

**The Metalibm project :
assisting the production
of high-performance, proven
elementary function code**

Florent de Dinechin
CITI/Socrates project



Outline

Introduction : performance versus accuracy

Elementary function schizophrenia

The art of implementing elementary functions

Correctly rounded functions computing just right

Open-source tools for FP coders

Three metalibm prototypes

Conclusion

Introduction : performance versus accuracy

Introduction : performance versus accuracy

Elementary function schizophrenia

The art of implementing elementary functions

Correctly rounded functions computing just right

Open-source tools for FP coders

Three metalibm prototypes

Conclusion

Common wisdom

The more accurate you compute, the more expensive it gets

In practice

- We (hopefully) notice it when our computation is **not accurate enough**.
- But do we notice it when it is **too accurate** for our needs?

Reconciling performance and accuracy ?

Or, regaining performance by computing just right ?

Double precision spoils us

The standard binary64 format (formerly known as double-precision) provides roughly **16** decimal digits.

Why should anybody need such accuracy?

Count the digits in the following

- Definition of the second : *the duration of **9,192,631,770** periods of the radiation corresponding to the transition between the two hyperfine levels of the ground state of the cesium 133 atom.*
- Definition of the metre : *the distance travelled by light in vacuum in $1/\mathbf{299,792,458}$ of a second.*
- Most accurate measurement ever (another atomic frequency) to 14 decimal places
- Most accurate measurement of the Planck constant to date : to 7 decimal places
- The gravitation constant G is known to 3 decimal places only

Parenthesis : then why binary64 ?

- This PC computes 10^9 operations per second (1 gigaflops)

An allegory due to Kulisch

- print the numbers in 100 lines of 5 columns double-sided :
1000 numbers/sheet
- 1000 sheets \approx a heap of 10 cm
- 10^9 flops \approx heap height speed of 100m/s, or 360km/h
- A teraflops (10^{12} op/s) prints to the moon in one second
- Current top 500 computers reach the petaflop (10^{15} op/s)
- each operation may involve a relative error of 10^{-16} ,
and they accumulate.

Doesn't this sound wrong ?

We would use these 16 digits just to accumulate garbage in them ?

Back to the point

... which was :

Mastering accuracy for performance

When implementing a “computing core”

- A goal : *never compute more accurately than needed*
- Two sub-goals
 - Know what accuracy you need
 - Know how accurate you compute

“Computing cores” considered so far : **elementary functions**, sums of products, linear algebra, Euclidean lattices algorithms.

By the way

“computing just right” implies “computing right” ...

A general technique for computing just right

I've seen it for orientation predicates, area of a triangle, elementary functions...

Fast in average, always accurate

1. use a quick and dirty routine
2. runtime-test if it was accurate enough
3. launch an expensive, accurate routine only when needed

If done well, **average time is close to that of the quick routine**

Only works if you know how to implement step 2

... requires to understand/master/engineer the accuracy of your code.

The rest of this talk is about this... elementary functions being used as an easy example.

Elementary function schizophrenia

Introduction : performance versus accuracy

Elementary function schizophrenia

The art of implementing elementary functions

Correctly rounded functions computing just right

Open-source tools for FP coders

Three metalibm prototypes

Conclusion

A producer/consumer mismatch

On the producer side

- Specification : a **mathematical function** function (e.g. \log), and a **floating-point format** (binary32 or binary64)
- Obviously, one should return the **correct rounding** of $\log(x)$ (IEEE-754-2008)
- Anything less means wasting bits, probably cheating for speed, and introducing system-specific idiosyncrasies.

On the consumer side

Code should

1. use all the available space but no more,
2. be as fast as possible, and
3. be “accurate enough”

Correctly rounded probably **vast overkill** for “accurate enough”.

Not only a quantitative matter



The overkill is not only about the accuracy (the number of correct bits).
Just two examples :

- In binary floating-point, $\sin(x)$ is very difficult to evaluate, just because 2 and π are irrational.
 - IEEE754-2008 introduced the new function $\sin\text{Pi}(x) = \sin(\pi x)$.
 - For most applications, it is just as good, only faster.
- Fixed-point versus floating-point (see next slide)

Fixed-point versus floating-point

What is wrong with the following computation in floating-point ?

$$s = \sum_{i=1}^n x_i \quad y = e^s$$

A fix-point number is a binary integer, scaled by some constant 2^m

- cheaper to compute on (addition defines the PC's cycle time)
- simpler to analyze (addition is exact until it overflows)
- most sensors provide fixed-point data
- many computations are inherently fixed-point (i.e. do not take very large values nor need arbitrary resolution around 0).

Exponential and logarithm functions

... map a fixed-point range to a floating-point-range

I'm schizophrenic, but I'm under treatment

- On one side, I'm advocating perfectly accurate (correctly rounded) functions.
- On the other side, I'm advocating "computing just right" using less accurate functions

This MetaLibm in the title

- should provide all these variants
- ... if you provide it with explicit accuracy specifications
- (yes I'm trying to twist your arm into computing just right)

The art of implementing elementary functions

Introduction : performance versus accuracy

Elementary function schizophrenia

The art of implementing elementary functions

Correctly rounded functions computing just right

Open-source tools for FP coders

Three metalibm prototypes

Conclusion

How does your PC compute elementary functions?

Rule of the game : use the hardware, i.e. $+$, $-$, \times
(and maybe $/$ and $\sqrt{\quad}$ but they are expensive).

- **Polynomial approximation** works on a small interval
- **Argument reduction** : using mathematical identities, transform large arguments in small ones

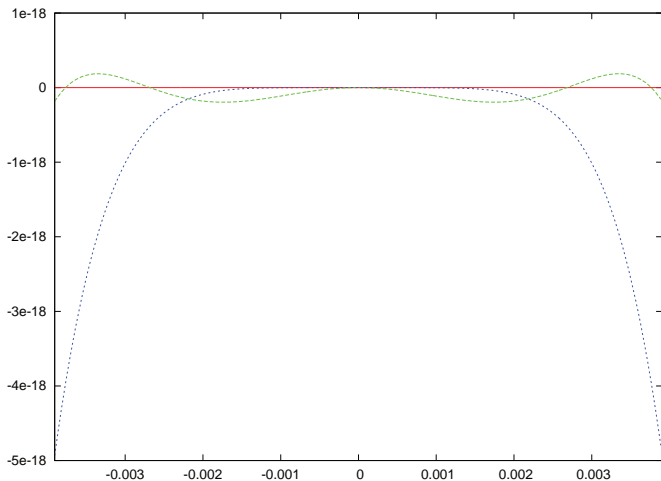
Simplistic example : an exponential

- identity : $e^{a+b} = e^a \times e^b$
- split $x = a + b$
 - a : k leading bits of x
 - b : lower bits of x $b \ll 1$
- tabulate all the e^a (2^k entries)
- use a Taylor polynomial for e^b

Know how accurate you compute

- Approximation errors
 - example : approximate a function f with a polynomial p :
 $\|p - f\|_\infty ?$
(see next slide)
 - in general : approximate an object by another one
- Rounding errors
 - for data, often called quantization errors ;
 - for operations, each individual error well specified by IEEE-754
 - but their accumulation difficult to manage
- In physics : time discretization errors, etc

Approximation of a function by a polynomial



$\|p - f\|_\infty$ for degree-5 Taylor and Remez approx. to exp on $[-2^{-8}, 2^{-8}]$

What is an error ? What is accuracy ?

The most important sentence of this talk

An error is a difference (absolute or relative) between two values, one being a reference for the other.

Examples :

- $\|p - f\|_\infty$ above is an upper bound of the approximation error
- error of the polynomial is with reference to the function (easy)
- error of the FP addition is with reference of the real sum (easy)
- error of one FP addition within the polynomial evaluation ?
(difficult because we have no direct reference in the function)
- yesterday : accuracy of the summation algorithms ?

Never say “the error of this term is ...” :

it doesn't mean anything without the reference.

*If you are not able to define the reference value,
you will not be able to know how accurate you compute*

Parenthesis : reproductibility and predictability

As soon as we are able to define the reference value,

who cares about exact reproductibility?

- What matters is to be able to reproduce enough significant digits.
- The compiler will not help you there :
No compiler has no access to the reference ! It is not in the code.

Important information NOT in the code

The context

$$y \in [-\pi/256, \pi/256]$$

What it is supposed to compute

a sine accurate to 2^{-60}

My programmer expertise

$y*(1+ts)$ is a bit less accurate than $y + y*ts$ in floating-point
... because $|t| < 2^{-14}$ because $|y| < 2^{-7}$

$$\begin{array}{r} \boxed{1} \\ + \quad \boxed{t} \\ = \quad \boxed{1+t} \end{array}$$

$$\begin{array}{r} \boxed{y} \\ + \quad \boxed{y*t} \\ = \quad \boxed{y+y*t} \end{array}$$

On the positive side : combining errors is easy

Since an error is a difference :

$$\begin{aligned} F(x) - f(x) &= F(x) - p(x) + p(x) - f(x) \\ &\quad (\text{rounding error} + \text{polynomial approximation error}) \\ |F(x) - f(x)| &\leq |F(x) - p(x)| + |p(x) - f(x)| \end{aligned}$$

... then recurse on $F(x) - p(x)$

Difficulties

- define “intermediate reference values”
- do not forget anything
- relative errors :

$$\frac{a - c}{c} = \frac{a - b}{b} + \frac{b - c}{c} + \frac{a - b}{b} \times \frac{b - c}{c}$$

Later in this talk : **Gappa**, a tool that helps you with all this

Correctly rounded functions computing just right

Introduction : performance versus accuracy

Elementary function schizophrenia

The art of implementing elementary functions

Correctly rounded functions computing just right

Open-source tools for FP coders

Three metalibm prototypes

Conclusion

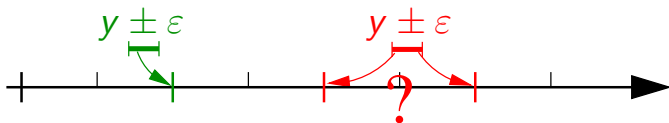
Know what accuracy you need ?

Correctly rounded elementary functions

- IEEE-754 floating-point single or double-precision
- **Elementary functions** : sin, cos, exp, log, implemented in the “standard mathematical library” (`libm`)
- **Correctly rounded** : As perfect as can be, considering the finite nature of floating-point arithmetic
 - same standard of quality as $+$, \times , $/$, $\sqrt{\quad}$
- Now recommended by the IEEE754-2008 standard, but long considered **too expensive**
because of the **Table Maker's Dilemma**

The Table Maker's Dilemma

- Finite-precision algorithm for evaluating $f(x)$
- Approximation + rounding errors \rightarrow overall error bound $\bar{\epsilon}$.
- What we compute : y such that $f(x) \in [y - \bar{\epsilon}, y + \bar{\epsilon}]$



Dilemma if this interval contains a midpoint between two FP numbers

The first digital signature algorithm

LOGARITHMICA.

Tabula inventum Logarithmorum inferiorem.

1	0,00	100001	0,00000,41429,2
2	0,30102,99915,6	100002	0,00000,86858,0
3	0,47712,12547,2	100003	0,00001,32286,4
4	0,60205,99903,3	100004	0,00001,77714,8
5	0,69897,00002,9	100005	0,00002,23142,3
6	0,77815,12553,8	100006	0,00002,68570,7
7	0,84509,80400,1	100007	0,00003,13998,1
8	0,90308,99869,9	100008	0,00003,59426,5
9	0,95424,25044,4	100009	0,00004,04854,9
10	1,00000,00000,0	100010	0,00004,50283,3
11	0,04139,26877,6	100011	0,00004,95711,7
12	0,07918,12460,5	100012	0,00005,41140,1
13	0,11918,21322,1	100013	0,00005,86568,5
14	0,16122,30165,8	100014	0,00006,32000,0
15	0,17609,12390,6	100015	0,00006,77432,4
16	0,20411,99826,6	100016	0,00007,22864,8
17	0,23041,89213,8	100017	0,00007,68297,2
18	0,25122,25071,0	100018	0,00008,13730,6
19	0,27877,36092,5	100019	0,00008,59163,0
20	0,30213,17272,8	100020	0,00009,04595,4
21	0,32186,01717,6	100021	0,00009,49998,8
22	0,33812,72247,1	100022	0,00010,00000,0
23	0,35103,33393,0	100023	0,00010,45003,4
24	0,36118,99997,7	100024	0,00010,90006,8
25	0,36909,80772,5	100025	0,00011,35010,2
26	0,37509,37772,9	100026	0,00011,80013,6
27	0,37944,77772,9	100027	0,00012,25017,0
28	0,38234,64272,6	100028	0,00012,70020,4
29	0,38409,40772,8	100029	0,00013,15023,8
30	0,38485,77215,1	100030	0,00013,60027,2
31	0,38471,37128,1	100031	0,00014,05030,6
32	0,38365,66172,6	100032	0,00014,50034,0
33	0,38179,79807,2	100033	0,00014,95037,4
34	0,37912,94705,5	100034	0,00015,40040,8
35	0,37565,07321,1	100035	0,00015,85044,2
36	0,37139,11662,4	100036	0,00016,30047,6
37	0,36634,12272,8	100037	0,00016,75051,0
38	0,36058,68702,1	100038	0,00017,20054,4
39	0,35413,23282,1	100039	0,00017,65057,8
40	0,34708,30165,8	100040	0,00018,10061,2
41	0,33943,79807,2	100041	0,00018,55064,6
42	0,33118,94705,5	100042	0,00019,00068,0
43	0,32243,07321,1	100043	0,00019,45071,4
44	0,31318,11662,4	100044	0,00019,90074,8
45	0,30343,12272,8	100045	0,00020,35078,2
46	0,29318,68702,1	100046	0,00020,80081,6
47	0,28243,23282,1	100047	0,00021,25085,0
48	0,27118,30165,8	100048	0,00021,70088,4
49	0,25943,79807,2	100049	0,00022,15091,8
50	0,24718,94705,5	100050	0,00022,60095,2
51	0,23443,07321,1	100051	0,00023,05098,6
52	0,22118,11662,4	100052	0,00023,50102,0
53	0,20743,12272,8	100053	0,00023,95105,4
54	0,19318,68702,1	100054	0,00024,40108,8
55	0,17843,23282,1	100055	0,00024,85112,2
56	0,16318,30165,8	100056	0,00025,30115,6
57	0,14743,79807,2	100057	0,00025,75119,0
58	0,13118,94705,5	100058	0,00026,20122,4
59	0,11443,07321,1	100059	0,00026,65125,8
60	0,09718,11662,4	100060	0,00027,10129,2
61	0,07943,12272,8	100061	0,00027,55132,6
62	0,06118,68702,1	100062	0,00028,00136,0
63	0,04243,23282,1	100063	0,00028,45139,4
64	0,02318,30165,8	100064	0,00028,90142,8
65	0,00343,79807,2	100065	0,00029,35146,2
66	0,00000,00000,0	100066	0,00029,80149,6
67	0,00000,00000,0	100067	0,00030,25153,0
68	0,00000,00000,0	100068	0,00030,70156,4
69	0,00000,00000,0	100069	0,00031,15159,8
70	0,00000,00000,0	100070	0,00031,60163,2
71	0,00000,00000,0	100071	0,00032,05166,6
72	0,00000,00000,0	100072	0,00032,50170,0
73	0,00000,00000,0	100073	0,00032,95173,4
74	0,00000,00000,0	100074	0,00033,40176,8
75	0,00000,00000,0	100075	0,00033,85180,2
76	0,00000,00000,0	100076	0,00034,30183,6
77	0,00000,00000,0	100077	0,00034,75187,0
78	0,00000,00000,0	100078	0,00035,20190,4
79	0,00000,00000,0	100079	0,00035,65193,8
80	0,00000,00000,0	100080	0,00036,10197,2
81	0,00000,00000,0	100081	0,00036,55200,6
82	0,00000,00000,0	100082	0,00037,00204,0
83	0,00000,00000,0	100083	0,00037,45207,4
84	0,00000,00000,0	100084	0,00037,90210,8
85	0,00000,00000,0	100085	0,00038,35214,2
86	0,00000,00000,0	100086	0,00038,80217,6
87	0,00000,00000,0	100087	0,00039,25221,0
88	0,00000,00000,0	100088	0,00039,70224,4
89	0,00000,00000,0	100089	0,00040,15227,8
90	0,00000,00000,0	100090	0,00040,60231,2
91	0,00000,00000,0	100091	0,00041,05234,6
92	0,00000,00000,0	100092	0,00041,50238,0
93	0,00000,00000,0	100093	0,00041,95241,4
94	0,00000,00000,0	100094	0,00042,40244,8
95	0,00000,00000,0	100095	0,00042,85248,2
96	0,00000,00000,0	100096	0,00043,30251,6
97	0,00000,00000,0	100097	0,00043,75255,0
98	0,00000,00000,0	100098	0,00044,20258,4
99	0,00000,00000,0	100099	0,00044,65261,8

15

- I want 12 significant digits
- I have an approximation scheme that provides 14 digits

• or,

$$y = \log(x) \pm 10^{-14}$$

• “Usually” that’s enough to round

$$y = x, \text{xxxxxxxxxxxx}17 \pm 10^{-14}$$

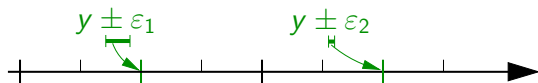
$$y = x, \text{xxxxxxxxxxxx}83 \pm 10^{-14}$$

• **Dilemma** when

$$y = x, \text{xxxxxxxxxxxx}50 \pm 10^{-14}$$

The first table-makers rounded these cases randomly, and recorded them to confound copiers.

Solving the table maker's dilemma



Ziv's onion peeling algorithm

1. Initialisation : $\varepsilon = \varepsilon_1$
2. Compute y such that $f(x) = y \pm \varepsilon$
3. Does $y \pm \varepsilon$ contain the middle point between two FP numbers?
 - If no, return $\text{RN}(y)$
 - If yes, dilemma ! Reduce ε , and go back to 2

It is a *while* loop... we have to show it terminates, a topic in itself.

Accuracy versus performance

CRLibm : 2-step approximation process

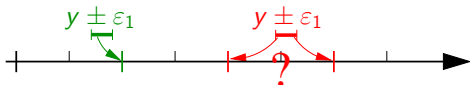
- first step **fast** but accurate to $\bar{\epsilon}_1$ sometimes not accurate enough
- (rarely) second step slower but **always accurate enough**

$$T_{\text{avg}} = T_1 + p_2 T_2$$

For each step, we need a **tight** bound on the error of the code :

$$\left| \frac{F(x) - f(x)}{f(x)} \right| \leq \bar{\epsilon}$$

- Overestimating $\bar{\epsilon}_2$ degrades T_2 ! (common wisdom)
- Overestimating $\bar{\epsilon}_1$ degrades p_2 !



First function development in Arénaire

First correctly rounded elementary function in CRLibm

- exp by David Defour
- worst-case time $T_2 \approx 10,000$ cycles
- complex, hand-written proof
- duration : a Ph.D. thesis (2002)

Conclusion was :

- performance and memory consumption of CR elem function is OK
- problem now is : performance and coffee consumption of the programmer (and that is because of the need for tight error bounds)

Latest CR function developments in our group

C. Lauter at the end of his PhD

- development time for sinpi, cospi, tanpi : 2 days
- worst-case time $T_2 \approx 1,000$ cycles

(but as a result of three more PhDs)

With MetaLibm prototype

- development time for 16 variants of sinpi, cospi, sincospi : 1 day

Open-source tools for FP coders

Introduction : performance versus accuracy

Elementary function schizophrenia

The art of implementing elementary functions

Correctly rounded functions computing just right

Open-source tools for FP coders

Three metalibm prototypes

Conclusion

The GMP family

- GMP (GNU Multiple Precision) and its beautiful C++ wrapper
 - integer arithmetic
 - best asymptotic algorithms + lower layers in hand-crafted assembly code
- MPFR : Multiple Precision Floating-point correctly Rounded
 - a floating-point layer on top of GMP
 - IEEE 754-like specification
 - Now a dependency of GCC, so you probably have it installed
 - Unfortunately, no hope of a beautiful C++ wrapper that would always take the proper decision as of **where to round**
- MPFI : interval arithmetic on top of MPFR

The Patriot bug

In 1991, a Patriot missile failed to intercept a Scud, and 28 people were killed.

- The code worked with time increments of 0.1 s.
- But 0.1 is not representable in binary.
- In the 24-bit format used, the number stored was 0.099999904632568359375
- The error was 0.0000000953.
- After 100 hours = 360,000 seconds, time is wrong by 0.34s.
- In 0.34s, a Scud moves 500m

In single, we don't have that many bits to accumulate garbage in them !

Test : which of the following increments should you use ?

10 5 3 1 0.5 0.25 0.2 0.125 0.1

Sollya (2)

Killer feature 2

multiple-precision, last-bit accurate evaluation of arbitrary expressions

```
1 fdedinec@krupnik: sollya
2 > e=exp(x) - (1+x+x^2/2+x^3/6);
3 > e(0.125);
4 Warning: rounding has happened. The value displayed is a
   faithful rounding of the true result.
5 1.04322334929834956738944784605392321697984118482926e-5
6 >
```

All these digits are meaningful! This is better than Maple.

Sollya (3)

Killer feature 3

guaranteed infinite norm $\|f(x)\|_\infty$ even in degenerate cases

- $\|f(x) - P(x)\|_\infty$ is a degenerate case...

Sollya (4)

Killer feature 4

Machine-efficient polynomial approximation

- Remez' minimax algorithm finds the best polynomial approximation **over the reals**
- But we need polynomials with **machine** coefficients
 - float, double, fixed-point, ...
- Rounding Remez coefficients does **not** provide the best polynomial among polynomial with machine coefficients.
- Sollya does (almost).
 - this saves a few bits of accuracy
 - especially relevant for small precisions (FPGAs)
 - that's how we get our polynomials

Nice number theory behind. And needs all the previous.

6 guaranteed log polynomials on one slide

A sollya script that computes approximations to the log of various qualities

```
f= log(1+y);
I=[-0.25;.5];
filename="/tmp/polynomials";
print("") > filename;
for deg from 2 to 8 do begin
  p = fpminimax(f, deg,[|0,23...|],I, floating, absolute);
  display=decimal;
  acc=floor(-log2(sup(supnorm(p, f, I, absolute, 2^(-40)))));
  print( "    // degree = ", deg,
        " => absolute accuracy is ", acc, "bits" ) >> filename;
  print("#if ( DEGREE ==", deg, ")") >> filename;
  display=hexadecimal;
  print("    float p = ", horner(p) , ";") >> filename;
  print("#endif") >> filename;
end;
```

Gamma : motivation

crlibm.pdf 5 years ago : 124 pages of this

```
1  yh2 = yh*yh;
2  ts = yh2 * (s3.d + yh2*(s5.d + yh2*s7.d));
3  Add12(*psh,*psl, yh, yl+ts*yh);
```

Upon entering DoSinZero, we have in $y_h + y_l$ an approximation to the ideal reduced value $\hat{y} = x - k \frac{\pi}{256}$ with a relative accuracy $\varepsilon_{\text{argred}}$:

$$y_h + y_l = \left(x - k \frac{\pi}{256}\right)(1 + \varepsilon_{\text{argred}}) = \hat{y}(1 + \varepsilon_{\text{argred}}) \quad (1)$$

with, depending on the quadrant, $\sin(\hat{y}) = \pm \sin(x)$ or $\sin(\hat{y}) = \pm \cos(x)$ and similarly for $\cos(\hat{y})$. This just means that \hat{y} is the ideal, errorless reduced value.

In the following we will assume we are in the case $\sin(\hat{y}) = \sin(x)$, (the proof is identical in the other cases), therefore the relative error that we need to compute is

$$\varepsilon_{\text{sinkzero}} = \frac{(*\text{psh} + *\text{psl})}{\sin(x)} - 1 = \frac{(*\text{psh} + *\text{psl})}{\sin(\hat{y})} - 1 \quad (2)$$

One may remark that we almost have the same code as we have for computing the sine of a small argument (without range reduction). The difference is that we have as input a double-double $y_h + y_l$, which is itself an inexact term.

At Line 4, the error of neglecting y_l and the rounding error in the multiplication each amount to half an ulp :

$$y_{h2} = y_h^2(1 + \varepsilon_{-53}), \text{ with } y_h = (y_h + y_l)(1 + \varepsilon_{-53}) = \hat{y}(1 + \varepsilon_{\text{argred}})(1 + \varepsilon_{-53})$$

Therefore

$$yh2 = \hat{y}^2(1 + \varepsilon_{yh2}) \quad (3)$$

with

$$\bar{\varepsilon}_{yh2} = (1 + \bar{\varepsilon}_{argred})^2(1 + \bar{\varepsilon}_{-53})^3 - 1 \quad (4)$$

Line 5 is a standard Horner evaluation. Its approximation error is defined by :

$$P_{ts}(\hat{y}) = \frac{\sin(\hat{y}) - \hat{y}}{\hat{y}}(1 + \varepsilon_{approx\,ts})$$

This error is computed in **Maple** as previously, only the interval changes :

$$\bar{\varepsilon}_{approx\,ts} = \left\| \frac{xP_{ts}(x)}{\sin(x) - x} - 1 \right\|_{\infty}$$

We also compute $\bar{\varepsilon}_{hornert\,ts}$, the bound on the relative error due to rounding in the Horner evaluation thanks to the `compute_horner_rounding_error` procedure. This time, this procedure takes into account the relative error carried by `yh2`, which is $\bar{\varepsilon}_{yh2}$ computed above. We thus get the total relative error on `ts` :

$$ts = P_{ts}(\hat{y})(1 + \varepsilon_{hornert\,ts}) = \frac{\sin(\hat{y}) - \hat{y}}{\hat{y}}(1 + \varepsilon_{approx\,ts})(1 + \varepsilon_{hornert\,ts}) \quad (5)$$

The final Add12 is exact. Therefore the overall relative error is :

$$\begin{aligned}
 \varepsilon_{\text{sinkzero}} &= \frac{((y_h \otimes t_s) \oplus y_l) + y_h}{\sin(\hat{y})} - 1 \\
 &= \frac{(y_h \otimes t_s + y_l)(1 + \varepsilon_{-53}) + y_h}{\sin(\hat{y})} - 1 \\
 &= \frac{y_h \otimes t_s + y_l + y_h + (y_h \otimes t_s + y_l) \cdot \varepsilon_{-53}}{\sin(\hat{y})} - 1
 \end{aligned}$$

Let us define for now

$$\delta_{\text{addsin}} = (y_h \otimes t_s + y_l) \cdot \varepsilon_{-53} \quad (6)$$

Then we have

$$\varepsilon_{\text{sinkzero}} = \frac{(y_h + y_l)t_s(1 + \varepsilon_{-53})^2 + y_l + y_h + \delta_{\text{addsin}}}{\sin(\hat{y})} - 1$$

Using (1) and (5) we get :

$$\varepsilon_{\text{sinkzero}} = \frac{\hat{y}(1 + \varepsilon_{\text{argred}}) \times \frac{\sin(\hat{y}) - \hat{y}}{\hat{y}} (1 + \varepsilon_{\text{approx}t_s})(1 + \varepsilon_{\text{hornert}t_s})(1 + \varepsilon_{-53})^2 + y_l + y_h + \delta_{\text{addsin}}}{\sin(\hat{y})} - 1$$

To lighten notations, let us define

$$\varepsilon_{\text{sin1}} = (1 + \varepsilon_{\text{approx}t_s})(1 + \varepsilon_{\text{hornert}t_s})(1 + \varepsilon_{-53})^2 - 1 \quad (7)$$

We get

$$\begin{aligned}\varepsilon_{\text{sinkzero}} &= \frac{(\sin(\hat{y}) - \hat{y})(1 + \varepsilon_{\text{sin1}}) + \hat{y}(1 + \varepsilon_{\text{argred}}) + \delta_{\text{addsin}} - \sin(\hat{y})}{\sin(\hat{y})} \\ &= \frac{(\sin(\hat{y}) - \hat{y}) \cdot \varepsilon_{\text{sin1}} + \hat{y} \cdot \varepsilon_{\text{argred}} + \delta_{\text{addsin}}}{\sin(\hat{y})}\end{aligned}$$

Using the following bound :

$$|\delta_{\text{addsin}}| = |(\mathbf{yh} \otimes \mathbf{ts} + \mathbf{y1}) \cdot \varepsilon_{-53}| < 2^{-53} \times |y|^3/3 \quad (8)$$

we may compute the value of $\bar{\varepsilon}_{\text{sinkzero}}$ as an infinite norm under **Maple**. We get an error smaller than 2^{-67} .

4 pages for 3 lines of code...

Two years of experience showed that nobody (including myself) should trust such a proof (and that nobody reads it anyway).

We wish we had an **automatic** tool that

- takes a set of C files,
- parses them,
- and outputs “The overall error of the computation is ...”.

It's hopeless, because the code doesn't include

- what it is supposed to compute
- the knowledge used to build it

Gappa

Written by Guillaume Melquiond, Gappa is a tool that

- takes an input that closely matches your C file,
- forces you to express **what this code is supposed to compute**
- ... and some numerical property to prove (expressed in terms of intervals)
- and eventually outputs a proof of this property suitable for checking by Coq or HOL Light

`gappa.gforge.inria.fr/`

Code generation for polynomial evaluation

- explores different parallelizations of a polynomial on a VLIW processor
- generates **code** and **Gappa proof of the evaluation error**

Used to generate the code for the division and square root of FLIP, a Floating-Point Library for Integer Processors (collaboration with ST Microelectronics)
Should help for vector code as well.

`cgpe.gforge.inria.fr/`

Three metalibm prototypes

Introduction : performance versus accuracy

Elementary function schizophrenia

The art of implementing elementary functions

Correctly rounded functions computing just right

Open-source tools for FP coders

Three metalibm prototypes

Conclusion

Challenges of elementary function coding

- Multiple targets
 - numerous architectures (Intel's x86/Xeon-Phi, Kalray's K1, ARM)
 - numerous constraints (throughput, latency, precision, memory consumption)
 - numerous programming styles (e.g. vector versus scalar)
- More function out there than the ones listed in C11 !

libm paradigm poorly addresses this complexity.

Metalibm : generate function code on demand for a given context

The MetaLibm open-ended vision

- For correctly rounded elementary functions we needed to automate the development of code+proof
- Now that this is (almost) done, we may *open up the set of functions/precisions/performance constraints*

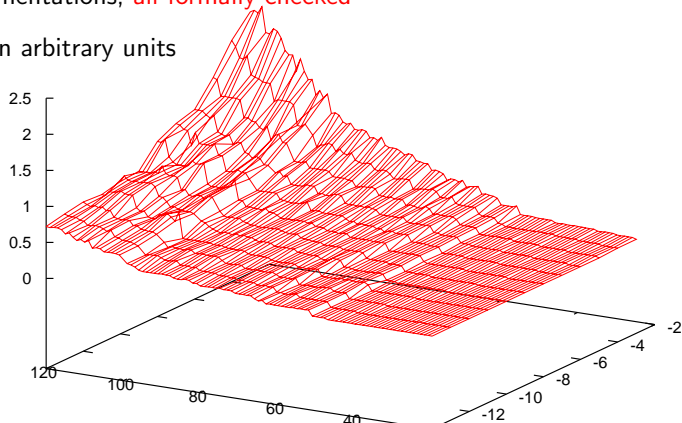
An ANR-funded project

- metalibm/OpenEnded
 - genericity in input
- metalibm/C11
 - focus on performance (match hand-coded libraries)
 - genericity in target processor
 - hand-code what we are unable (yet) to automate : range reductions, floating-point trickery, ...
- FPGAs, DSP filters for good measure

1/ Christoph Lauter's metalibm

- Example : $\log(1 + x)$
- Two parameters
 - k from 1 to 13, defines table size
 - target accuracy, between 20 and 120 bits
- 1203 implementations, **all formally checked**

z axis : timings in arbitrary units



Current experiments with this prototype

Connect it to the Dynamic Dictionary of Mathematical Functions

`http://ddmf.msr-inria.inria.fr/`

... to obtain code for a function defined by a **differential equation** +
limit conditions

2/ MetalibmC11 : an ad-hoc approach

Philosophy : take good working C code, wrap it in `printfs`, then generalize it.

- low abstraction, but succes guaranteed
- ... with perf identical to hand-written code

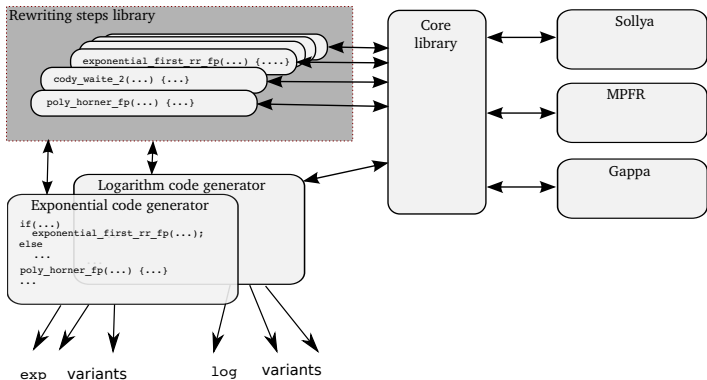
All scripted in Python.

All this is work in progress

- A `Processor` class and its subclasses
 - encapsulates processor-specific code generation and tricks
 - still tinkering a lot there
- A `Format` class and its subclasses
- A `Polynomial` class that manages both approximation and evaluation
- A `CFunction` class for libm functions
 - automatically generates test programs

Metaexp in one slide

- Already 8 useful implementations, each for arbitrary accuracy p (float/double, subnormals or not, Estrin or Horner)
- A case study for structuration as a metaskelton



- No Gappa generation yet

Metalog in one slide

- experiment with optimized for latency / optimized for throughput
 - using autovectorisation with gcc 4.7
 - works for single but not for double
(no %ymm1 in the generated assembly?!?)
Either AVX doesn't replicate all SSE2 functions, or GCC is not ready
- I'm not sure I understand how a degree-20 Horner polynomial is evaluated in 37 cycles
- Estrin evaluation would be useful here
 - but current implementation not modular enough
 - short-term TODO
- No Gappa generation yet

A glance at generated code

```
/* Exceptional case filtering, vectorizable */
minfty.ui = 0xff800000; /* minus infinity */
nan.ui = 0x7fc00000; /* nan */
ret_minfty = ((xx.ui & 0x7fffffff) == 0) ? minfty.f : 0.0f; /* x == +/-0 ?
ret_nan = (xx.ui > 0x80000000) ? nan.f : 0.0f; /* x<0 ?*/
x_is_inf_or_nan = ((xx.ui & 0x7fffffff) >= 0x7f800000) ? xx.f : 0.0f; /*
exn = ret_minfty + ret_nan + x_is_inf_or_nan; /* 0.0 if normal or subnormal
/* Now remains to add exn somewhere where it will propagate to the result
x_subnormal = (xx.ui < 0x00800000) && (xx.ui > 0);
subnormal_scale = x_subnormal ? 0x1.p48f : 1.0f; /* scale mantissa*/
e_x = x_subnormal ? -127-48 : -127; /* ... and initialize exponent*/
xx.f *= subnormal_scale;
/* Now decompose x into fraction and exponent */
e_x += ((xx.i) >> 23) & ((1<<8)-1); /* extract exponent*/
fraction.i = (xx.i & 0x007fffff); /* extract fraction bits*/
adjust = (fraction.i>>22); /* first non-implicit bit of the fraction, tell
fraction.i = fraction.i |0x3f800000; /* add the exponent of one */
fraction.i -= adjust << 23; /* if m>1.5, divide fraction by 2 (exact opera
e_x += adjust; /* and update exponent so we still have x = 2^e_x * fra
```



```
/* Now back to floating-point */
y = fraction.f - 1.0f; /* Sterbenz-exact; may cancel but we don't care */
y += exn; /* exn is either 0.0, or an inf or NaN that will propagate to the
/* Now y in [-0.25, 0.5], and we must evaluate log(1+y) */
/* Horner evaluation */
y2 = y*y;
p9 = c9;
p8 = c8 + y*p9;
p7 = c7 + y*p8;
p6 = c6 + y*p7;
p5 = c5 + y*p6;
p4 = c4 + y*p5;
p3 = c3 + y*p4;
p2 = c2 + y*p3;
p = y + y2*p2;
r = e_x*log_2 + p;
return r;
```

Metatrigpi in one slide

- $\sin(\pi x)$ and $\cos(\pi x)$ recommended by IEEE 754-2008
 - No costly range reduction
 - Correct rounding proven feasible
- `sincospif(float x, float *s, float *c)`
computes both in one function
- `sincospio2f(float x, float *s, float *c)`
computes $\sin(\frac{\pi}{2}x)$ and $\cos(\frac{\pi}{2}x)$ even faster

Developed in one day, with tests and all.

Weakness of previous attempts

- Lauter's too generic, will miss many function-specific tricks we need for performances.
- Dinechin's not abstract enough, misses "what's being computed". As a consequence, automatic proof generation mostly hopeless.

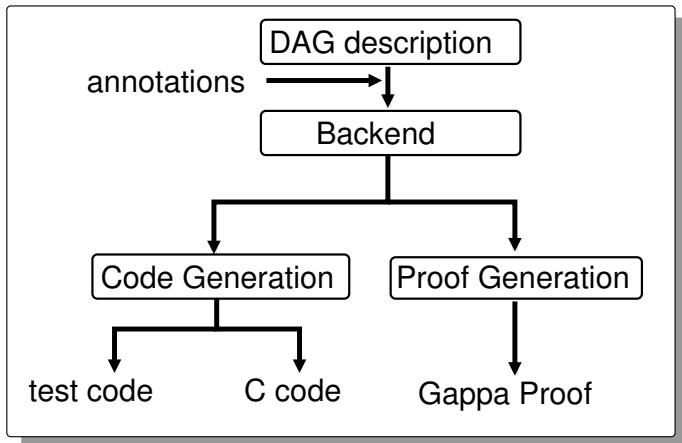
3/ Le troisième sera le bon

New Metalibm developed by N. Brunie

- DAG representation of the function, and of the implementation
- abstract target description

disconnect description, optimization and code+proof generation

MetLibm process



DAG description

- describe function implementation with operation nodes
 - **VARIABLE, Constant**
 - **TableLoad, TableLoad_HL**
 - **Addition, Multiplication, Modulo ...**
 - **Test, SpecOp, ExponentInsertion ...**
 - **ConditionBlock**
- DAG can be built two ways :
 - by composing constructors
 - by overloaded python expressions

example of DAG description

```
vx = VARIABLE("x", precision = fformat)

# reduced argument
red_x = NearestInt(vx / log(2), precision = int32, tag = "red_x")

# HIGH and LOW part log(2) generation
log2_hi = round(log(2), fformat.sollya_name - 10, RN)
log2_lo = round(log(2) - log2_hi, fformat.sollya_name, RN)
r = (vx - (red_x * log2_hi)) - red_x * log2_lo
r.set_attributes(tag = "r", exact = True)

red_int = Interval(-log(2)/2, log(2)/2)

poly = Polynomial.generate_fpminimax(exp(x), 5, red_int, [ML_Binary64]*6,
poly_scheme = PolySchemeGenerator.generate_horner(poly, r)

result = Return(ExponentInsertion(red_x) * poly_scheme)

backend_scheme = Backend(processor).backend_float(result, ML_Binary64)
source_code = CodeGenerator(processor).generate_expr(backend_scheme)
```

The annotation system

- to facilitate generated code reading
- to optimize DAG
- to enforce numerical constraints

Some examples of annotations :

- **tag, debug**
- **precision**
- **likely**
- **exact, interval**

Backend processing

The main operations performed by the backend are :

- instanciating every undetermined format
- introducing necessary conversions
- optimizations at the level of abstract operations

Other tasks include :

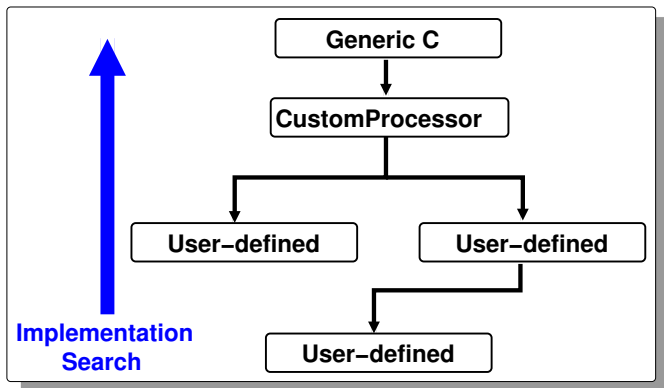
- dynamic support library expansion
- pre-vectorization processing

memoization is used for resource sharing

Code Generation

- generated from fully type-instantiated description
- constants, tables and core code generation
- several targets are available
- support for processor-specific code generation
- Gappa proof generation

Processor-specific code generation



Challenges of vectorization

- introduction of vector ISA : x86's MMX, SSE (1, 2, 3, 4, 4.1, 4.2), AVX (1 and 2), Xeon-Phi, ARM's NEON
- heavily reliant on branch uniformity
- compiler support can be unreliable
- generate branchless code
 - blending returns
 - extract most common path
 - generating callouts
- use intrinsics to force use of vector instructions

All this still preliminary

Challenges :

- Genericity versus ad-hoc efficiency (it is easy to write a generator that works for only one function)
- Draw the line between the compiler and Metalibm. Do not reimplement a full compiler!
- Open the code! (some at Intel, some at Kalray, all locked for the moment)

Conclusion

Introduction : performance versus accuracy

Elementary function schizophrenia

The art of implementing elementary functions

Correctly rounded functions computing just right

Open-source tools for FP coders

Three metalibm prototypes

Conclusion

Main messages

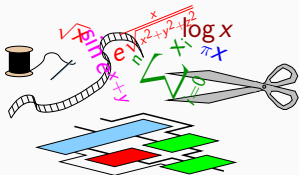
- If you're computing accurately enough, you're probably computing too accurately.

- Are you able to express what your code is supposed to compute?
If yes,
 - we can help you sort out the gory floating-point issues
 - we can provide functions computing just right for you

My other research project

Computing just right for FPGAs

... but I was given another advertising slot for this.



<http://flopoco.gforge.inria.fr/>

Thank you for your attention

Introduction : performance versus accuracy

Elementary function schizophrenia

The art of implementing elementary functions

Correctly rounded functions computing just right

Open-source tools for FP coders

Three metalibm prototypes

Conclusion