharing contents with $ys : s$ may not be destroyed in e

let $t = (\mathit{Node} \ i \ x \ d)@r \ \mathbf{in}$
case! $t \ \mathbf{of} \ \dots i \dots$

A recursive substructure of a destroyed DS may not be referenced in e

let $xs = ys \ \mathbf{in} \ \mathbf{case!} \ xs \ \mathbf{of} \ \dots \quad (ys \in P), \quad \text{A read-only parameter may not be destroyed}$

Rules of the Type System (3)

$$\frac{
 \begin{array}{l}
 \forall z \in R \cup \{x\}, i \in \{1..n\}. z \notin \text{fv}(e_i) \\
 (\forall i \in \{1..n\}). \Sigma(C_i) = \sigma_i \\
 R = \text{sharerec}(x, \text{case! } x \text{ of } \overline{C_i \overline{x_{ij}^{n_i}} \rightarrow e_i}) - \{x\} \\
 \Gamma_R = \{x : T !@ \rho\} \cup \{y : \text{dgr}(\text{ty}(y)) \mid y \in R\}
 \end{array}
 \quad
 \begin{array}{l}
 (\forall i \in \{1..n\}). \overline{s_{ij}^{n_i}} \rightarrow \rho \rightarrow T @ \rho \trianglelefteq \sigma_i \\
 (\forall i \in \{1..n\}. \forall j \in \{1..n_i\}). \text{inh!}(t_{ij}, s_{ij}, T !@ \rho) \\
 (\forall i \in \{1..n\}). \Gamma^P + \overline{[x_{ij} : t_{ij}]^{n_i}} \vdash e_i : s
 \end{array}
 }{
 \Gamma_R \otimes \Gamma^P \vdash \text{case! } x \text{ of } \overline{C_i \overline{x_{ij}^{n_i}} \rightarrow e_i} : s
 } \text{ [CASE!]}$$

Incorrect typings

case! t of

Node i x d → ... t ...

A destroyed DS and its *sharerec* may not be referenced again

Correct typings

case! t of

Node i x d → case! i of

Empty → d!

The substructures of a destroyed DS inherit the condemned property and may be either destroyed or reused

case t of

Node i x d → case! t of ...

A condemned DS may be read before being destroyed

Safe Types Inference Algorithm

- The first phase is a **Hindley-Milner inference algorithm** with some added complications arising from typing region variables and doing region inference.
- Next, the abstract syntax tree (AST) is traversed bottom-up collecting variables names in four sets. The rules $e \vdash_{inf} (D, R, S, N)$ classify variables in scope in e respectively as condemned, in-danger, safe and don't-know.

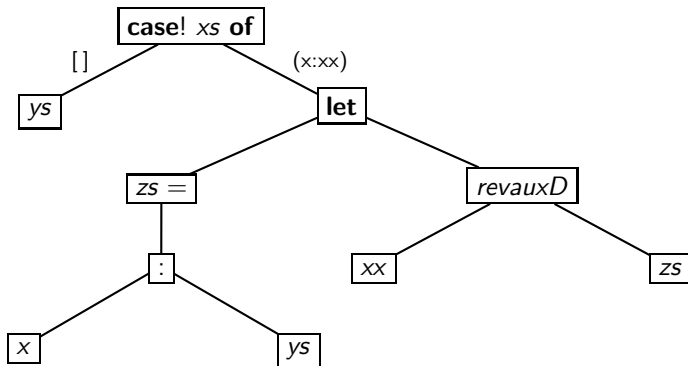
Safe Types Inference Algorithm

- The first phase is a **Hindley-Milner inference algorithm** with some added complications arising from typing region variables and doing region inference.
- Next, the abstract syntax tree (AST) is traversed bottom-up collecting variables names in four sets. The rules $e \vdash_{inf} (D, R, S, N)$ classify variables in scope in e respectively as condemned, in-danger, safe and don't-know.
- At nodes such as **let** and **case**, the information coming from different subexpressions is **checked for consistency**, so detecting possible typing errors.
- Also, some variables classified as **don't-know** in one branch may have acquired a different mark in another branch. This information is propagated top-down by the rules $(D, R, S) \vdash_{check} e$ in order to check that the new marks are consistent in an expression e .

Safe Types Inference Algorithm

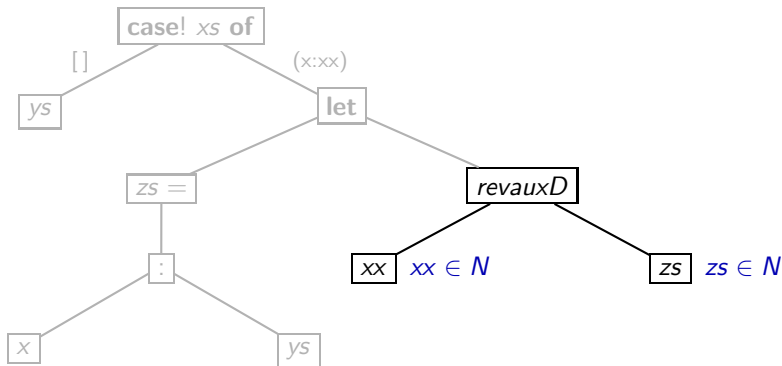
- The first phase is a **Hindley-Milner inference algorithm** with some added complications arising from typing region variables and doing region inference.
- Next, the abstract syntax tree (AST) is traversed bottom-up collecting variables names in four sets. The rules $e \vdash_{inf} (D, R, S, N)$ classify variables in scope in e respectively as condemned, in-danger, safe and don't-know.
- At nodes such as **let** and **case**, the information coming from different subexpressions is **checked for consistency**, so detecting possible typing errors.
- Also, some variables classified as **don't-know** in one branch may have acquired a different mark in another branch. This information is propagated top-down by the rules $(D, R, S) \vdash_{check} e$ in order to check that the new marks are consistent in an expression e .
- So, the algorithm **oscillates** between the bottom-up and the top-down modes. Also, a **fix point computation** is done at each recursive function definition. Its arguments are initially classified as don't-know and its body analysed until all argument marks stabilise. In-danger and don't-know arguments are not allowed in a valid stable type.

Example: *revauxD* - First iteration



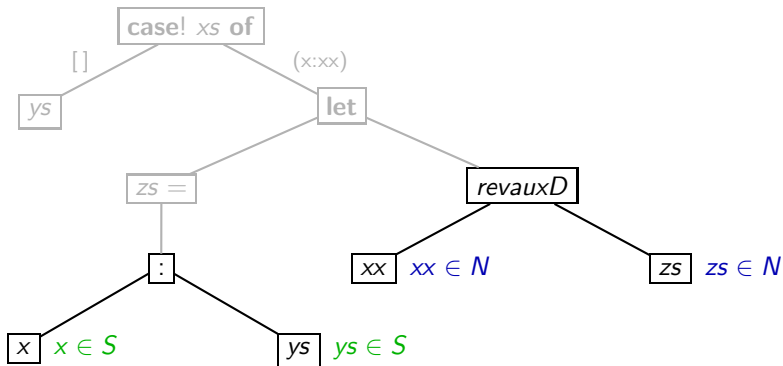
- We review the example of reversing a list. The definition of *revAuxD* is represented by means of its abstract syntax tree.

Example: *revauxD* - First iteration



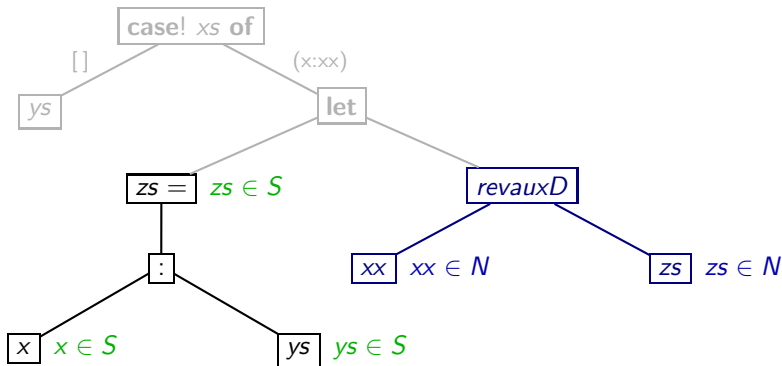
- Initially all parameter positions are marked as **don't-know**.
- Therefore, the actual parameters **zs** and **xx** in function application belong to the set N .

Example: *revauxD* - First iteration



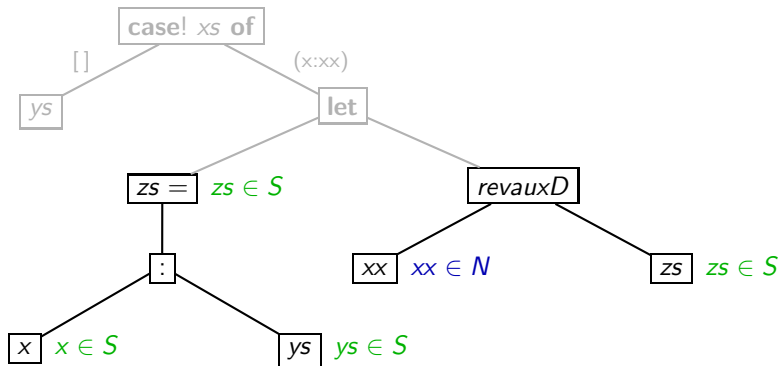
- Variables x and ys are used to construct a data structure. They are marked as **safe**.

Example: *revauxD* - First iteration



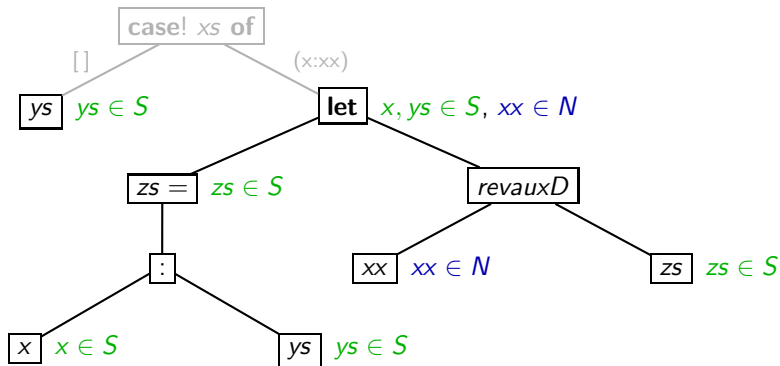
- From the information of the main expression of the let binding, we see that **zs** is not used destructively there. Hence it can be marked as **safe**.
- However, a top-down traversal of the function application is needed in order to check that the corresponding occurrence of **zs** may be marked as **safe**.

Example: *revauxD* - First iteration



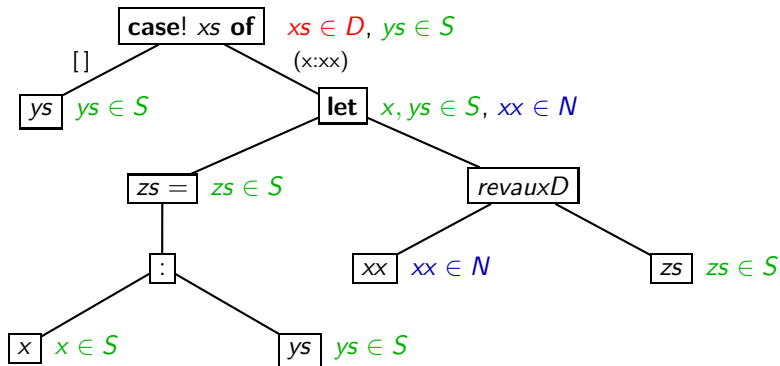
- From the information of the main expression of the let binding, we see that **zs** is not used destructively there. Hence it can be marked as **safe**.
- However, a top-down traversal of the function application is needed in order to check that the corresponding occurrence of **zs** may be marked as **safe**.

Example: *revauxD* - First iteration



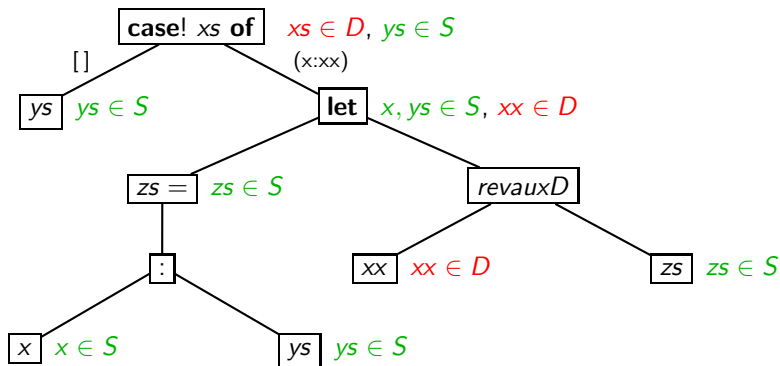
- In the **case!** alternative guarded by [], variable *ys* is marked as **safe**.
- Notice that the marks of *ys* are consistent in both alternatives and hence \sqcup is well-defined.

Example: *revauxD* - First iteration



- In the **case!** expression `xs` is inferred as **condemned**, since we are destroying its associated data structure.
- However, a top-down traversal of the second alternative is needed to check that `xx` may be marked as **condemned**, since it is a recursive child of `xs`.

Example: *revauxD* - First iteration



- The following sets are inferred:

$$D = \{xs\} \quad R = \{\} \quad S = \{ys\} \quad N = \{\}$$

- The type signature of *revAuxD* is updated: The first position is now **condemned** and the second one is **safe**.

Termination Analysis

- In **Safe**, heap memory consumption depends on the length of recursive calls chains.
- Being able to giving **bounds** to those chains length becomes essential in computing space bounds.
- We have investigated how to analyse termination and complexity bounds of Safe programs by using standard **term rewriting techniques**
- In particular, we have investigated how to use **termination proofs** which combine the **dependency pairs** approach, together with **polynomial interpretations**, to obtain suitable bounds to the length of chains of recursive calls in Safe programs.

Transformation from Core-SAFE to TRS

$$\begin{aligned}
trP(\overline{def_i^n}) &\stackrel{\text{def}}{=} \bigcup_{i=1}^n trF(def_i) \\
trF(f \overline{x_i^n} = e) &\stackrel{\text{def}}{=} f(x_1, \dots, x_n) \rightarrow trR(e, fv(e), []) \\
trR(c, V, C) &\stackrel{\text{def}}{=} c \Leftarrow C \\
trR(x, V, C) &\stackrel{\text{def}}{=} x \Leftarrow C \\
trR(C_r \overline{a_i^n}, V, C) &\stackrel{\text{def}}{=} C_r(a_1, \dots, a_n) \Leftarrow C \\
trR(f \overline{a_i^n}, V, C) &\stackrel{\text{def}}{=} f(a_1, \dots, a_n) \Leftarrow C \\
trR(k : \mathbf{case} \ x \ \mathbf{of} \ \overline{C_i \overline{x_{ij}^{n_i}} \rightarrow e_i^n}, V, C) &\stackrel{\text{def}}{=} \\
&\quad \{ \mathbf{case}_k(x, var(V)) \Leftarrow C \} \cup \\
&\quad \{ \mathbf{case}_k(C_i(x_{i1}, \dots, x_{in_i}), var(V)) \rightarrow trR(e_i, fv(e_i), []) \mid i \in \{1..n\} \} \\
trR(\mathbf{let} \ x_1 = e_1 \ \mathbf{in} \ e_2, V, C) &\stackrel{\text{def}}{=} trR(e_2, fv(e_2), C ++ [, trL(e_1, fv(e_1)) \rightarrow x_1]) \\
trL(e, V) &\stackrel{\text{def}}{=} trR(e, V, []) \text{ if } e \in \{c, x, C_r \overline{a_i^n}, f \overline{a_i^n}, \mathbf{case}\} \\
trL(\mathbf{let} \ x_1 = e_1 \ \mathbf{in} \ e_2, V) &\stackrel{\text{def}}{=} [trL(e_1, fv(e_1)) \rightarrow x_1,] ++ trL(e_2, fv(e_2))
\end{aligned}$$

Transformation of `splitD`

Then, we obtain a deterministic 3-CTRS \mathcal{R} . Let us call $U(\mathcal{R})$ to the TRS resulting from the standard transformation to unconditional TRSs:

```

splitD(n,xs) -> case1(n,n,xs)
case1(0,n,xs) -> Tup(Nil,xs)
case1(S(x),n,xs) -> case2(xs,n)
case2(Nil,n) -> Tup(Nil,Nil)
case2(Cons(x,xx),n) -> U1(n-1,x,xx)
U1(n',x,xx) -> U2(splitD(n',xx),x)
U2(z,x) -> U3(case3(z),z,x)
U3(xs1,z,x) -> U4(case4(z),x,xs1)
U4(xs2,x,xs1) -> U5(Cons(x,xs1),xs2)
U5(xs1',xs2) -> Tup(xs1',xs2)
case3(Tup(ys1,ys2)) -> ys1
case4(Tup(zs1,zs2)) -> zs2

```

The Dependency Pairs Approach

- Dependency pairs are rewrite rules obtained from a Term Rewriting System which show how function calls evolve in rewriting computations.

Definition

Given a TRS $\mathcal{R} = (\mathcal{C} \uplus \mathcal{D}, R)$ the set $DP(\mathcal{R})$ of *dependency pairs* for \mathcal{R} is given as follows: if $f(t_1, \dots, t_m) \rightarrow r \in R$ and $r = C[g(s_1, \dots, s_n)]$ for some defined symbol $g \in \mathcal{D}$, and context $C[\cdot]$, and $s_1, \dots, s_n \in \mathcal{T}(\mathcal{F}, \mathcal{X})$, then $f^\sharp(t_1, \dots, t_m) \rightarrow g^\sharp(s_1, \dots, s_n) \in DP(R)$, where f^\sharp and g^\sharp are new fresh symbols associated to f and g respectively.

- Dependency pairs can be presented as a (dependency) graph whose cycles represent *recursive calls*.
- Termination of rewriting computations can be analysed by considering the cycles in the dependency graph.
- This corresponds to the standard view of recursion as the source of any non-terminating behaviour.

Polynomial Bounds

Proposition (polynomial bounds)

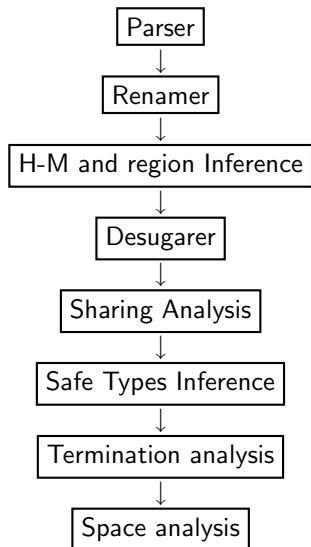
If $\llbracket f^\# \rrbracket$ is the polynomial interpreting the symbol $f^\#$ associated to a n -ary function symbol f in a Core-SAFE program and x_1, \dots, x_n are interpreted as the sizes of the input arguments to f , then the number $N(x_1, \dots, x_n)$ of recursive calls to f with arguments t_1, \dots, t_n of sizes x_1, \dots, x_n , respectively, is bounded by $\llbracket f^\# \rrbracket(x_1, \dots, x_n)$, i.e. $N(x_1, \dots, x_n) \leq \llbracket f^\# \rrbracket(x_1, \dots, x_n)$

Safe function	Polynomial inferred
$length(x)$	$x + 1$
$splitD(n, x)$	$n + x + 1$
$mergeD(x, y)$	$x + y + 1$
$msortD(x)$	No proof obtained
$insert(x, t)$	$t + 1$

Memory Consumption Analysis

- Once we have obtained (or manually annotated) a polynomial for every **Safe** function, we foresee that inferring **polynomial space bounds** will be relatively easy.
- Operational semantics will be enriched with space cost annotations both for heap and stack consumption, i.e. a **cost model** of the core language will be developed. Cell deallocation will give rise to negative consumption.
- Every output region and the **self** region will constitute an independent **cost centre**. Cost will be charged at construction and at function applications.
- In function applications we will make use of the previously inferred costs for the function. A function may charge costs to each output region and to its **self** region.
- The cost finally charged for a recursive call to every region will be multiplied by the polynomial representing the bound on the number of recursive calls of the current function.
- Least upper bounds will be computed at **case** expressions.

Implementation



- A complete front-end for SAFE has been developed by using standard tools (about 5,000 Haskell lines).
- The sugared language contains **data** declarations, infix operators, **where** clauses, guards, etc. The desugarer transforms it to core language.
- The Hindley-Milner infers the usual polymorphic types for functions and also the regions needed for each function. The safe types analysis makes use of the underlying H-M types.
- The **sharing analysis** decorates the program with sharing information from which *shareall* and *shar-erec* are extracted.
- Then, **safe types** are inferred and the AST is decorated with them.
- Finally **termination** and **space analyses** are done and the AST is decorated with polynomials for termination and space bounds.

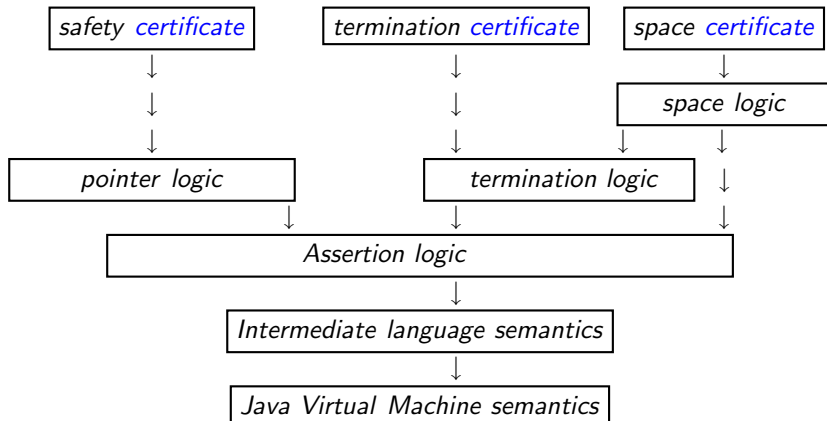
The Proof Carrying Code Paradigm

- There is a **producer** that generates code, and also a proof or **certificate** related to this code, asserting that certain properties are satisfied.
- The certificate may be generated by hand, automatically or by using a **proof assistant**. The paradigm allows any combination of those.
- There is a **recipient** who gets the code and the proof and **checks** that the proof is correct. The recipient only trusts in his/her **proof checker**.
- Once the code is proof-checked, the recipient runs it **safely** in his machine. This approach is superior to having runtime checks.
- For the paradigm to work, certain properties are desirable:
 - ❶ The code should be as standard as possible. **Java bytecode** or **C** are good candidates.
 - ❷ The certificate should be **specific** to this code.
 - ❸ Certificates should be as **small** as possible.
 - ❹ The proof checking should be not only automatic but also **efficient**.

The Proof Assistant Isabelle/HOL

- Isabelle is developed at University of Cambridge (Larry Paulson) and Technical University of Munich (Tobias Nipkow). See <http://isabelle.in.tum.de>
- Isabelle is a **generic proof assistant**. It allows mathematical formulas to be expressed in a formal language and provides tools for proving those formulas in a logical calculus.
- Provides support for nominal datatypes (binding structures) due to the HOL-Nominal logic.
- It has an inductive package for **inductive predicates** and **inductive sets**.
- Allows **general recursive function** definitions.
- Has a **code generator** for a subset of HOL, targeting SML, Haskell, and OCaml.
- Interactive successful sessions are recorded in **theory** files, which can be seen as real **proofs**, since they can be replayed when desired.

Our approach



- **Certificates** are high-level objects inferred by the compiler: **types**, **polynomials**
- The PCC infrastructure consists of layers of Isabelle/HOL **definitions** and **theorems** pushing up the level at which the proof is done.

A high-level assertion in pointer logic

Definition

We say that the Safe expression e satisfies the derived assertion $A = \llbracket L, \Gamma \rrbracket$ denoted $e : \llbracket L, \Gamma \rrbracket$ if:

① $L = fv(e)$ and $E \vdash h, e \Downarrow h', v$

② For every pair of variables $x \in dom(E), z \in L,$

$$\Gamma[z] = d \wedge recReach(E, z, h) \cap closure(E, x, h) \neq \emptyset \rightarrow x \in dom(\Gamma) \wedge \Gamma[x] \neq s$$

③ Whenever the initial configuration (Γ, E, h, L) is *good*, the final configuration (v, h') is also good. (Γ, E, h, L) is good whenever:

① $S \cap R = \emptyset$ where $S \stackrel{\text{def}}{=} \bigcup_{x \in L, \Gamma[x]=s} \{closure(E, x, h)\}$ and

$$R \stackrel{\text{def}}{=} \bigcup_{x \in L, \Gamma[x]=d} \{p \in live(E, L, h) \mid p \rightarrow_h^* recReach(E, x, h)\}.$$

② $closed(E, L, h)$

(v, h') is good if $closed(v, h')$.

Proof Rules in Pointer Logic

$$\frac{e_1 : \llbracket L_1, \Gamma_1 \rrbracket \quad e_2 : \llbracket L_2, \Gamma_2 + [x_1 : s] \rrbracket \quad \text{def}(\Gamma_1 \triangleright^{L_2} \Gamma_2)}{\text{let } x_1 = e_1 \text{ in } e_2 : \llbracket L_1 \cup (L_2 - \{x_1\}), \Gamma_1 \triangleright^{L_2} \Gamma_2 \rrbracket} \text{LET1}$$

$$\frac{e_1 : \llbracket L_1, \Gamma_1 \rrbracket \quad e_2 : \llbracket L_2, \Gamma_2 + [x_1 : d] \rrbracket \quad \text{def}(\Gamma_1 \triangleright^{L_2} \Gamma_2)}{\text{let } x_1 = e_1 \text{ in } e_2 : \llbracket L_1 \cup (L_2 - \{x_1\}), \Gamma_1 \triangleright^{L_2} \Gamma_2 \rrbracket} \text{LET2}$$

$$\frac{\forall i. e_i : \llbracket L_i, \Gamma_i \rrbracket \wedge x \in L_i \quad \Gamma = \bigotimes_i (\Gamma_i / \{\bar{x}_{ij}\}) \quad L = \bigcup_i (L_i - \{\bar{x}_{ij}\})}{\text{case } x \text{ of } \overline{C_i \bar{x}_{ij}} \rightarrow e_i : \llbracket L, \Gamma \rrbracket} \text{CASE}$$

Conclusion

- PCC is an emerging technology that can dramatically increase our confidence in programs, both in **safety critical applications** and in frameworks where **untrusted code** must be accepted.
- The combination of powerful **static analyses** at source level with **automatic generation of certificates** simplifies life at the producer's side.
- Writing and analysing programs in a **declarative language** simplifies both program production and compiler construction.
- Producing proofs at **bytecode level** would imply a large amount of work. Moreover, certificates and checking time **will increase dramatically**.
- Developing layers of **specialised logics** is also a large amount of work but has the advantage that complex theorems are **proved only once**.