



RAPPORT DE STAGE D'OPTION SCIENTIFIQUE

Titre

NON CONFIDENTIEL

Option :	INFORMATIQUE
Champ de l'option :	Math-Informatique
Directeur de l'option :	Olivier Bournez
Directeur de stage :	Olivier Bournez
Dates du stage :	7 avril - 22 août 2014
Nom et adresse de l'organisme :	SRI International Computer Science Laboratory (CSL) 333 Ravenswood Avenue Menlo Park, CA 94025-3493 United States

July 12, 2014

Contents

1	Introduction	2
2	SRI	2
2.1	The HACMS Project	2
2.2	PVS	2
2.3	Translating PVS	2
2.3.1	Parsing and typechecking PVS	2
2.3.2	Other translator	2
3	Translating PVS Syntax	2
3.1	PVS Syntax	2
3.2	Translator architecture	2
3.3	A few translation rules	5
4	Types	6
4.1	PVS Types	6
4.2	Translating types	6
4.3	Translating PVS syntax	7
5	Difficulties and successes	8
5.1	if expressions	8
5.2	Integer, rationnals	8
5.3	Garbage collection	8
5.4	Performance	9
5.5	Update expressions	9
5.5.1	Pointer counting	9
5.5.2	Using a different data structure	10
5.5.3	Flow analysis on the PVS code	10
5.5.4	Analysis of the C code	10
5.5.5	Algorithm	12
5.5.6	Combination of solutions	12
6	Conclusion	13
6.1	What's left to be done ?	13
6.2	My stay at SRI	13

1 Introduction

2 SRI

2.1 The HACMS Project

2.2 PVS

2.3 Translating PVS

2.3.1 Parsing and typechecking PVS

These two task we leave to PVS native parser and typechecker.

The parser generates objects representing the expressions of the theory.

We only convert a subset of PVS. This subset is defined by a subset of expression objects we can translate. The objective is, of course, to be able to translate the maximum of (if not all) PVS expression objects.

2.3.2 Other translator

- Common Lisp (native) - Clean - Yices

3 Translating PVS Syntax

3.1 PVS Syntax

We describe here the syntax of PVS and the objects system used to represent them in Lisp. Some slots of the classes are voluntarily omitted. For a full description of PVS parser representation, refer to [5].

3.2 Translator architecture

Describe here the Lisp functions and data structures

Skeleton

Expected input

Output objects

Assertions that we (try to) maintain

Main steps :

- Typechecking: The PVS typechecker perform a type analysis on the PVS code to associate a PVS type to each expression. This might generates some proof obligations (TCC).
- Lexical and syntactic analysis: The PVS parser transforms PVS code into a Lisp internal representation.
- Translation: The translator generates a different representation from PVS expressions and functions declarations. Typically, an expression e is translated into a tuple of four elements (t, n, i, d) , where t represents a C type used to describe the expression, n is a string representing the expression, i is a list of instructions supposed to be executed prior to using n (initialisation of n) and d is a list of instructions to be executed when n isn't needed anymore (destruction of n).
- Optimizations

<i>Expr</i>	::=	<i>Number</i> <i>Name</i> <i>Expr Arguments</i> <i>Expr Binop Expr</i> <i>Unaryop Expr</i> <i>Expr</i> ‘ { <i>Id</i> <i>Number</i> } (<i>Expr</i> ⁺) (# <i>Assignment</i> ⁺ , #) <i>IfExpr</i> LET <i>LetBinding</i> ⁺ IN <i>Expr</i> <i>Expr</i> WHERE <i>LetBinding</i> ⁺ <i>Expr</i> WITH [<i>Assignment</i> ⁺ ,]
<i>Number</i>	::=	<i>Digit</i> ⁺
<i>Id</i>	::=	<i>Letter IdChar</i> ⁺
<i>IdChar</i>	::=	<i>Letter</i> <i>Digit</i>
<i>Letter</i>	::=	A ... Z
<i>Digit</i>	::=	0 ... 9
<i>Arguments</i>	::=	(<i>Expr</i> ⁺)
<i>IfExpr</i>	::=	IF <i>Expr</i> THEN <i>Expr</i> { ELSIF <i>Expr</i> THEN <i>Expr</i> } * ELSE <i>Expr</i> ENDIF
<i>Name</i>	::=	true false number_field_pred real_pred integer_pred integer? rational_pred floor ceiling rem ndiv even? odd? cons car cdr cons? null null? restrict length member nth append reverse
<i>Binop</i>	::=	= \= OR \/ AND & /\ IMPLIES => WHEN IFF <=> + - * / < <= > >=
<i>Unaryop</i>	::=	NOT -
<i>Assignment</i>	::=	<i>AssignArg</i> ⁺ { := -> } <i>Expr</i>
<i>AssignArg</i>	::=	(<i>Expr</i> ⁺) ‘ <i>Id</i> ‘ <i>Number</i>
<i>LetBinding</i>	::=	{ <i>LetBind</i> (<i>LetBind</i> ⁺) } = <i>Expr</i>
<i>LetBind</i>	::=	<i>Id</i> [: <i>TypeExpr</i>]

Figure 1: Syntax of the PVS subset of the translator

expr \subset syntax [<i>abstract class</i>] <i>type</i> the type of the expression
name \subset syntax [<i>mixin class</i>] <i>id</i> the identifier <i>actuals</i> a list of actual parameters <i>resolutions</i> singleton This is a mixin for names, i.e., name-exprs , type-names , etc.
name-expr \subset name expr [<i>class</i>]
number-expr \subset expr [<i>class</i>] <i>number</i> a nonnegative integer
tuple-expr \subset expr [<i>class</i>] <i>exprs</i> a list of expressions
application \subset expr [<i>class</i>] <i>operator</i> an expr <i>argument</i> an expr (maybe a tuple-expr)
field-application \subset expr [<i>class</i>] <i>id</i> identifier <i>actuals</i> . a list of actuals <i>argument</i> the argument A field application is the internal representation for record extraction, e.g., r ‘a
record-expr \subset expr [<i>class</i>] <i>assignments</i> a list of assignments
lambda-expr \subset binding-expr [<i>class</i>] This is the subclass of binding-expr used for LAMBDA expressions.
if-expr \subset application [<i>class</i>] When an application has an operator that resolves to the if_def it is changed to this class.
update-expr \subset expr [<i>class</i>] <i>expression</i> . an expr <i>assignments</i> a list of assignments An update expression of the form e WITH [x := 1, y := 2] , maps to an update-expr instance, where the expression is e, and the assignments slot is set to the list of generated assignment instances.
assignment \subset syntax [<i>class</i>] <i>arguments</i> . the list of arguments <i>expression</i> the value expression Assignments occur in both record-exprs and update-exprs. The arguments form is a list of lists. For example, given the assignment ‘a(x, y)‘1 := 0, the arguments are ((a) (x y) (1)) and the expression is 0.

Figure 2: (Partial) CLOS representation of PVS expressions

- Code generation: C code is generated.

We first define a function T to translate an expression e .

$$T(e) = (T^t(e) , T^n(e) , T^i(e) , T^d(e))$$

```

T(2) = ( int,"2", [], [] )
T(4294967296) = ( mpz_t, ?,
                  [mpz_init(?); |
                   mpz_set_str(?, "4294967296");],
                  [mpz_clear(?);])
T(lambda(x:below(10)):x) = ( int*, ?,
                             [? = malloc(10 * sizeof(int)); |
                              int i;|
                              for(i = 0; i < 10; i++)
                                ?[i] = i;]
                             [free(?);])

```

Figure 3: Translation examples: number expressions

It may occur that $T^n(e) = ?$. In that case, the symbol $?$ appearing in $T^i(e)$ and $T^d(e)$ needs to be replaced by a proper variable name.

We then define two other operators:

- R which take an expression and a type and may add an extra conversion in the instructions to make sure its result has the expected type. Also the result of this function has a proper name.
- S which take an expression, a type and a name. It makes sure that the given variable (type + name) is set to a value representing the expression.

3.3 A few translation rules

Translation rules :

```

number-expr "2"
(C-int, "2", [], [])

number-expr "12315468453213"
(C-mpz, nil,
  [mpz_t ~a; | mpz_t_init("12315468453213"); ],
  [mpz_clear ~a;])

application "f(e1, e2)"
(C-mpz, nil,
  [ instr(e1) | instr(e2)
    | mpz(~a); | f(~a, e1, e2) ]
  [mpz_clear(~a);])

```

<i>TypeExpr</i>	::=	<i>Name</i> <i>EnumerationType</i> <i>Subtype</i> <i>TypeApplication</i> <i>FunctionType</i> <i>TupleType</i> <i>CotupleType</i> <i>RecordType</i>
<i>EnumerationType</i>	::=	{ <i>IdOps</i> }
<i>Subtype</i>	::=	{ <i>SetBindings</i> <i>Expr</i> } (<i>Expr</i>)
<i>TypeApplication</i>	::=	<i>Name Arguments</i>
<i>FunctionType</i>	::=	[FUNCTION ARRAY] [-[<i>IdOp</i> :] <i>TypeExpr</i> ⁺ , -> <i>TypeExpr</i>]
<i>TupleType</i>	::=	[-[<i>IdOp</i> :] <i>TypeExpr</i> ⁺]
<i>CotupleType</i>	::=	[-[<i>IdOp</i> :] <i>TypeExpr</i> ⁺ ₊]
<i>RecordType</i>	::=	[# <i>FieldDecls</i> ⁺ , #]
<i>FieldDecls</i>	::=	<i>Ids</i> : <i>TypeExpr</i>

Figure 4: Fragment of the PVS type system

4 Types

4.1 PVS Types

A PVS theory can be typechecked using the emacs interface `M-x typecheck` or with Lisp function `(tc name-theory)`. This first runs the PVS parser on the code and generates CLOS objects to represent it. Then, the PVS typechecker is run on this internal representation of the theory and tries to give a type to all expressions generating TCC when needed.

Here we describe how PVS types are represented in Lisp.

4.2 Translating types

PVS types:boolean, number, number_field, real, rational, integer, $A \rightarrow B$, restricted types below(10) := { $x : \text{int} | 0 \leq x < 10$ }) enum datatype

This requires a type analysis to decide on the type of a PVS expression. For example the PVS `int` type can be represented by the `int`, `unsigned long` or `mpz_t` C types. In that case, we study the range of the expression to decide which types are allowed to represent it. Then we take the context in which the expression appears to decide. For instance in

```
incr(x:below(10)):int = x+1
```

the `x` expression, result of the function `incr` can always be represented by an `int` or `unsigned long` in C but we choose here to represent it using a `mpz_t`.

Auxiliary type system : C-type with a flag : mutable (meaning that the expression it describes only has one pointer pointing to it.

```
int a = 2;      a : int[mutable]
int* a = malloc( 10 * sizeof(int*) );
```

destructive addition:

type-expr \subset syntax	[<i>abstract class</i>]
.....	
type-name \subset type-expr name <i>adt</i>	[<i>class</i>]
.....	
subtype \subset type-expr <i>supertype</i> <i>predicate</i>	[<i>class</i>]
.....	
funtype \subset type-expr <i>domain</i> <i>range.</i>	[<i>class</i>]
.....	
tupletype \subset type-expr <i>types</i>	[<i>class</i>]
.....	
recordtype \subset type-expr <i>fields</i>	[<i>class</i>]
.....	

Figure 5: (Partial) CLOS representation of PVS types

```
d_add(*mpz_t res, mpz_t[mutable] a, long b) {
  mpz_add(a, a, b);
  (*res) = a;
}
```

Rq : `d_add` is given a mutable `mpz_t`, meaning that it can modify it and is responsible for freeing it. It is also responsible for allocating memory for the result. Here it uses the memory to assign `res`.

Use an auxiliary language :

```
( expr, C-type[mutable] )
```

Conversions and copies create mutables types (at a cost) : `a[mutable]_from_b`

[2]

C types:[3]

```
// integer and floating point types
[unsigned] char, int, long, double
type* //arrays
char* // strings
struct types // structures with fields
enum types
short int, float, union, size_t // etc...
```

Listing 1: C types

Translation rules :

We can only translate a subset of all PVS types. What's missing ?

4.3 Translating PVS syntax

We can only translate a subset of PVS syntax. What's missing ?

subrange(a, b)	<code>int</code> <i>// if small enough</i> <code>unsigned long</code> <i>// if too big or needed for function call</i> <code>mpz_t</code> <i>// else</i>
int	<code>mpz_t</code>
rat	<code>mpq_t</code>
[below(a) -> Type]	<code>(Ctype)*</code>
T : TYPE = [# x _i : t _i #]	<code>struct CT {</code> <code>...</code> <code>Ct_i x_i;</code> <code>...</code> <code>};</code> <i>// These types must be declared</i>
[Range -> Domain]	C closure parameterized by the Domain return type.

Figure 6: Translation rules for PVS types

5 Difficulties and successes

5.1 if expressions

Represented by if-expr

5.2 Integer, rationals

In PVS, the `integer` represent the whole set \mathbb{Z} of all relative numbers (and `rational` also describe \mathbb{Q}). In C, we have finite types `int`, `long`, ...

We need the GMP library which introduces the types `mpz_t` and `mpq_t`. These types are arrays and should be used just as integer (not as pointers except they still need to be freed).

5.3 Garbage collection

We implement a very simple "Reference Counting Garbage Collector" as described in [4].

We maintain a hashtable of pointer counters. Each pointer in the code is a key in the hashtable to which we associate an `int` counter as value.

Pointers only occurs in arrays or struct.

Arrays are created in the code.

```
T* a = b;
```

becomes

```
T* a = GC( b ); // Should not happen often...
```

```
t[0] = b; // with b of type T*
```

becomes

```
GC_free( t[0] );
t[0] = (T*) GC( b );
```

Examples

```
int* f() {
    int* res;
    res = (int*) GC_alloc( 10 * sizeof(int) );
    [... init res ...]
    return res; // pointer count = 1
}
```

```

void main() {
    int* a = f(); // pointer counter of a = 1
    int** b = (int**) GC_alloc( sizeof( int* ) );
    GC_free( b[0] ); // useless
    b[0] = (int*) GC( a ); // pointer counter of a = 2
    printf("f(0) = %s", b[0][0]);
    GC_free(b); // frees b, pointer count of a = 1
    GC_free(a); // frees a
}

```

5.4 Performance

5.5 Update expressions

Update expressions are represented by PVS as `update-expr` objects.

$$E := t \text{ with } [e1 := e2]$$

Problem : t is an expression typed as a function. Therefore it might be represented in C as an array (if domain type is `below(n)`). We want to know if we can update t in place to obtain a C object representing E or if we have to make a copy of t .

We consider three solutions to this problem.

5.5.1 Pointer counting

We keep track of the number of pointer pointing to an array or a struct.

This requires to build our own C struct (heavy)

```

struct array_int {
    int pointer_count = 1;
    int *data;
};

```

When we update the struct, if the pointer is 0, we update in place.

Besides every update require now to read the structure and make a test (small compared to a copy but no so small compared to a single in place update)

Besides, the creation / destruction gets more complicated

Passing argument to function :

```

array_int f(array_int arg) {
    arg.pointer_count++; // Since now f also have a pointer to the struct
    if (arg.pointer_count == 1) {
        arg.data[0] = 0;
        return arg;
    } else {
        array_int res;
        res.data = malloc( 10 * sizeof(int*) );
        copy(res, arg); // Very long...
        res.pointer_count--; // This function is about to lose its pointer to res
        return res;
    }
}

void main() {
    array_int t;
    init(t); // somehow...

    t.pointer_counter--; // We assure we won't use the pointer "t" to the array anymore
}

```

```

    array_int r = f(t);
    t = null; // This way we guarantee the variable "t" won't be used later in the code

    [...]
}

```

This add quite some code compared to the simple :

```

array_int f(array_int arg) {
    arg[0] = 0;
    return arg;
}

void main() {
    array_int t;
    init(t); // somehow...
    array_int r = f(t);

    [...]
}

```

5.5.2 Using a different data structure

PVS uses arrays in a very particular way, we might then represent them with an other structure than just only a C array. For example :

```

struct r_list_int {
    int key;
    int value;
    r_list_int tl;
};

struct array_int {
    int *data;
    r_list_int replacement_list;
};

```

Each structure represent the array `data` with the modifications contained in the linked list `r_list_int`

Problems : Just as the previous solution : - add some extra code - add some extra computation (runtime tests, reading the replacement list) - require to create as many structures and associated functions as there are range types for the manipulated arrays - Very dependent on the GC

5.5.3 Flow analysis on the PVS code

An other optimization would be to perform a analysis on the PVS variables to make sure an update Pavol [1] suggests three analysis

5.5.4 Analysis of the C code

Trying to avoid copying arrays by analyzing the C code generated.

2 different functions (destructive and non destructive)

Algorithm :

Always have a "non destructive" version of any function. A "cautious" version that never modify the arguments in place and always make copies when necessary (when a "mutable" version of an array is necessary (for instance updates)).

In destructive versions of all functions : Flag all array arguments to "mutable". Then for each of these arguments : - If it never occurs destructively, then remove flag (function just read the arg) - If it occurs destructively, it can never occur at all AFTER. - Need to define the order of

evaluation of expression (easy rules on simple expressions) -; Need to be able to detect occurrences of a name-expr -; Otherwise, unflag the arg

A variable V of type array is created in these cases:

- $V = \lambda x. e(x) : V$ has bang type
- $\text{update}(V, T, \text{key}, \text{value}) : V$ has bang type because this is basically a copy and a destructive update.
- $f(V, \dots) = \dots : \text{type of } V \text{ depends on } f.$

In these case, it has always bang type. Or it can be set to an other referenced object.

- $V = T \rightarrow V$ (should have bang type iff T has !type too and never occurs afterwards). Happens in
- $V = T[i] \rightarrow$ depends on the target type of T .
- $V = T.\text{field} \rightarrow$ depends on the type of the field.

At first all updates are non destructive.

First pass : All array variables (actuals and local variables) found in the code are flagged. Local variables are flagged according to the previous rules and actuals are flagged **mutable** in destructive version and **not mutable** in non-destructive versions. In functions returning an array (or record type), the variable result is also flagged.

Other passes : Reading the code backwards, for every occurrence T of a variable flagged **mutable**:

- If it is found in a $V = T$ instruction, then we give the bang type to V and remove bang type from T so that previous occurrences of T won't assume the uniqueness of the reference. This adds a new variable to the set of bang variables, hence the need to make several passes.
- If it is used in a $V = \text{copy}(T)$ instruction, then we replace it with a $V = T$ instruction and do as previous.
- If it is found in an $\text{update}(V, T, i, e)$, then turn that into $V = T; \text{destr_update}(V, i, e)$.
- If it is a function call
- If we reach the declaration of a variable that is marked **mutable**, this means this variable is never read. In that case, we actually don't need it (unflag it I guess...).

At the end, when we have reached the transitive closure of this definition, if we reach the beginning of the function and some arguments are still bang, this means their bangness is never used, put the flag on that argument to **non mutable** and remove the instructions freeing that variable (reminder : mutable arguments of a functions are freed inside the function or are used in a mutable way and appear somewhere in the result (trapped in closures) or are freed in other function calls.

5.5.5 Algorithm

Initialization: We initialize a set M of mutable variables to all array arguments of a function f . We also initialize a set F of variables to free to M since f has the responsibility to free all variables flagged as mutable arguments.

We read the code backwards. T_i refer to variables that are in the set M . S_i refer to variables that are not in the set M .

$S = T$	$M \leftarrow M \cup \{S\} - \{T\}$ $F \leftarrow F \cup \{S\} - \{T\}$
$S = \text{update}(S_2, \text{key}, \text{value})$	$M \leftarrow M \cup \{S\}$ $F \leftarrow F \cup \{S\}$
$S = \text{update}(T, \text{key}, \text{value})$	$M \leftarrow M \cup \{S\} - \{T\}$ $F \leftarrow F \cup \{S\}$
$S = g(T_i, S_i)$	If the arguments of g don't allow g to be called destructively: $M \leftarrow M - \{T_i\}$ $M \leftarrow M \cup \{S\}$ if return type of g is mutable $F \leftarrow F \cup \{S\}$
$S = g(T_i, S_i)$	Otherwise: $\rightarrow S = g_d(T_i, S_i)$ $M \leftarrow M \cup \{S\} - \{T\}$ $F \leftarrow F \cup \{S\}$
$S = S_2[i]$	$M \leftarrow M \cup \{S\} - \{T\}$ $F \leftarrow F \cup \{S\}$

All arguments of the function are flagged **mutable**

What is a destructive occurrence :

$$E := f(t \text{ with } [e1 := e2] , t(0))$$

order of eval : $e1$ and $e2$ (t can occur non destr) t (expression of an update : destr) $t(0)$ (occurrence of t (even non destr))

$f(x:\text{Arr}):int = g(h(t), t)$ is destructively translated to

```

1  int f_d(int* t) {    // t has type ! since this is destructive f
2    int* arg1 = h(t); // h can't be called destructively because
3                      // even though t is !, it appears later (line 4)
4    int* arg2 = t;     // t is ! and never appears later => arg2 is !
5    return g( arg1, arg2); // arg2 is ! but g can only be called
6  }                   // destructively if arg1 is

```

Listing 2: Example

if g has type $[\text{Array!} \rightarrow ?]$ then t can't be destructive

if g has type $[\text{Array} \rightarrow ?]$ then t can be destructive

First algorithm:

Need multiple passes as the flags disappear

5.5.6 Combination of solutions

We use the C code analysis to write some updates as destructive. However a few updates remain non destructive. For example:

If a function is called but requires its two argument to be **mutable** and only the first is **mutable**. Then the non-destructive version is called and the first argument gets copied even though it was **mutable**.

If we perform an update on $T[i]$, our analysis doesn't tell if $T[i]$ is **mutable** or **non-mutable**.

update(A, key, value)	A[key] = value;	A must be mutable
set(A, <i>expr</i>)	A = <i>expr</i> ;	
declare(A, <i>expr</i> (<i>i</i>))	<pre> A = malloc(1 * sizeof(T)); int i; for(i = 0; i < 1; i++) A[i] = i + 1; </pre>	
copy(A, B)	<pre> A = malloc(1 * sizeof(T)); int i; for(i = 0; i < 1; i++) A[i] = B[i]; </pre>	
init(A)	int* A;	
free(A)	free(A);	
base(str, (A, B, ...))	int aux = A[0] + B[1];	A and B are only read.
return	return result;	

Figure 7: C instructions

value(<i>cste</i>)	42	
variable(<i>type</i> , <i>name</i>)	name	
call(f, <i>exprs</i>)	f(<i>expr</i> ₁ , ..., <i>expr</i> _{<i>n</i>})	

Figure 8: C expressions

To prevent that, we also perform a GC check. An update is actually a test whether an object is **mutable** or not and the appropriate update.

6 Conclusion

6.1 What's left to be done ?

Use a C structure to represent a closure

```

struct r_list_int {
    int (*body)(void* env, void* args);
    void* env;
    void* args;
};

```

6.2 My stay at SRI

Besides the conception and implementation of the PVS to C compiler, my stay at SRI International was rich in ??? events.

The first days of my stay were the occasion to discover PVS and Coq as I started working on a translator Coq to PVS. With Robin, we also wrote as an exercise a basic linear algebra library.

I discovered Lisp the hard way while learning how the back end of PVS worked. I've wrote a Lisp parser to help me see clear in the huge code (classes definitions, inheritances and organization, function dependances, ...)

I've had the chance to attend to many interesting seminars. The SRI also organized a Summer School to which we were allowed to attend and which was very interesting.

Shankar never hesitated to include us in many project

I've been included in the HACMS project which was very interesting. With other: Correcting translator PVS to SMT-LIB

Draft

Discovering PVS : Translating Coq proofs to PVS PVS library for basic linear algebra

Robin project, HACMS Contest week-end 14-15 June Summer School Parsing Lisp code -i generate HTML architecture fileCorrecting translator PVS to SMT-LIB [1]

References

- [1] Pavol Černý. Static analyses for guarded optimizations of high level languages.
- [2] Jean-Christophe Filliâtre. Langages de programmation et compilation. MPRI, 2013-14. <https://www.lri.fr/~filliatr/ens/compil/>.
- [3] Eric Huss. The c library reference guide. *Webmonkeys: A Special Interest*, 2004.
- [4] Richard Jones and Rafael Lins. Garbage collection: Algorithms for automatic dynamic memory management, 1996. *John Wiley & Sons Ltd., England*.
- [5] N. Shankar and S. Owre. *PVS API Reference*. Computer Science Laboratory, SRI International, Menlo Park, CA, September 2003.