

Worst Cases for Correct Rounding of the Elementary Functions in Double Precision

Vincent Lefèvre

INRIA, Projet Spaces, LORIA, Campus Scientifique
B.P. 239, 54506 Vandoeuvre-lès-Nancy Cedex, FRANCE

`Vincent.Lefevre@loria.fr`

Jean-Michel Muller

CNRS, Projet CNRS/ENS Lyon/INRIA Arnaire
LIP, Ecole Normale Supérieure de Lyon
46 Allée d'Italie, 69364 Lyon Cedex 07, FRANCE

`Jean-Michel.Muller@ens-lyon.fr`

August 14, 2003

Abstract

We give the results of our search for the worst cases for correct rounding of the major elementary functions in double precision floating-point arithmetic. These results allow the design of reasonably fast routines that will compute these functions with correct rounding, at least in some interval, for any of the four rounding modes specified by the IEEE-754 standard. They will also allow one to easily test libraries that are claimed to provide correctly rounded functions.

Keywords Elementary functions, floating-point arithmetic, computer arithmetic, Table Maker's Dilemma.

1 Introduction

In general, the result of an arithmetic operation on two floating-point (FP) numbers is not exactly representable in the same FP format: it must be *rounded*. In a FP system that follows the IEEE 754 standard [2, 6], the user can choose an *active rounding mode* from:

- rounding towards $+\infty$, or upwards: $RU(x)$ is the smallest FP number that is greater than or equal to x ;
- rounding towards $-\infty$, or downwards: $RD(x)$ is the largest FP number that is less than or equal to x ;
- rounding towards 0: $RZ(x)$ is equal to $RU(x)$ if $x < 0$, and to $RD(x)$ otherwise;
- rounding to the nearest even: $RN(x)$ is the FP number that is closest to x . If x is exactly halfway between two consecutive FP numbers, $RN(x)$ is the one whose last mantissa bit is a zero.

The standard requires that the system should behave as if the results of the operations $+$, $-$, \div , \times and \sqrt{x} were first computed with “infinite precision”, and then rounded accordingly to the active rounding mode. Operations that satisfy this property are called correctly (or exactly) rounded.

Correctly rounding these operations is easily feasible. Consider for instance division. One can show that the quotient of two n -bit numbers either is exactly representable as an n -bit number, or cannot contain a string of more than $n - 1$ consecutive zeros or ones in its binary representation. From this, we deduce that rounding an approximation to such a quotient with $2n$ bits of accuracy will always be equivalent to rounding the exact result. Information on correct rounding of the algebraic functions¹ can be found in references [3, 8].

Unfortunately, there is no such requirement for the elementary functions², because the accuracy with which these functions must be computed to make sure that we can round them correctly was

¹Function f is algebraic if there exists a 2-variable polynomial P with integer coefficients such that $y = f(x) \Leftrightarrow P(x, y) = 0$.

²By *elementary functions* we mean the radix 2, e and 10 logarithms and exponentials, and the trigonometric and hyperbolic functions.

not known for double and double extended precisions (for single precision, since checking 2^{32} input numbers is rather quickly done, there already exist libraries that provide correct rounding. See for instance [14]).

Requiring correctly rounded results would not only improve the accuracy of computations: it would help to make numerical software more portable. Moreover, as noticed by Agarwal et al. [1], correct rounding facilitates the preservation of useful properties such as monotonicity, symmetry³ and important identities. See [13] for more details.

Before going further, let us start with definitions.

Definition 1 (Infinite mantissa) We call Infinite mantissa of a nonzero real number x the number

$$\mathcal{M}_\infty(x) = x/2^{\lfloor \log_2 |x| \rfloor}.$$

$\mathcal{M}_\infty(x)$ is the real number x' such that $1 \leq x' < 2$ and $x = x' \times 2^k$, where k is an integer.

For instance,

$$\mathcal{M}_\infty(1/3) = 4/3 = 1.01010101 \dots$$

(in binary) and

$$\mathcal{M}_\infty(\pi) = \pi/2 = 1.1001001000011111101101 \dots$$

Definition 2 (Mantissa distance) If x is a FP number, then $\mathcal{M}_\infty(x)$ is the mantissa of its FP representation.

If a and b belong to the same "binade" (they have the same sign and satisfy $2^p \leq |a|, |b| < 2^{p+1}$, where p is an integer), we call their Mantissa distance the distance between their infinite mantissas, that is, $|a - b|/2^p$.

For instance, the mantissa distance between $7/2$ and $8/3$ is $5/12 = 0.416666 \dots$

Let f be an elementary function and x a FP number. Unless x is a very special case – e.g., $\log(1)$, $\arccos(1)$ or $\sin(0)$ –, $y = f(x)$ cannot be exactly computed. The only thing we can do is to compute an approximation y^* to y . If we wish to provide correctly rounded functions, we have to know what the accuracy of this approximation should be to make sure that rounding y^* is equivalent to rounding y . In other words, from y^* and the known bounds on the approximation, the only information we have is that y belongs to some interval Y .

³With the RN and RZ modes.

Definition 3 (breakpoint) Let us call \diamond the rounding function. We call a **breakpoint** a value z where the rounding changes, that is, if t_1 and t_2 are real numbers satisfying $t_1 < z < t_2$ then $\diamond(t_1) < \diamond(t_2)$.

For “directed” rounding modes (i.e., towards $+\infty$, $-\infty$ or 0), the breakpoints are the FP numbers. For rounding to the nearest (RN) mode, they are the exact middle of two consecutive FP numbers.

If Y contains a breakpoint, then we cannot provide $\diamond(y)$: the computation must be carried again with a larger accuracy. There are two ways of solving that problem:

- iteratively increase the accuracy of the approximation, until interval Y no longer contains a breakpoint⁴. The problem is that it is difficult to predict how many iterations will be necessary;
- compute, only once and in advance, the smallest nonzero mantissa distance between the image⁵ of a FP number and a breakpoint. This makes it possible to deduce the accuracy with which $f(x)$ must be approximated, in the worst case, to make sure that rounding the approximation is equivalent to rounding the exact result.

The first solution was suggested by Ziv [15]. It has been implemented in a library available through the internet⁶. The last iteration uses 768 bits of precision. Although there is no formal proof that this suffices (the results presented in this paper actually give the proof for the functions and domains considered here), elementary probabilistic arguments [4, 5, 13] show that requiring a larger precision is *extremely* unlikely.

We decided to implement the second solution, since the only way to implement the first one safely is to overestimate the accuracy that is needed in the worst cases.

For some algebraic functions, it is possible to directly build worst cases as a function of the number n of mantissa bits of the floating-point format being considered. For instance, for the reciprocal

⁴This is not possible if $f(x)$ is equal to a breakpoint. And yet, $x = 0$ is the only FP input value for which $\sin(x)$, $\cos(x)$, $\tan(x)$, $\arctan(x)$ and e^x have a finite radix-2 representation – and the breakpoints *do* have finite representations –, and $x = 1$ is the only FP input value for which $\ln(x)$, $\arccos(x)$ and $\cosh^{-1}(x)$ have a finite representation. Concerning 2^x and 10^x , 2^x has a finite representation if and only if x is an integer, and 10^x has a finite representation if and only if x is a nonnegative integer. Also, $\log_2(x)$ (resp. $\log_{10}(x)$) has a finite representation if and only if x is an integer power of 2 (resp. 10). All these cases are straightforwardly handled separately, so we do not discuss them in the rest of the paper.

⁵We call *image* of x the number $f(x)$, where f is the elementary function being considered.

⁶<http://www.alphaWorks.ibm.com/tech/mathlibrary4java>.

Table 1: Worst cases for the exponential function between 1 and 2 and small values of n . They have been computed through exhaustive searching

n	worst case x	$\exp(x)$
5	29/16	110.0010000000110...
6	57/32	101.1110111111100...
7	47/32	100.010110000000011...
8	47/32	100.010110000000011...
9	135/128	10.1101111011111011...
10	937/512	110.00111011111111010...
11	1805/1024	101.110100111111111011...
12	1805/1024	101.110100111111111011...
13	5757/4096	100.000100111101111111000...
14	12795/8192	100.11000100100011111110101...

function $1/x$, a worst case is $x = 1 - 2^{-n}$. Unfortunately, when we deal with the transcendental functions, there is no known way of directly getting worst cases: the only known solution is to try all possible input numbers. Table 1 gives the worst cases for the exponential function between 1 and 2 and small values of n . Nobody succeeded in finding some “regularity” in such a table, that could have helped to predict worst cases for larger values of n . The basic principle of our algorithm for searching the worst cases was outlined in [11]. We now present properties that have allowed us to hasten the search, as well as the results obtained after having run our algorithms on several workstations, and consequences of our results. The results we have obtained are worst cases for the Table Maker’s Dilemma, that is, FP numbers whose image is closest to a breakpoint. For instance, the worst case for the natural logarithm in the full double precision range is attained for

$$x = 1.011000101010100010000110000100110110001010 \\ 0110110110 \times 2^{678}$$

whose logarithm is

$$\log x = \overbrace{111010110.0100011110011110101 \dots 110001}^{53 \text{ bits}} \\ \underbrace{000000000000000000 \dots 0000000000000000}_{65 \text{ zeros}} 1110\dots$$

This is a “difficult case” in a directed rounding mode, since it is very near a FP number. One of the two worst cases for radix-2 exponentials in the full double precision range is

$$\frac{1.1110010001011001011001010010011010111111}{100101001101} \times 2^{-10}$$

whose radix-2 exponential is

$$\frac{\overbrace{1.0000000001010011111111000010111 \dots 0011}^{53 \text{ bits}}}{\underbrace{011111111111111111 \dots 11111111111111110100 \dots}_{59 \text{ ones}}}$$

It is a difficult case for rounding to the nearest, since it is very close to the middle of two consecutive FP numbers.

2 Our Algorithms for Finding the Worst Cases

2.1 Basic principles

The basic principles are given in [11], so we only quickly describe them and focus on new aspects. Assume we wish to find the worst cases for function f in double precision. Let us call **test number** (TN) a number that is representable with 54 bits of mantissa (it is either a FP number or the exact middle of two consecutive FP numbers). The TNs are the values that are breakpoints for one of the rounding modes. Finding worst cases now reduces to the problem of finding FP numbers x such that $f(x)$ is closest (for the mantissa distance) to a TN. We proceed in two steps: we first use a fast “filtering” method that eliminates all points whose distance to the closest breakpoint is above a given threshold. The value of the threshold is chosen so that this filtering method does not require highly accurate computations, and so that the number of values that remain to be checked after the filtering is so small that an accurate computation of the value of the function at each remaining value is possible. Details on the choice of parameters are given in [10].

In [11], we suggested to perform the filtering as follows:

- first, the domain where we look for worst cases is split into “large subdomains” where all input values have the same exponent;

- each large subdomain is split into “small subdomains” that are small enough so that in each of them, within the accuracy of the filtering, the function can be approximated by a linear function. Hence in each small subdomain, our problem is to find a point on a grid that is closest to a straight line. We solve a slightly different problem: given a threshold ϵ we just try to know if there can be a point of the grid at distance less than ϵ from the straight line. ϵ is chosen so that for one given small subdomain this event is very unlikely.
- using a variant to the Euclidean algorithm suggested by V. Lefèvre [9], we solve that problem. If we find that there can be a point of the grid at distance less than ϵ from the straight line, we check all points of the small subdomain.

2.2 Optimization: f and f^{-1} simultaneously

Let us improve that method. Instead of finding **floating-point numbers** x such that $f(x)$ is closest to a test number, we look for **test numbers** x such that $f(x)$ is closest to a TN. This makes it possible to compute worst cases for f and for its inverse f^{-1} in one pass only, since the image $f(a)$ of a breakpoint a is near a breakpoint b if and only if $f^{-1}(b)$ is near a . One could object that by checking the TNs instead of checking the double precision FP numbers only, we double the number of points that are examined. So getting in one pass the results for two functions (f and f^{-1}) seems to be a no-win no-loss operation. This is not quite true, since there are sometimes *much fewer* values to check for one of the two functions than for the other one.

Consider as an example the radix-2 exponential and logarithm, with input domain $I = [-1, 1]$ for 2^x , which corresponds to input domain $J = [1/2, 2]$ for $\log_2(y)$. The following two strategies would lead to the same final result: the worst cases for 2^x in I and for $\log_2(y)$ in J .

1. check 2^x for every test number x in I ;
2. check $\log_2(y)$ for every test number y in J .

If we use the first strategy, we need to check all TNs of exponent between⁷ -53 and -1 . There are

⁷For numbers of smaller absolute value, there is no longer any problem of implementation: their radix-2 exponential is 1 or $1^- = 1 - \text{ulp}(1/2)$ or $1^+ = 1 + \text{ulp}(1)$ depending on their sign and the rounding mode.

106×2^{53} such numbers. With the second strategy, we need to check all positive TNs of exponent equal to -1 or 0 , that is, 2×2^{53} numbers. Hence, with the second strategy, we check approximately 53 times fewer values than with the first one.

If we separately check all FP numbers in I and all FP numbers in J , we check $106 \times 2^{52} + 2 \times 2^{52}$ numbers.

Hence, in the considered domain, it is better to check $\log_2(y)$ for every TN y in $[1/2, 2]$. In other domains, the converse holds: when we want to check both functions in the domain defined by $x > 1$ (for 2^x) or $y > 2$ (for $\log_2(y)$), we only have to consider 10 values of the exponent if we check 2^x for every TN in the domain, whereas we would have to consider 1022 values of the exponent if we decided to check $\log_2(x)$ for the TNs in the corresponding domain.

The decision whether it is better to base our search for worst cases on the examination of f in a given domain I or f^{-1} in $J = f(I)$ can be helped by examining $T_f(x) = |x \times f'(x)/f(x)|$ in I . If $T_f(x) \gg 1$, then I contains fewer test numbers than J , so it is preferable to check f in I . If $T_f(x) \ll 1$, it is preferable to check f^{-1} in J . When $T_f(x) \approx 1$, a more thorough examination is necessary. Of course, in all cases, another important point is which of the two functions is better approximated by a polynomial of small degree.

2.3 Optimization: special input values

For most functions, it is not necessary to perform tests for the input arguments that are extremely close to 0. For example, consider the exponential of a very small positive number ϵ , on a FP format with p -bit mantissas, assuming rounding to nearest. If $\epsilon < 2^{-p}$ then (since ϵ is a p -bit number), $\epsilon \leq 2^{-p} - 2^{-2p}$. Hence,

$$e^\epsilon \leq 1 + (2^{-p} - 2^{-2p}) + \frac{1}{2}(2^{-p} - 2^{-2p})^2 \dots < 1 + 2^{-p}.$$

therefore $\exp(\epsilon) < 1 + (1/2) \text{ulp}(1)$. Thus, the correctly rounded value of $\exp(\epsilon)$ is 1. A similar reasoning can be done for other functions and rounding modes. Some results are given in Tables 2 and 3.

There are also some special cases for which we have a small number of values only to check, so

Table 2: Some results for small values in double precision, assuming **rounding to the nearest**. These results make finding worst cases useless for negative exponents of large absolute value.

This function	can be replaced by	when
$\exp(\epsilon), \epsilon \geq 0$	1	$\epsilon < 2^{-53}$
$\exp(\epsilon), \epsilon \leq 0$	1	$ \epsilon \leq 2^{-54}$
$\ln(1 + \epsilon)$	ϵ	$ \epsilon < \sqrt{2} \times 2^{-53}$
$2^\epsilon, \epsilon \geq 0$	1	$\epsilon < 1.4426 \dots \times 2^{-53}$
$2^\epsilon, \epsilon \leq 0$	1	$ \epsilon < 1.4426 \dots \times 2^{-54}$
$\sin(\epsilon), \arcsin(\epsilon), \sinh(\epsilon), \sinh^{-1}(\epsilon)$	ϵ	$ \epsilon \leq 1.4422 \dots \times 2^{-26}$
$\cos(\epsilon)$	1	$ \epsilon < \sqrt{2} \times 2^{-27}$
$\cosh(\epsilon)$	1	$ \epsilon < 2^{-26}$
$\tan(\epsilon), \tanh(\epsilon), \arctan(\epsilon), \tanh^{-1}(\epsilon)$	ϵ	$ \epsilon \leq 1.817 \dots \times 2^{-27}$

that it can be done without using our programs. Consider, for the double-precision format, the cosine of numbers of absolute value less than 2^{-25} . Since $|x| < 2^{-25}$ implies $\cos(x) > 1 - 2^{-51}$, the only TNs in the area where $\cos(x)$ lies are the values $1 - k \times 2^{-54}$, for $k = 0, 1, \dots, 8$.

- the case $k = 0$, related to the TN 1, corresponds to the special case $\cos(0) = 1$. It is solved as follows:

- with rounding to the nearest, if $x \leq \text{RN}(2^{-27}\sqrt{2}) = \frac{6369051672525773}{604462909807314587353088}$ then return 1;
- with rounding towards $+\infty$, if $x \leq 2^{-26}$ then return 1.

- for each other value of k (from 1 to 8),

- compute

$$\begin{cases} t_k^- &= \text{RD}(\arccos(1 - k \times 2^{-54})) \\ t_k^+ &= \text{RU}(\arccos(1 - k \times 2^{-54})) \end{cases}$$

- compute

$$|\cos(t_k^-) - (1 - k \times 2^{-54})|$$

and

$$|\cos(t_k^+) - (1 - k \times 2^{-54})|$$

Table 3: Some results for small values in double precision, assuming **rounding towards** $-\infty$. These results make finding worst cases useless for negative exponents of large absolute value. x^- is the largest FP number strictly less than x .

This function	can be replaced by	when
$\exp(\epsilon), \epsilon \geq 0$	1	$\epsilon < 2^{-52}$
$\exp(\epsilon), \epsilon < 0$	$1^- = 1 - 2^{-53}$	$ \epsilon \leq 2^{-53}$
$\ln(1 + \epsilon), \epsilon \neq 0$	ϵ^-	$ \epsilon < \sqrt{2} \times 2^{-52}$
$2^\epsilon, \epsilon \geq 0$	1	$\epsilon < 1.4426 \dots \times 2^{-52}$
$2^\epsilon, \epsilon < 0$	$1^- = 1 - 2^{-53}$	$ \epsilon < 1.4426 \dots \times 2^{-53}$
$\sin(\epsilon), \sinh^{-1}(\epsilon), \epsilon > 0$	ϵ^-	$\epsilon \leq 1.817 \dots \times 2^{-26}$
$\sin(\epsilon), \sinh^{-1}(\epsilon), \epsilon \leq 0$	ϵ	$ \epsilon \leq 1.817 \dots \times 2^{-26}$
$\arcsin(\epsilon), \sinh(\epsilon), \epsilon \geq 0$	ϵ	$\epsilon \leq 1.817 \dots \times 2^{-26}$
$\arcsin(\epsilon), \sinh(\epsilon), \epsilon < 0$	ϵ^-	$\epsilon \leq 1.817 \dots \times 2^{-26}$
$\cos(\epsilon), \epsilon \neq 0$	$1^- = 1 - 2^{-53}$	$ \epsilon < 2^{-26}$
$\cosh(\epsilon), \epsilon \neq 0$	1	$ \epsilon < \sqrt{2} \times 2^{-26}$
$\tan(\epsilon), \tanh^{-1}(\epsilon), \epsilon \geq 0$	ϵ	$\epsilon \leq 1.4422 \dots \times 2^{-26}$
$\tan(\epsilon), \tanh^{-1}(\epsilon), \epsilon < 0$	ϵ^-	$ \epsilon \leq 1.4422 \dots \times 2^{-26}$
$\tanh(\epsilon), \arctan(\epsilon), \epsilon > 0$	ϵ^-	$\epsilon \leq 1.4422 \dots \times 2^{-26}$
$\tanh(\epsilon), \arctan(\epsilon), \epsilon \leq 0$	ϵ	$ \epsilon \leq 1.4422 \dots \times 2^{-26}$

to know if $x = t_k^-$ or t_k^+ is a hard-to-round input value for the cosine function.

The worst case among these ones appears for $k = 2$. It corresponds to

$$\cos(2^{-26}) = 0. \overbrace{111111111 \dots 1111}^{53 \text{ bits}} 0 \underbrace{00 \dots 00000000}_{54 \text{ zeros}} 1010 \dots$$

A similar study can be done for the hyperbolic cosine function. We get very similar results:

- in RN mode, if $x \leq \frac{9007199254740991}{604462909807314587353088} = 2^{-26} - 2^{-79}$ then return 1;
- in RZ mode, if $x \leq \frac{1592262918131443}{75557863725914323419136} = \text{RD}(2^{-26}\sqrt{2})$ then return 1;
- for the other values, the worst case for $x < 2^{-25}$ is

$$\cosh(2^{-26}) = \overbrace{1.00000000 \dots 0000}^{53 \text{ bits}} 1 \underbrace{00 \dots 00000000}_{55 \text{ zeros}} 1 \dots$$

2.4 Normal and subnormal numbers

Our algorithms assume that input and output numbers are normalized FP values. Hence, we have to check whether there exist normalized FP numbers x such that $f(x)$ is so small that we should return a subnormal number. To do that, we use a method based on the continued fraction theory, suggested by Kahan [7], and originally designed for finding the worst cases for range reduction. It gives the normalized FP number that is closest to an integer nonzero multiple of $\pi/2$. This number is $\alpha = 16367173 \times 2^{72}$ in single precision, and $\beta = 6381956970095103 \times 2^{797}$ in double precision. Therefore, $A = |\cos(\alpha)| \approx 1.6 \times 10^{-9}$ and $B = |\cos(\beta)| \approx 4.7 \times 10^{-19}$ are lower bounds on the absolute value of the sine, cosine and tangent of normalized single precision (for A) and double precision (for B) FP numbers. These values are larger than the smallest normalized FP numbers. Hence the sine, cosine and tangent of a normalized FP number can always be rounded to normalized FP numbers.

3 Implementation of the Method

3.1 Overview of the implementation

The tests are implemented in three steps:

1. As said above, the first step is a *filter*. It amounts to testing if 32 (in general) consecutive bits are all zeros⁸ thus keeping one argument out of 2^{32} , in average. This step is very slow and needs to be parallelized.
2. The 2nd step consists in reducing the number of worst cases obtained from the first step and grouping all the results together in the same file. This is done with a slower but more accurate test than in the 1st step. As the number of arguments has been drastically reduced, this step is performed on a single machine.

⁸These are the bits following the first 54 bits of the mantissa, unless the exponent of the output values changes in the tested domain.

3. The 3rd step is run by the user to restrict the number of worst cases. Results on the inverse function are also obtained. This step is very fast.

Our programs are written in C, and use the MPN routines of GMP.

3.2 Details on the first step

Let us give more details about the first step. The user chooses a function f , an exponent, a mantissa size (usually 53), and the first step starts as follows.⁹

- First, the tested interval is split into 2^{13} subintervals J_i containing 2^{40} TNs and f is approximated by polynomials P_i of degree d_i (~ 4 to 20) on J_i . For each i , we start with $d_i = 1$, and increase d_i until the approximation is accurate enough. P_i is expressed modulo the distance between two consecutive TNs, as we only need to estimate the bits following the rounding bit.
- Then, each J_i is split into subintervals $K_{i,j}$ containing 2^{15} arguments and P_i is approximated by degree-2 polynomials $Q_{i,j}$ on $K_{i,j}$, with 64-bit precision.
- On $K_{i,j}$: $Q_{i,j}$ is approximated by a degree-1 polynomial (by ignoring the degree-2 coefficient) and the variant of the Euclidean algorithm is used. If it fails, that is, if the obtained distance is too small, then:
 - $K_{i,j}$ is split into 4 subintervals $L_{i,j,k}$.
 - For each k : the Euclidean algorithm is used on $L_{i,j,k}$, and if it fails, the arguments are tested the one after the other, using two 64-bit additions for each argument.

The first step requires much more time than the other ones, thus it is parallelized (we use around 40 workstations, in background). As the calculations in different intervals are totally independent, there is no need for communications between the different machines. The workstations have primary users. We must not disturb them. So, the programs were written so that they can run with a low priority, automatically stop after a given time, and automatically detect when a machine is used and stop if this is the case.

⁹The numbers given here are just those that are generally chosen; other values may be chosen for particular cases.

4 Results: e^x and $\ln(x)$

For these functions, there is no known way of deducing the worst cases in a domain from the worst cases in another domain. And yet, we have obtained the worst cases for all possible double precision FP inputs. They are given in Tables 4 and 5. From these results we can deduce the following properties.

Property 1 (Computation of exponentials) *Let y be the exponential of a double-precision number x . If $|x| < 2^{-54}$ then the results given Tables 2 and 3 apply, so that correctly rounding e^x is easy. Otherwise, let y^* be an approximation to y such that the mantissa distance¹⁰ between y and y^* is bounded by ϵ .*

- for $|x| \geq 2^{-30}$, if $\epsilon \leq 2^{-53-59-1} = 2^{-113}$ then for any of the 4 rounding modes, rounding y^* is equivalent to rounding y ;
- for $2^{-54} \leq |x| < 2^{-30}$, if $\epsilon \leq 2^{-53-104-1} = 2^{-158}$ then rounding y^* is equivalent to rounding y ;
- for $|x| < 2^{-54}$, correctly rounding e^x consists in returning:
 - 1 in RN mode;
 - $1 + 2^{-52}$ in RU mode if $x > 0$;
 - 1 in RD or RZ mode if $x \geq 0$;
 - 1 in RU mode if $x \leq 0$;
 - $1 - 2^{-53}$ in RD or RZ mode if $x < 0$.

Property 2 (Computation of logarithms) *Let y be the natural (radix- e) logarithm of a double-precision number x . Let y^* be an approximation to y such that the mantissa distance between y and y^* is bounded by ϵ . If $\epsilon \leq 2^{-53-64-1} = 2^{-118}$ then for any of the 4 rounding modes, rounding y^* is equivalent to rounding y .*

¹⁰If one prefers to think in terms of relative error, one can use the following well-known properties: if the mantissa distance between y and y^* is less than ϵ then their relative distance $|y - y^*|/|y|$ is less than ϵ . If the relative distance between y and y^* is less than ϵ_r , then their mantissa distance is less than $2\epsilon_r$.

the bits that represent $p - 1$ (the same as for y), followed by the bits that represent $\log_2(m)$. But the bits that represent $1 - \log_2(m)$ are obtained by complementation¹¹ of the bits that represent $\log_2(m)$. Hence, there is a chain of k consecutive 1s (or 0s) after bit 54 of $\mathcal{M}_\infty(y)$ if and only if there is a chain of k consecutive 0s (or 1s) after bit 54 of $\mathcal{M}_\infty(y')$. Therefore, x is a worst case for input values < 1 if and only if x' is a worst case for input values > 1 . This is illustrated in Table 7: the infinite mantissa of the worst case for $x > 1$ starts with the same bit chain (1000000000) as the mantissa of the worst case for $x < 1$, then the bits that follow are complemented (1000100011111110...000110 0 0⁵⁵1100... for the case $x < 1$ and 011101110000001...111001 1 1⁵⁵0011... for $x > 1$).

Using these properties, we rather quickly obtained the worst cases for the radix-2 logarithm of all possible double precision input values: it sufficed to run our algorithm for the input numbers of exponents 0, 1, 2, 4, 8, 16, ... 512.

These results, given in Table 7, make it possible to deduce the following property.

Property 4 (Computation of radix-2 logarithms) *Let y be the radix-2 logarithm $\log_2(x)$ of a double-precision number x . Let y^* be an approximation to y such that the mantissa distance between y and y^* is bounded by ϵ . If $\epsilon \leq 2^{-53-55-1} = 2^{-109}$ then for any of the 4 rounding modes, rounding y^* is equivalent to rounding y .*

6 Results: Trigonometric Functions

The results given in Tables 8 to 13 give the worst cases for functions sin, arcsin, cos, arccos, tan and arctan. For these functions, we have worst cases in some bounded domain only, because they are more difficult to handle than the other functions. And yet, it is sometimes possible to prune the search. Consider the arc-tangent of large values. The double precision number that is closest to $\pi/2$ is

$$\alpha = 884279719003555/2^{49}.$$

¹¹1 is replaced by 0 and 0 is replaced by 1.

Table 14: Worst cases for the hyperbolic cosine function in the range $[1/64, 32]$.

Interval	Worst Case (binary)
$[\frac{1}{64}, \frac{1}{2}]$	$\cosh(1.0001011111011000101010011111001000000110001000010111 \times 2^{-6})$ = 1.0000000000001001100011110101111100001001101111100011 1 1 ⁵⁴ 0010...
	$\cosh(1.1011111100000011000001011110001011000110110000110111 \times 2^{-3})$ = 1.0000011000011111010011000011100111100001011011110010 0 0 ⁵⁴ 1101...
$[\frac{1}{2}, 1]$	$\cosh(1.0000001110010010001111110010101101000111110000000111 \times 2^{-1})$ = 1.0010000110011100000110011000100111100011001101110010 1 0 ⁵⁴ 1011...
	$\cosh(1.101001100000001100011100110101011111001001110111010 \times 2^{-1})$ = 1.0101101111111111000001000001101100100110000011111101 1 1 ⁵⁴ 0001...
[1, 2]	$\cosh(1.0001000001001011011001001000111100010001001110100001)$ = 1.1001111011111101110010100110001010110111000000001001 1 1 ⁵⁵ 0100...
[2, 32]	$\cosh(1.1110101001011111001011110010111001001011000011000101 \times 2^1)$ = 10111.000100001101101100001100110100001111111011010101 1 0 ⁵⁷ 1110...

Table 15: Worst cases for the hyperbolic arccosine function in the range $[\cosh(1/64), \cosh(32)] \approx [1.000122, 3.948148 \times 10^{13}]$.

Interval	Worst Case (binary)
$[\cosh(1/64), \cosh(2)]$	$\cosh^{-1} 1.1001111011111101110010100110001010110111000000001010$ = 1.0001000001001011011001001000111100010001001110100001 0 0 ⁵⁵ 1001...
$[\cosh(2), \cosh(32)]$	$\cosh^{-1} 1.0010100101111101111000110101110100000010111010010000 \times 2^{13}$ = 1001.1101101010110001011010010011001010001101001011101 0 1 ⁶¹ 0001...

References

- [1] R. C. Agarwal, J. C. Cooley, F. G. Gustavson, J. B. Shearer, G. Sliselman, and B. Tuckerman. New scalar and vector elementary functions for the IBM system/370. *IBM J. of Research and Development*, 30(2):126–144, March 1986.
- [2] American National Standards Institute and Institute of Electrical and Electronic Engineers. IEEE standard for binary floating-point arithmetic. *ANSI/IEEE Standard, Std 754-1985*, New York, 1985.
- [3] C. S. Iordache and D. W. Matula, *Infinitely Precise Rounding for Division, Square root, and Square Root Reciprocal*, Proc. 14th IEEE Symp. on Computer Arithmetic, 1999, pp. 233-240.
- [4] C. B. Dunham. Feasibility of “perfect” function evaluation. *SIGNAL Newsletter*, 25(4):25–26, October 1990.
- [5] S. Gal and B. Bachelis. An accurate elementary mathematical library for the IEEE floating point standard. *ACM Trans. on Math. Software*, 17(1):26–45, March 1991.
- [6] D. Goldberg. What every computer scientist should know about floating-point arithmetic. *ACM Computing Surveys*, 23(1):5–47, March 1991.
- [7] W. Kahan. Minimizing q^*m-n , text accessible electronically at <http://http.cs.berkeley.edu/~wkahan/>. At the beginning of the file “nearpi.c”, 1983.
- [8] T. Lang and J.M. Muller. Bound on run of zeros and ones for algebraic functions. In Burgess and Ciminiera, editors, *Proc. of the 15th IEEE Symposium on Computer Arithmetic (Arith-15)*. IEEE Computer Society Press, 2001.
- [9] V. Lefèvre. *Developments in Reliable Computing*, chapter An Algorithm That Computes a Lower Bound on the Distance Between a Segment and \mathbb{Z}^2 , pages 203–212. Kluwer, Dordrecht, Netherlands, 1999.

- [10] V. Lefèvre. *Moyens Arithmétiques Pour un Calcul Fiable*. PhD thesis, École Normale Supérieure de Lyon, Lyon, France, 2000.
- [11] V. Lefèvre, J.M. Muller, and A. Tisserand. Toward correctly rounded transcendentals. *IEEE Trans. Computers*, 47(11):1235–1243, November 1998.
- [12] D. Stehlé, V. Lefèvre and P. Zimmermann. Worst Cases and Lattice Reduction. In Bajard and Schulte, editors, *Proceedings of the 16th IEEE Symposium on Computer Arithmetic (Arith 16)*. IEEE Computer Society Press, 2003.
- [13] J.M. Muller. *Elementary Functions, Algorithms and Implementation*. Birkhauser, Boston, 1997.
- [14] M. J. Schulte and E. E. Swartzlander. Hardware designs for exactly rounded elementary functions. *IEEE Trans. Computers*, 43(8):964–973, August 1994.
- [15] A. Ziv. Fast evaluation of elementary mathematical functions with correctly rounded last bit. *ACM Trans. on Math. Software*, 17(3):410–423, September 1991.