

A CASE STUDY ON A NUMERICAL CODE IN C: LOG-SUM-EXP

PAUL BONNOT

JOINT WORK WITH RAPHAËL RIEU-HELFT (TRUSTINSOFT)
AND CLAUDE MARCHÉ (INRIA-SACLAY&LMF)

NOVEMBER 21ST, 2022

- Goal: formally verify two C functions defined on floating-point numbers : LSE and MLSE.
 - ▶ Absence of overflow
 - ▶ Bounds on the rounding errors w.r.t. to pure real computations
- Results obtained:
 - ▶ Formal proof of LSE in WhyML and in C.
 - ▶ Formal proof of MLSE in WhyML and (partly) in C.

LSE FUNCTION

```
double lse(double a[], size_t size) {  
    int i;  
    double s = 0.0;  
    for (i = 0; i < size; i++) {  
        s += exp_approx(a[i]);  
    }  
    return log_approx(s);  
}
```

Let a be a vector of size n . Then

$$LSE(a) = \log \sum_{i=0}^{n-1} \exp(a_i)$$

Objective: bound the error of FP computation compared to the real computation with errors A and B depending on input size.

$$\left| \widehat{LSE}(a, n) - LSE(a, n) \right| \leq A |LSE(a, n)| + B$$

MLSE FUNCTION

Let a be a vector of n numbers and x a number between 0 and 1

$$MLSE(a, n, x) = \log(n) + \frac{x^2}{\log(4)} + \sum_{i=1}^n -\log \left(\sum_{j=1}^n \exp \left(\frac{-(a_i + x - a_j)^2}{2} \right) \right)$$

We look for values A' and B' such that:

$$\left| \widehat{MLSE}(a, n, x) - MLSE(a, n, x) \right| \leq A' |MLSE(a, n, x)| + B'$$

THE TOOLS USED

ACSL: Formal specification of C code

Example [Boldo & Marché, 2011] :

```
1  /*@ requires  $|x| \leq 2^{-5}$ ; */
2  /*@ ensures  $|\text{result} - \cos(x)| \leq 2^{-23}$ ; */
3  float my_cos(float x) {
4      /*@ assert  $\left|1 - \frac{x^2}{2} - \cos(x)\right| \leq 2^{-24}$ ; */
5      return 1 - x * x * 0.5;
6  }
```

- From C to WhyML : *TrustInSoft Analyzer* with the J^3 plug-in
- From WhyML to SMT : *Why3*
- Solvers :
 - ▶ Generalistic SMT solvers: *Alt-Ergo*, *CVC4*, *CVC5*
 - ▶ Specialized solver: *Gappa*, for rounding errors, for absence of overflows
 - ▶ Specialized solver: *DReal*, for inequalities over the reals

THE APPROACH USED

1. Analysis of the summation of vectors
 - 1.1 On paper
 - 1.2 Proof formalization in WhyML
2. LSE analysis
 - 2.1 On paper
 - 2.2 Proof formalization in WhyML
3. MLSE analysis
 - 3.1 On paper
 - 3.2 Proof formalization in WhyML
4. Formal proof on C code with ACSL annotations

ROUNDING ERRORS ON SUMMATION OF VECTORS

ROUNDING ERRORS ON SUMMATION OF VECTORS

SUM OF DOUBLES, ON PAPER

REMINDER

For rounding mode *RNE* and format `double`, for all x :

$$|\text{round}(x) - x| \leq |x|\varepsilon + \eta$$

with $\varepsilon = 2^{-53}$ and $\eta = 2^{-1075}$.

We define the sum of floating point numbers like this :

- $x \oplus y = \text{round}(x + y)$
- $\bigoplus_{i=m}^n a_i = 0$ if $n \leq m$;
- $\bigoplus_{i=m}^n a_i = (\bigoplus_{i=m}^{n-1} a_i) \oplus a_n$ if $m < n$.

Better bound for addition:

$$|(x \oplus y) - (x + y)| \leq |x + y|\varepsilon$$

SUM OF DOUBLES

Result 1

Assuming $n \leq \frac{1}{2\varepsilon}$ and that for all i , $|a_i| \leq 2^{970}$:

$$\left| \bigoplus_{i=1}^n a_i - \sum_{i=1}^n a_i \right| \leq 2\varepsilon n \sum_{i=1}^n |a_i|$$

Absolute error

$\sum_{i=1}^n |a_i|$ instead of $|\sum_{i=1}^n a_i|$

SUM OF DOUBLES

We first prove an intermediate bound by induction over n :

$$\left| \bigoplus_{i=1}^n a_i - \sum_{i=1}^n a_i \right| \leq \sum_{i=1}^n |a_i| (\varepsilon n + \varepsilon^2 n^2)$$

We then use the fact that $n \leq \frac{1}{2\varepsilon}$ to prove the final bound

No overflow on inductive case : $(\bigoplus_{i=m}^{n-1} a_i) \oplus a_n \leq \text{maxFloat}$

Lemma

For all $c \geq 0$ and for all vector a such as $|a_i| \leq c$, we have:

$$nc(-1 - 2\varepsilon n) \leq \bigoplus_{i=1}^n a_i \leq nc(1 + 2\varepsilon n)$$

With $\varepsilon = 2^{-53}$, $n \leq 2^{51}$ and $c \leq 2^{970}$ we don't have an overflow.

ROUNDING ERRORS ON SUMMATION OF VECTORS

WHYML FORMALIZATION

WHYML FORMALIZATION

```
1  let rec function sum_double (f:int -> double) (a b:int)
2  =
3  if (b <= a) then
4      0
5  else
6      (sum_double f a (b - 1))  $\oplus$  f (b - 1)
```

We use, from WhyML IEEE-float library:

- type `double`, which includes values for infinities and NaN
- \oplus for addition, without check for overflow

WHYML FORMALIZATION

```
1  constant max_val = 2970
2  constant max_size = 251
3
4  let ghost sum_double_err (f:int -> double) (a b:int)
5    variant { b - a }
6    requires { 0 ≤ b - a ≤ max_size ∧
7              ∀i. a ≤ i < b → |f i| ≤ max_val }
8    ensures {
9      |result - sum_real f a b| ≤
10      (sum_real (fun i -> |f i|) a b) × (ε(b-a) + ε2(b-a)2)
11    }
12 = ...
```

Formalizes the result on the bounds, with a pre-condition for absence of overflow

WHYML FORMALIZATION: THE PROOF

```
1 let ghost sum_double_err (f:int -> double) (a b:int)
2   ...
3 = let s = sum_double f a (b - 1) in
4   assert IH {
5     |s - sum_real a (b - 1)| ≤
6     (sum_real (fun i -> |f i|) a (b-1)) × ε(b-a-1)
7     + ε2(b-a-1)2)
8   };
9   sum_real_bounds (-max_val) max_val f a (b-1);
10  sum_real_bounds 0 max_val (fun i -> |f i|) a (b-1);
11  assert s_bound { |s| ≤ 21023 };
12  let s' = s ⊕ f (b-1) in ...
```

- Assertion IH: inductive hypothesis
- Calls to ghost function `sum_real_bounds`: put a constant bound on `sum_real`
- Assertion `s_bound` proves the absence of overflows

ERROR BOUNDS ON LSE

ERROR BOUNDS ON LSE

ON PAPER

CONTRACTS OF FLOATING-POINT VERSIONS OF \exp AND \log IN WHYML

Notations: E_{\log} (resp. E_{\exp}) relative error of $\widehat{\log}$ (resp. $\widehat{\exp}$)

```
constant exp_error:real
```

```
axiom exp_error_bounds =  $2^{-53} \leq \text{exp\_error} \leq 2^{-7}$  (* was  $2^{-40}$  *)
```

```
constant exp_max:real = 672 (* was 25 *)
```

```
val function exp_approx (x:double) : double
```

```
  requires {  $|x| \leq \text{exp\_max}$  }
```

```
  ensures {  $|\text{result} - \exp(x)| \leq \text{exp\_error} \times \exp(x)$  }
```

```
constant log_error:real
```

```
constant max_size:real =  $2^{51}$  (* was 100 *)
```

```
axiom log_error_bounds =  $2^{-53} \leq \text{log\_error} \leq 2^{-2}$  (* was  $2^{-36}$  *)
```

```
val function log_approx (x:double) : double
```

```
  requires {  $0 < x \leq 2 \times \text{max\_size} \times (1 + \text{exp\_error}) \times \exp(\text{max\_val})$  }
```

```
  ensures {  $|\text{result} - \log(x)| \leq \text{log\_error} \times |\log(x)|$  }
```

ERROR BOUND ON LSE

$\widehat{\text{LSE}}(a, n)$: the floating-point implementation of LSE

Result 2

Assuming that $n \leq 2^{51}$ and for all i , $|a_i| \leq 672$:

$$\begin{aligned} & \left| \widehat{\text{LSE}}(a, n) - \text{LSE}(a, n) \right| \\ & \leq E_{\log} |\text{LSE}(a, n)| - \log(1 - (E_{\text{exp}} + 2\epsilon n(1 + E_{\text{exp}}))) \times (1 + E_{\log}) \end{aligned}$$

Proof: combining rounding errors of float sum and these of $\widehat{\text{exp}}$ and $\widehat{\text{log}}$

The relative error is E_{\log} .

We can approximate the constant error to get

$$2(E_{\text{exp}} + 2\epsilon n) \times (1 + E_{\text{exp}})$$

ERROR BOUNDS ON LSE

WHYML FORMALIZATION

WHYML FORMALIZATION OF LSE

```
1 let lse (a:array double) : double
2   requires {  $\forall i. 0 \leq i < a.length \rightarrow |a[i]| \leq 672$  }
3   requires {  $0 < a.length \leq 2^{51}$  }
4   ensures {
5     let exact = log(sum_real (fun i -> exp(a[i])) 0 n) in
6     let err = exp_error + 2ε a.length (1 + exp_error) in
7     |result - exact| ≤
8       log_error * |exact| + log (1 - err) * (1 + log_error)
9   }
10  =
11  let ref s = 0 in
12  for i = 0 to a.length - 1 do
13    invariant {
14      s = sum_double (fun i -> exp_approx (a[i])) 0 i }
15    s <- s + exp_approx a[i]
16  done;
17  log_approx s
```

ERROR BOUND ON MLSE

ERROR BOUND ON MLSE

ON PAPER

ERROR BOUND ON MLSE

Reminder:

$$\text{MLSE}(a, n, x) = \log(n) + \frac{x^2}{\log(4)} + \sum_{i=1}^n -\log \left(\sum_{j=1}^n \exp \left(\frac{-(a_i + x - a_j)^2}{2} \right) \right)$$

We define :

$$\text{MLSE}_{||}(a, n, x) = \log(n) + \frac{x^2}{\log(4)} + \sum_{i=1}^n \left| -\log \left(\sum_{j=1}^n \exp \left(\frac{-(a_i + x - a_j)^2}{2} \right) \right) \right|$$

Result 3

$$\begin{aligned} \left| \widehat{\text{MLSE}}(a, n, x) - \text{MLSE}(a, n, x) \right| &\leq \varepsilon |\text{MLSE}(a, n, x)| \\ &\quad + \text{MLSE}_{||}(a, n, x) \times (2E_{\log} + 2^{-48}n) \\ &\quad - 6n \log(1 - \max(2^{-45}, E_{\exp})) \\ &\quad + 2E_{\log} |\log(n)| + 2^{-1072} \end{aligned}$$

PROVING C CODE

β is still a prototype under development:

- No ghost functions
 - ▶ Consequence: manual applications of lemmas in Why3 IDE
- Encoding layer between C and WhyML
 - ▶ Consequence: additional transformations need to be performed on Why3 IDE before applying lemmas

We lose in automation

REUSING WHY3 THEORIES

β^3 attributes to map an existing Why3 symbol to an ACSL symbol

```
/*@ axiomatic SumFloat __attribute__((j3_theory ("sum_double.Sum"))) {  
    logic double sum_float(double f[], integer a,  
        integer b) __attribute__((j3_symbol("Sum.sum")));  
} /*
```

We can do part of the proof in Why3

Not working currently

CONCLUSIONS

RESULTS OBTAINED

- Verification of floating-point summation in WhyML with the bound:

$$\left| \bigoplus_{i=1}^n a_i - \sum_{i=1}^n a_i \right| \leq 2\epsilon n \sum_{i=1}^n |a_i|$$

- Verification of LSE function in WhyML and in C with the bound :

$$\begin{aligned} & \left| \widehat{\text{LSE}}(a, n) - \text{LSE}(a, n) \right| \\ & \leq E_{\log} |\text{LSE}(a, n)| - \log(1 - (E_{\text{exp}} + 2\epsilon n(1 + E_{\text{exp}}))) \times (1 + E_{\log}) \end{aligned}$$

- Verification of MLSE function in WhyML with the bound :

$$\begin{aligned} & \left| \widehat{\text{MLSE}}(a, n, x) - \text{MLSE}(a, n, x) \right| \leq \epsilon |\text{MLSE}(a, n, x)| \\ & \quad + \text{MLSE}_{||}(a, n, x) \times (2E_{\log} + 2^{-48}n) \\ & \quad - 6n \log(1 - \max(2^{-45}, E_{\text{exp}})) \\ & \quad + 2E_{\log} |\log(n)| + 2^{-1072} \end{aligned}$$

- Technical details :
 - ▶ Proof of C versions with J^3 by supporting attributes
 - ▶ Support for anonymous functions in TIS-kernel and J^3
- Better automation for rounding errors combination ?
- Better handling of higher order in MLSE ?