Jean-Michel Muller

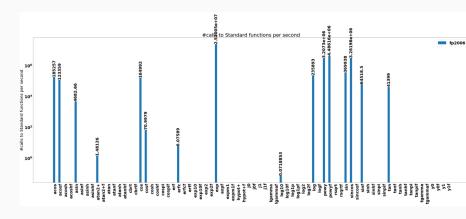
CNRS - Laboratoire LIP

http://perso.ens-lyon.fr/jean-michel.muller/

- real name: elementary transcendental functions. The ones than can be built from the complex exponential and logarithm;
- the ones listed by IEEE-754 are:

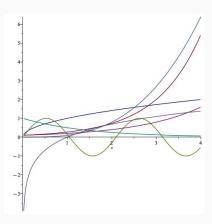
$$\begin{array}{c} {\rm e}^{x},{\rm e}^{x}-1,2^{x},2^{x}-1,10^{x},10^{x}-1,\\ {\rm log}(x),{\rm log}_{2}(x),{\rm log}_{10}(x),{\rm log}(1+x),{\rm log}_{2}(1+x),{\rm log}_{10}(1+x),\\ \sqrt{x^{2}+y^{2}},1/\sqrt{x},(1+x)^{n},x^{n},x^{1/n}(n\ {\rm is\ an\ integer}),x^{y},\\ {\rm sin}(\pi x),{\rm cos}(\pi x),{\rm tan}(\pi x),{\rm arcsin}(x)/\pi,{\rm arccos}(x)/\pi,{\rm arctan}(x)/\pi,{\rm arctan}(y/x)/\pi,\\ {\rm sin}(x),{\rm cos}(x),{\rm tan}(x),{\rm arcsin}(x),{\rm arccos}(x),{\rm arctan}(x),{\rm arctan}(y/x),\\ {\rm sinh}(x),{\rm cosh}(x),{\rm tanh}(x),{\rm arcsinh}(x),{\rm arccosh}(x),{\rm arctanh}(x). \end{array}$$

- the C Standard defines a rather similar list;
- a few functions (sin, cos, exp, log, ...) are very frequently called, and are
 used as building blocks for implementing the other functions → efficient
 and accurate implementation of these functions is necessary.



Number of calls per second of various functions in a CERN proton collision application.

- core set of "atomic functions": exp, log, sin, . . .
- highest possible quality: reproducibility, proven error bounds, etc.
- best possible quality: correct rounding.



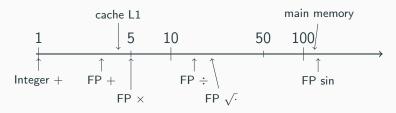
Work still needs to be done...

library	GNU libo	· IML	AMD	Newlib	OpenLibm	Musl	Apple	LLVM	MSVC	FreeBSD	ArmPL	CUDA	ROCm
version	2.40	2024.0.2	4.2	4.4.0	0.8.3	1.2.5	14.5	18.1.8	2022	14.1	24.04	12.2.1	5.7.0
acos	0.523	0.531	1.36	0.930	0.930	0.930	1.06		0.934	0.930	1.52	1.53	0.772
acosh	2.25	0.509	1.32	2.25	2.25	2.25	2.25		3.22	2.25	2.66	2.52	0.661
asin	0.516	0.531	1.06	0.981	0.981	0.981	0.709		1.05	0.981	(2.69)	1.99	0.710
asinh	1.92	0.507	1.65	1.92	1.92	1.92	1.58		2.05	1.92	2.04	2.57	0.661
atan	0.523	0.528	0.863	0.861	0.861	0.861	0.876		0.863	0.861	2.24	1.77	1.73
atanh	1.78	0.507	1.04	1.81	1.81	1.80	2.01		2.50	1.81	3.00	2.50	0.664
cbrt	3.67	0.523	1.53e22	0.670	0.668	0.668	0.729		1.86	0.668	1.79	0.501	0.501
cos	0.516	0.518	0.919	0.887	0.834	0.834	0.948	$_{\rm Inf}$	0.897	0.834		(1.52)	0.797
cosh	1.93	0.516	1.85	(2.67)	1.47	1.04	0.523		1.91	1.47	1.93	1.40	0.563
erf	1.43	0.773	1.00	1.02	1.02	1.02	6.41)	4.62	1.02	2.29	1.50	1.12
erfc	5.19	0.826		4.08	4.08	3.72	10.7		8.46	4.08	1.71	4.51	4.08
exp	0.511	0.530	1.01	0.949	0.949	0.511	0.521	0.500	1.50	0.949	0.511	0.928	0.929

Largest errors in ulps for double-precision calculation of some math functions. ulp (x) is the distance between two FP numbers in the neighborhood of x (so the largest values should be 0.5 – which is the case with +, -, \times , \div , and $\sqrt{\cdot}$).

(Extracted from Gladman, Innocente, Mather, and Zimmermann, Accuracy of Mathematical Functions..., Aug. 2024)

• hardware arithmetic of our processors: \pm , \times , ab + c (FMA), and \div (costs more than + and \times);

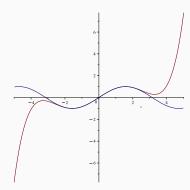


Typical current latencies in number of cycles (logarithmic scale). These figures vary from one processor to another, but the orders of magnitude remain similar.

- functions of one variable one can build from \pm and \times : polynomials;
- No choice: approximate the functions by piecewise polynomials.

Approximating the functions by polynomials?

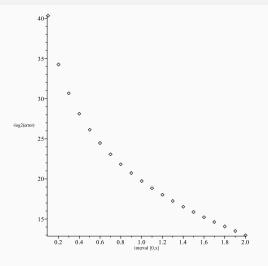
- polynomial approximation: valid in small interval only;
- need to reduce the initial argument to that interval;
- three steps:
 - range reduction;
 - polynomial evaluation;
 - · reconstruction.



The sine function and its best degree-5 approximation in $[-\pi/2, \pi/2]$.

(clearly valid only in $[-\pi/2, \pi/2]$)

The advantage of small intervals



 $-\log_2(error)$ of the best degree-5 approximation to e^t in [0,x] as a function of x.

How do you compute an exponential?

I want to compute e^x , I have a polynomial approximation to exp in $\left[-\frac{\ln(2)}{2},+\frac{\ln(2)}{2}\right]$.

- Range reduction: $x \to y = x k \log(2)$, with $k \in \mathbb{Z}$ and $y \in \left[-\frac{\ln(2)}{2}, +\frac{\ln(2)}{2}\right]$
- Polynomial evaluation: get approximation z to e^y using the polynomial;
- Reconstruction: obtain $e^x = z \cdot 2^k$.

How do you compute a cosine?

I want to compute $\cos(x)$. I have polynomial approximations to sin and \cos in $\left[-\frac{\pi}{4},\frac{\pi}{4}\right]$.

Argument reduction:
$$x \to y = x - k\frac{\pi}{2}$$
, with $k \in \mathbb{Z}$ and $y \in \left[-\frac{\pi}{4}, \frac{\pi}{4}\right]$

Polynomial evaluation: compute
$$\begin{cases} c = \cos(y) \text{ (if } k \text{ is even), or } \\ s = \sin(y) \text{ (if } k \text{ is odd).} \end{cases}$$

Reconstruction:
$$\operatorname{obtain} \operatorname{cos}(x) = \left\{ \begin{array}{ll} c & \text{if} & k \bmod 4 = 0 \\ -s & \text{if} & k \bmod 4 = 1 \\ -c & \text{if} & k \bmod 4 = 2 \\ s & \text{if} & k \bmod 4 = 3; \end{array} \right.$$

"Naive" range reduction: example with the constant $=\frac{\pi}{2}$

Naive reduction: we set $C = RN(\pi/2)$, and we successively compute

- $k = \lceil x/C \rceil$,
- y = RN(x kC) with an FMA, or RN(x RN(kC)) if no FMA..

Error analysis assuming FMA instruction is available

$$k = \lceil x/C \rfloor \implies \frac{x}{C} - \frac{1}{2} \le k \le \frac{x}{C} + \frac{1}{2}$$

$$\Rightarrow |x - kC| \le \frac{c}{2}.$$

$$|\operatorname{RN}(x - kC) - (x - k\frac{\pi}{2})| \le |\operatorname{RN}(x - kC) - (x - kC)| + |(x - kC) - (x - k\frac{\pi}{2})|$$

$$\le u \cdot |x - kC| + k \cdot |C - \frac{\pi}{2}|$$

$$\le u \frac{c}{2} + ku.$$

Hence:

- the absolute error is $\leq u(k+\frac{C}{2})$ (the bound grows linearly with x, since k is proportional to x);
- the relative error depends on how small $\left|x-k\frac{\pi}{2}\right|$ can be.
- \rightarrow continued fractions.

In prectice the naive reduction can be very inaccurate

In binary64 arithmetic (p = 53):

- if x = 355:
 - y is 7230134.89 ulp away from the exact reduced argument if we do not use an FMA, and
 - 4084406.89 ulp away with an FMA.
- If x = 37362253, even with an FMA, y is 440183437673129 ulp away from the exact value.

(remember: for us good accuracy means final error not much larger than 0.5 ulp !)

Continued fraction convergents to $\frac{\pi}{2}$:

```
1, 2, \frac{3}{2}, \frac{11}{7}, \frac{344}{219}, \frac{355}{226}, \frac{51819}{32989}, \frac{52174}{33215}, \frac{260515}{165849}, \frac{573204}{364913}, \frac{4846147}{3085153}, \frac{5419351}{3450066}, \frac{37362253}{23785549}, \dots
```

Cody and Waite reduction

Idea: approximate $\frac{\pi}{2}$ by two FP numbers C_1 and C_2 such that

- C_1 fits in p-m bits, where m is the max. number of bits of k for which we want an accurate result $\rightarrow kC_1$ will be a FP number
- $C_2 = RN(\frac{\pi}{2} C_1)$, so that $C_1 + C_2$ represents $\frac{\pi}{2}$ with significantly more than p bits of precision.
- More precisely: $\left| \frac{\pi}{2} C_1 C_2 \right| < 2^{-2p+m}$.

Then we compute (here with an FMA)

$$RN(RN(x-kC_1)-kC_2).$$

Cody and Waite reduction

We compute

$$RN(RN(x-kC_1)-kC_2).$$

- kC_1 is a FP number;
- as x and kC_1 are very near, their difference is a FP number (Sterbenz) $\rightarrow \text{RN}(x kC_1) = x kC_1$. Hence we obtain

$$\mathsf{RN}\,(x-k(C_1+C_2)).$$

same error analysis as the naive reduction:

$$\left| \mathsf{RN} \left(x - k (C_1 + C_2) \right) - \left(x - k \frac{\pi}{2} \right) \right| < u \frac{C}{2} + k u^2 2^m.$$

• With m = 26, if x = 355, y is 0.108 ulp away from the exact result (still bad for x = 37362253, but not as much).

Generalizations, variants

• More than 2 constants: one may have C_2 fit in p-m bits too, and approximate $\frac{\pi}{2}$ by a sum $C_1+C_2+C_3$. One then evaluates

$$RN(((x-kC_1)-kC_2)-kC_3);$$

- systematic study of the numbers $\gamma \approx \frac{\pi}{2}$ such that $x-k\gamma$ is a FP number;
- express the reduced argument in double-word arithmetic.

After range reduction: using a polynomial approximation to the function

Two issues:

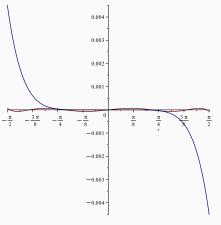
- choosing the right polynomial;
- evaluating it quickly and accurately.

Function f, approximated by $p(x) = a_0 + a_1x + a_2x^2 + a_3x^3 + \cdots + a_nx^n$; frequently evaluated as

$$a_0 + x \cdot (a_1 + x \cdot (a_2 + x \cdot (\cdots (a_{n-1} + x \cdot a_n))) \cdots)$$

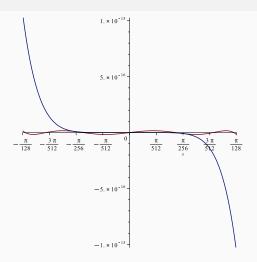
(called Horner's scheme (1819). Was already known by Newton, and probably by chinese mathematicians centuries before Horner)

Forget about Taylor series



- error of degree-5 Taylor series of sin vs error of best "minimax" degree-5 approximation in $\left[-\frac{\pi}{2},+\frac{\pi}{2}\right]$;
- Taylor series are local best approximations: they cannot compete on a whole interval.

Taking a much smaller domain does not change the difference



Error of degree-5 Taylor series of sin vs error of best "minimax" degree-5 approximation in $\left[-\frac{\pi}{128}, +\frac{\pi}{128}\right]$.

Minimax approximation of functions by polynomials

- $\mathcal{P}_n = \{ \text{polynomials of degree } \leq n \text{ with coefficients } \in \mathbb{R} \};$
- L^{∞} norm (also called supremum norm):

$$||g||_{\infty} = \sup_{x \in [a,b]} |g(x)|;$$

- function f, interval [a, b];
- first, we look for $P^* \in \mathcal{P}_n$ such that

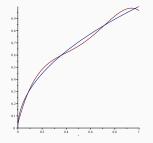
$$||f-P^*||_{\infty}=\min_{Q\in\mathcal{P}_n}||f-Q||_{\infty}.$$

Continuous world: everything was done in the 19th and early 20th centuries

Theorem 1 (Weierstrass, 1885)

Let f be a continuous function on [a,b]. For any $\epsilon>0$ there exists a polynomial p such that $\|p-f\|_{\infty,[a,b]}\leq \epsilon$.

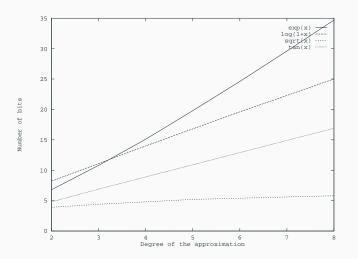
- → so it's not hopeless... but the theorem gives no clue on the necessary degree n to reach a given accuracy;
- for some functions it can be quite large.



The square root and its best degree-4 approximation in [0, 1].

(with the exponential function, the two curves would be undistinguishable at this scale)

Some functions may need high-degree polynomials



 $-\log_2(error)$ of the minimax polynomial approximations to various functions on [0,1], as a function of the degree.

Continuous world: everything was done in the 19th and early 20th centuries

Theorem 2 (Kirchberger 1902 – frequently attributed to Chebyshev)

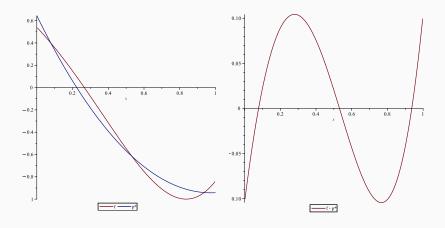
 $P^* \in \mathcal{P}_n$ is the minimax degree-n approximation to f on [a,b] if and only if there exist at least n+2 points

$$a < x_0 < x_1 < x_2 < \cdots < x_{n+1} < b$$

such that:

$$p^*(x_i) - f(x_i) = (-1)^i [p^*(x_0) - f(x_0)] = \pm ||f - p^*||_{\infty}.$$

Example: $f(x) = \cos(x + e^x)$, degree 2, interval [0,1]



Proof that the condition is sufficient let PE In that satisfies the condition Imagin have is a Qe Pa such that 11 Q-flo < 11 P-flos We necessarily have $|Q(x_i) - f(x_i)| < |P(x_i) - f(x_i)|$ home of P(xi) > g(xi) then Q(xi) < P(xi)

) if P(xi) < g(xi) then Q(xi) > P(xi) as the sign of P-J changes between α_{i} and α_{i+1} , so does the says of α_{i} -P→ Q-P has a noot between x; and x;+1

m+2 points → m+1 roots
while is imposable, since Q-P is a ronzero phynomial of degree < m.

Continuous world: everything was done in the 19th and early 20th centuries

Remez algorithm (1934): iteratively builds the set of points x_0, \ldots, x_{n+1} of Kirchberger's theorem.

- Start from an initial set of points $x_0, x_1, \ldots, x_{n+1}$ in [a, b]. (can be arbitrary, but $x_i = \frac{a+b}{2} + \frac{(b-a)}{2} \cos\left(\frac{i\pi}{n+1}\right), 0 \le i \le n+1$, called Chebyshev points, is in general a good choice)
- Consider the linear system of equations

$$\begin{cases}
p_0 + p_1 x_0 + p_2 x_0^2 + \dots + p_n x_0^n - f(x_0) &= +\epsilon \\
p_0 + p_1 x_1 + p_2 x_1^2 + \dots + p_n x_1^n - f(x_1) &= -\epsilon \\
p_0 + p_1 x_2 + p_2 x_2^2 + \dots + p_n x_2^n - f(x_2) &= +\epsilon \\
\dots &\dots &\dots \\
p_0 + p_1 x_{n+1} + \dots + p_n x_{n+1}^n - f(x_{n+1}) &= (-1)^{n+1} \epsilon.
\end{cases}$$
(1)

n+2 equations, with n+2 unknowns: $p_0, p_1, p_2, \ldots, p_n$ and ϵ . Nonzero determinant (Vandermonde matrix) \rightarrow exactly one solution $(p_0, p_1, \ldots, p_n, \epsilon)$.

Continuous world: everything was done in the 19th and early 20th centuries

- \rightarrow gives a polynomial $P(x) = p_0 + p_1 x + \cdots + p_n x^n$.
- We now compute the set of points y_i in [a, b] where P f has its extremes, and we start again (step 2), replacing the $x_i's$ by the y_i 's.

In practice: extremely fast convergence.

Illustration: degree-4 approximation to sin(exp(x)) in [0,2]

- we start with 0, 0.1909830057, 0.6909830062, 1.309016994,
 1.809016994, 2 (heuristic: Chebyshev points);
- the corresponding linear system is

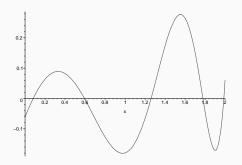
```
 \begin{cases}      \rho_0 & -0.8414709848 \\      \rho_0 & +0.1909830057\rho_1 \\      & +0.003647450847\rho_2 \\      & +0.00133038977\rho_4 \\      & +0.099830062\rho_1 \\      & +0.4774575149\rho_2 \\      & +0.2279656785\rho_4 \\      & +0.29150289\rho_3 \\      & +0.29150289\rho_3 \\      & +0.29150289\rho_3 \\      & +0.2910694\rho_1 \\      & +1.713525491\rho_2 \\       & +2.936169607\rho_4 \\      & +0.5319820928 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\      & +0.77912944 \\
```

• solving this system gives the polynomial:

$$P^{(1)}(x) = 0.7808077493 + 1.357210937x$$
$$-0.7996276765x^2 - 2.295982186x^3 + 1.189103547x^4.$$

Illustration: degree-4 approximation to sin(exp(x)) in [0, 2]

• difference $P^{(1)}(x) - \sin(\exp(x))$:

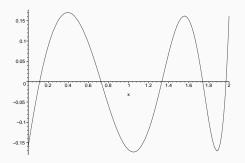


- extremes of $P^{(1)}(x) \sin(\exp(x))$ in [0,2]: 0, 0.3305112886, 0.9756471625, 1.554268282, 1.902075854, 2.
- Solving the linear system associated to this list of points gives the polynomial:

$$P^{(2)}(x) = 0.6800889007 + 2.144092090x$$
$$-1.631367834x^{2} - 2.226220290x^{3} + 1.276387351x^{4}.$$

Illustration: degree-4 approximation to sin(exp(x)) in [0,2]

• difference $P^{(2)}(x) - \sin(\exp(x))$:



- the extreme values of $|P^{(2)}(x) \sin(\exp(x))|$ are very close together: $P^{(2)}$ already "almost" satisfies the condition of Kirchberger's theorem;
- extremes of $P^{(2)}(x) \sin(\exp(x))$ in [0,2]: 0, 0.3949555564, 1.048154245, 1.556144609, 1.879537115, 2 \rightarrow the linear system associated to these points gives a polynomial $P^{(3)}$, etc.

But the world of Floating-Point numbers is not continuous. . .

If we want ultimate accuracy, i.e., correct rounding (or even just an error not larger than $\approx 1\,\text{ulp}$), this will not work:

- by approximating the coefficients of the minimax polynomial by FP numbers, we already loose a significant amount of accuracy;
- need to take into account the error due to evaluating the polynomial.

Example: function log(1+x) in [0,1]

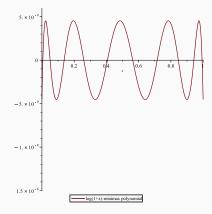
- Goal: binary32/single precision arithmetic ($\beta = 2$, p = 24);
- Withe the Remez algorithm, we get a degree-9 polynomial;

```
P^*(x) = 4.5312626764178853045083975073772119446839784174464956459529643270987 \cdot 10^{-9} \\ +0.999999025258539069369331280889443153540235373687738534675639041879165305 \cdot x \\ -0.4999653017333068626463502632864949694222430843003074565833086084942028155 \cdot x^2 \\ +0.332849806189948672280944634938400232746236110583281184032142726147068820 \cdot x^3 \\ -0.246517965788047515383441886060448202153937598635915727063060602574084424 \cdot x^4 \\ +0.185154569966668361477975368769654445113328710403271118930768041265705015 \cdot x^5 \\ -0.126057455452415889884958065630540564789570519172717287039473166654700319 \cdot x^6 \\ +0.0672945529992441805619775364650010239787939256621502918569305409951238042 \cdot x^8 \\ -0.02345535069051077538074324049760472948991145589046990280611685936442020761 \cdot x^8 \\ +0.00384529980982606902249675587076617855256867203322449811715889629373844490 \cdot x^9 .
```

• error $\approx 4.53 \times 10^{-9}$.

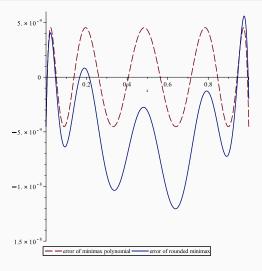
Example: function log(1+x) in [0,1]

• error curve $\log(1+x) - P^*(x)$: nice illustration of Kirchberger's theorem



- Problem: P* has real coefficients;
- I want to implement the function in binary32 arithmetic ($\beta = 2$, p = 24). What happens if I round each coefficient of P^* to the nearest FP number?

The disaster...



Error multiplied by ≈ 2.66

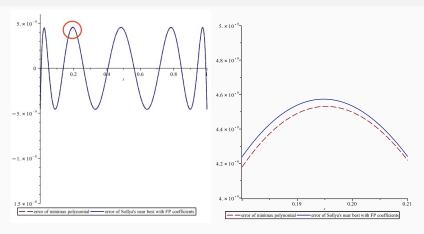
Near best polynomial approximations with FP coefficients

- algorithms that give near-best polynomials approximations under the constraint that the coefficients be FP numbers;
- more complex than Remez yet reasonable cost (and done once for all): Brisebarre and Chevillard, Efficient Polynomial L[∞]-Approximations, 18th IEEE Symposium on Computer Arithmetic, 2007. https://hal.inria.fr/inria-00119513
- still very active domain: various other constraints on coefficients, norms other than L^{∞} , rational approximations, taking into account the evaluation error when choosing the polynomial...
- Sollya software (Chevillard, Joldes, Lauter). Widely used for designing function libraries.

https://www.sollya.org

It also provides a Certified supnorm: proven bound on $||f - p||_{\infty}$.

With our example of function log(1+x) in [0,1]



Sollya prompt: P3 = fpminimax(log(1+x),9,[|24...|],[0;1],absolute); Error Sollya / Error minimax ≈ 1.012

ightarrow we recover almost all the loss due to the discretization of the polynomial.

Polynomial evaluation error

- Horner scheme
- use of interval arithmetic to know where each intermediate variable can lie;
- standard model to bound the error of each operation;
- can be very accurate if initial interval cut into many subintervals (a separate study for each subinterval);
- Gappa (https://gappa.gitlabpages.inria.fr), designed by G.
 Melquiond does this automatically and can generate a formal proof.

But what should be recommended?

- reproducibility, portability \rightarrow the result of $\sin(x)$, $\cos(x)$, $\exp(x)$... should be uniquely specified (no "fuzzy" specification of the form "error $< \epsilon$ ");
- specifying the result of an algorithm? Dangerous:
 - might make progress of elementary function algorithms difficult;
 - if somebody comes up with functions "better than the standard",
 the standard is dead

 \rightarrow what should be required is that we return the best possible result, i.e., correct rounding of the exact result.

However: difficult because of the table maker's dilemma.



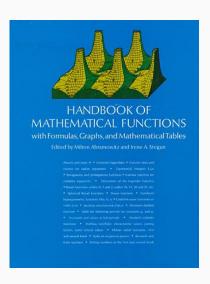


Table 4.2		NATUDAL	LOGARITHMS		
Table 4.2		NATURAL	LOGARTIMS		
\boldsymbol{x}	$\ln x$	x	$\ln x$	\boldsymbol{x}	$\ln x$
0.901 -0.10 0.902 -0.10 0.903 -0.10	0536 05156 578263 0425 00213 737991 0314 07589 195134 0203 27255 651516 0092 59185 899606	0.951 -0.05 0.952 -0.04 0.953 -0.04	129 32943 875505 024 12164 367467 919 02441 907717 814 03753 279349 709 16075 338505	1.000 1.001 1.002 1.003 1.004	0.00000 00000 000000 0.00099 95003 330835 0.00199 80026 626731 0.00299 55089 797985 0.00399 20212 695375
0.906 -0.09 0.907 -0.09 0.908 -0.09	9982 03352 822109 9871 59729 391577 9761 28288 670004 9651 09003 808438 9541 01848 046582	0.956 -0.04 0.957 -0.04 0.958 -0.04	604 39385 014068 499 73659 307358 395 18875 291828 290 75010 112765 186 42040 986988	1.005 1.006 1.007 1.008 1.009	0.00498 75415 110391 0.00598 20716 775475 0.00697 56137 364252 0.00796 81696 491769 0.00895 97413 714719
0.911 -0.09 0.912 -0.09 0.913 -0.09	9431 06794 712413 9321 23817 221787 9211 52889 078057 9101 93983 871686 3992 47075 279870	0.961 -0.03 0.962 -0.03 0.963 -0.03	082 19945 202551 978 08700 118446 874 08283 164306 770 18671 840115 666 39843 715914	1.010 1.011 1.012 1.013 1.014	0.00995 03308 531681 0.01093 99400 383344 0.01192 85708 652738 0.01291 62252 665463 0.01390 29051 689914
0.916 -0.08 0.917 -0.08 0.918 -0.08	3883 12137 066157 3773 89143 080068 3664 78067 256722 3555 78883 616466 3446 91566 264500	0.966 -0.03 0.967 -0.03 0.968 -0.03	562 71776 431511 459 14447 696191 355 67835 288427 252 31917 055600 149 06670 913708	1.015 1.016 1.017 1.018 1.019	0.01488 86124 937507 0.01587 33491 562901 0.01685 71170 664229 0.01783 99181 283310 0.01882 17542 405878
0.921 -0.08 0.922 -0.08 0.923 -0.08	3338 16089 390511 3229 52427 268302 3121 00554 255432 3012 60444 792849 7904 32073 404529	0.971 -0.02 0.972 -0.02 0.973 -0.02	045 92074 847085 942 88106 908121 839 94745 216980 737 11967 961320 634 39753 396020	1.020 1.021 1.022 1.023 1.024	0.01980 26272 961797 0.02078 25391 825285 0.02176 14917 815127 0.02273 94869 694894 0.02371 65266 173160

Consider the binary64 FP number ($\beta = 2, p = 53$)

$$x = \frac{8520761231538509}{2^{62}}$$

We have

$$2^{x} = \Big(9018742077413030.99999999999999998805240837303\cdots\Big) \times 2^{-53}$$

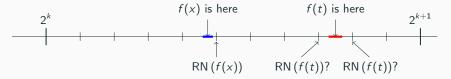
So what ?

Hardest-to-round case for function 2^x and binary64 FP numbers.

Correct rounding of the elementary functions

- radix 2, precision p;
- FP number x and integer m (with m > p) \rightarrow one can compute an approximation y to f(x) whose error on the significand is $\leq 2^{-m}$.
- can be done with a possible wider format, or using algorithms such as TwoSum, TwoMult, Double-Word (or triple-word) arithmetic, etc.
- getting a correct rounding of f(x) from y: not possible if f(x) is too close to a breakpoint: a point where the rounding function changes.

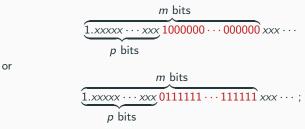
Correct rounding of the elementary functions



From the knowledge that f(x) lies in the blue interval, we can deduce the value of RN (f(x)). However, knowing that f(t) lies in the red interval does not allow us to know if RN (f(t)) is the FP number below f(t) or the FP number above it.

We are in trouble when f(x) has the form

RN rounding function,

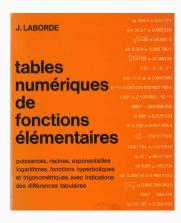


other rounding functions,

or



An example from the sixties



- no pocket calculators → tables of functions + interpolation were very useful;
- when designing his tables, Laborde was, for rare input values, unable to decide in which direction a result needed to be rounded;
- for these values, he used to choose a random rounding direction, and he would carefully record them to detect plagiarism of his tables.

Lindemann's theorem (1882)

- algebraic number: root of a nonzero polynomial with integer coefficients;
- the algebraic numbers are a field (proof not so easy); the $\sqrt{}$ of an algebraic number is algebraic (straightforward: $\sum a_i X^i \to \sum a_i X^{2i}$);
- transcendental number = not algebraic. Historical examples: $\sum_{k=0}^{\infty} 10^{-k!}$, π , e;
- ullet the FP numbers are rational \Rightarrow they are algebraic.

Theorem 3

If $z \in \mathbb{C}$ is a nonzero algebraic number then e^z is transcendental.

Consequences of Lindemann's theorem

- The sine, cosine of a nonzero algebraic number is transcendental.
 - x algebraic nonzero \rightarrow ix algebraic nonzero \rightarrow $u = e^{ix}$ transcendental;
 - if sin(x) was algebraic, then w = 2i sin(x) = u 1/u would be algebraic too;

(note that
$$w = u - 1/u$$
 implies that $u^2 - wu - 1 = 0$)

- therefore $(w \pm \sqrt{w^2 + 4})/2$ would be algebraic too;
- but $u \in (w \pm \sqrt{w^2 + 4})/2 \rightarrow$ contradiction;
- obviously similar for hyperbolic sine and cosine;
- true for the tangent function, since $\sin(x) = \tan(x)/\sqrt{\tan^2 x + 1}$;
- obviously true for their inverses: log, arcsin, arctan, ...

Finding *m* beyond which there is no problem ?

- function f: sin, cos, arcsin, arccos, tan, arctan, exp, log, sinh, cosh,
- Lindemann's theorem \rightarrow except for straightforward cases (e^0 , ln(1), $sin(0), \ldots$), if x is a FP number, there exists an m, say m_x , s.t. rounding the m_x -bit approximation \Leftrightarrow rounding f(x);
- finite number of FP numbers $\to \exists m_{\max} = \max_x (m_x)$ s.t. $\forall x$, rounding the m_{\max} -bit approximation to f(x) is equivalent to rounding f(x);
- ullet this reasoning does not give any hint on the order of magnitude of m_{\max} . Could be huge.

A bound derived from a result due to Baker (1975)

- $\alpha = i/j$, $\beta = r/s$, with $i, j, r, s < 2^p$;
- $C = 16^{200}$:

$$|\alpha - \log(\beta)| > (p2^p)^{-Cp\log p}$$

Application: To evaluate In et exp in double precision (p = 53) with correct rounding, it suffices to compute an approximation accurate to around

$$10^{244}$$
 bits

0000ps...

Some improvement

Definition 4 (Weil height)

Let α be an algebraic number of degree n and $P(x) = \sum_{i \leq n} p_i x^i$ be its minimal polynomial. Let $P(x) = p_n \cdot \prod_{i \leq n} (x - \alpha_i)$ be the factorization of P over the complex numbers. Then the Weil height of α is

$$H(lpha) = \left(p_n \cdot \prod_{i \leq n} \max(1, |lpha_i|)\right)^{rac{1}{n}}.$$

Note: if $\alpha = b/a \in \mathbb{Q}$ with gcd(a, b) = 1 then $H(\alpha) = max\{|a|, |b|\}$.

Some improvement

Theorem 5 (Y. Nesterenko and M. Waldschmidt, specialized here to the rational numbers)

Let α and α' be rational numbers. Let θ be an arbitrary non-zero real number. Let A,A', and E be positive real numbers with

$$E \ge e$$
, $A \ge \max(H(\alpha), e)$, $A' \ge H(\alpha')$.

Then

$$\begin{split} &\left| e^{\theta} - \alpha \right| + \left| \theta - \alpha' \right| \geq \\ &\exp \left\{ -211 \cdot \left(\ln A' + \ln \ln A + 2 \ln (E \cdot \max\{1, |\theta|\}) + 10 \right) \right. \\ &\cdot \left(\ln A + 2E |\theta| + 6 \ln E \right) \cdot \left(3.7 + \ln E \right) \cdot \left(\ln E \right)^{-2} \right\}. \end{split}$$

^aHere, $e = 2.718 \cdots$ is the base of the natural logarithms.

Some improvement

For the evaluation of exponentials in binary64 (p = 53), we find that m = 7,290,678 suffices.

Impossible \rightarrow too expensive

In recent years, significant further improvements of the theoretical bounds, but they remain quite large. . .

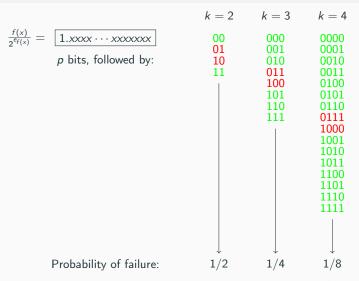
Best current theoretical results

Estimates of current theoretical upper bounds or the hardness to round for exp, trigonometric and hyperbolic functions in the binary128 format. For each function f, we report the values θ_f such that, over a given binade, the hardness to round f is less than $\theta_f \cdot 113$.

Binade	exp	sin	cos	sinh	cosh	tan	cot	tanh	coth
[1/8, 1/4)	226	1371	1359	688	682	303	302	303	298
[1/4, 1/2)	297	2070	2058	1062	1057	410	410	409	405
[1/2,1)	403	3288	3281	1698	1695	604	604	600	598
[1, 2)	593	5678	6481	2889	2889	1194	1196	931	929
[2,4)	920	11408	10266	5285	5285	1854	1855	1507	1504
[4,8)	1485	20395	20395	10155	10155	3361	3360	2634	2631

But in practice it's much less. Let us see why.

Approximation with error $\leq 2^{-p-k+1}$ on the significand



k bits \rightarrow probability of failure 2^{1-k}

Rule of thumb

If we approximate the significand of f(x) with error $\leq 2^{-p-k+1}$ (roughly speaking, if we compute a p+k-bit approximation to f(x)), the probability of not being able to deduce RN (f(x)) is around 2^{1-k} .

• exceptions to that rule: if x is tiny, not all bit strings are possible in sin(x), exp(x), etc. just after the first p bits. For instance,

$$sin(1.xxxx \cdots x1 \times 2^{-p}) = 1.xxxx \cdots x011111111111 \cdots \times 2^{-p}.$$

- in practice this is not a problem, just choose polynomial approximations where the lowest order term is exactly x for sin or sinh or $\log(1+x)$, 1 for $\exp(x)$, etc. They will automatically deliver correct rounding when x is tiny enough;
- the rule is essential for designing efficient algorithms;

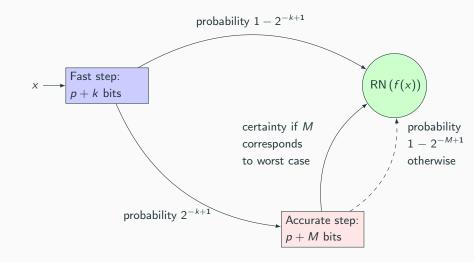
Expected vs actual worst cases

- Rule of thumb \rightarrow if w is the word size, as there are $\approx 2^w$ FP numbers, with a p+k-bit approximation there is a total of 2^{w+1-k} failures \rightarrow vanishes as soon as $k \approx w+1$.
- frequently less in practice: correlations (e.g. log₂), function not defined in full range (exp because of overflow, arcsin, etc.);
- actual values (excluding tiny trivial input values):

format	p+w+1	exp	log ₂	arcsin
binary32	57	52	51	$54 (x > 2^{-23})$
binary64	118	113	108	126 ($ x > 2^{-25}$)

 \rightarrow the rule of thumb is not that bad.

2-step process (Ziv's strategy) – typically, $k \approx 10$



Algorithm that compute the hardest-to-round cases

- active domain since ≈ 2000 ;
- best algorithms based on lattice reduction;
- time of computation of hardest-to-round cases remains an exponential function of p;
- p = 24 (binary32) is easy, p = 53 (binary64) is feasible but very costly, larger formats out of reach.

Algorithm that compute the hardest-to-round cases

Table 1: Worst cases for exponentials of binary64 FP numbers.

Interval	worst case (binary)
$[-\infty, -2^{-30}]$	$\exp(-1.1110110100110001100011101111110110110001001111$
	= 1.11111111111111111111111111111111111
[-2 ⁻³⁰ , 0)	$\exp(-1.000000000000000000000000000000000000$
	$= 1.111111111111111 \cdots 1111111111111111100 0 0^{100}1010 \times 2^{-1}$
(0, +2 ⁻³⁰]	$\exp(1.11111111111111111111111111111111111$
	= 1.000000000000000000000000000000000000
[2 ⁻³⁰ , +∞]	$\exp(1.011111111111111111111111111111111111$
	= 1.0000000000000000000000000000000010111111
	$\exp(1.100000000000001011111111111111111111$
	= 1.00000000000000000000000000000110000000
	$\exp(1.100111101001110010111011111111111010110000$
	= 1.00000000000000000000000000000011001111010
	exp(110.000011110101001011110011011110101111011001111
	= 110101100.0101000010110100000010011100100

Results

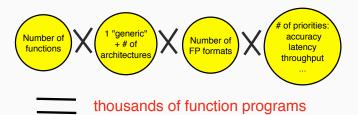
 Table 2: Worst cases for logarithms of binary64 FP numbers.

Interval	worst case (binary)
[2 ⁻¹⁰⁷⁴ , 1)	$\log(1.111010100111000111011000010111001110$
[2 ,1)	= -101100000.0010100101101010011001101010100001011111
	$\log(1.1001010001110110111000110000010011001$
	= -100001001.1011011000001100101011110100011110110
	$\log(1.00100110111010011100010011010011001001$
	= -10100000.101010110010110000100101111001101000010000
	$\log(1.0110000100111001010101011110111001000000$
	= -10111.1111000000101111100110111101011110110
(1, 2 ¹⁰²⁴]	$\log(1.01100010101010001000011000010011011000101$
	= 111010110.010001111001111010111101001111100100

The CORE-Math library

- developed in Nancy by Paul Zimmermann and colleagues;
- all binary32 and binary64 functions from the C23 standard except compound and Igamma, with correct rounding;
- binary32 acos, acosh, asin, asinh, atan, atan2, atanh, cbrt, cosh, erf, erfc, expm1, exp2m1, exp10m1, tgamma, lgamma, log10, log1p, log2p1, log10p1, sinh, tan, tanh functions integrated into GNU libc
- code and an extensive bibliography are available from https://core-math.gitlabpages.inria.fr/

Libraries of math functions



- impossible to debug, maintain, keep consistent, improve. . .
- and physicists would like many other functions

First solution: computer-assisted library design

Metalibm project (http://www.metalibm.org), launched by Florent de Dinechin. Two versions

- fully automated for the end user;
- assistance for the specialist.

Metalibm builds upon tools such as Sollya and Gappa.

But this is not the ultimate goal

Generation of functions at compile-time

- take into account the exact context: underlying architecture, accuracy requirements, priorities (latency/throughput);
- possibly, information on input domain (→ simplify/avoid range reduction), or special cases (e.g., infinities, NaNs known not to happen);
- compound functions: if you need

$$E_4(x) = \frac{x}{e^x - 1} - \ln(1 - e^{-x}),$$

then you directly generate $E_4(x)$ instead of generating exp, In and combining them.

- formal proof absolutely necessary (no library to heavily test beforehand);
- need collaboration of people from computer arithmetic, mathematics, computer algebra, compilation, formal proof...