

# An Implementation of Sin and Cos Using Gal's Accurate Tables

Pascal Leroy (phl)

2025-09-26

This document describes the implementation of functions `Sin` and `Cos` in Principia. The goals of that implementation are to be portable (including to machines that do not have a fused multiply-add instruction), achieve good performance, and ensure correct rounding.

## Overview

The implementation follows the ideas described by [GB91] and uses accurate tables produced by the method presented in [SZ05]. It guarantees correct rounding with a high probability. In circumstances where it cannot guarantee correct rounding, it falls back to the (slower but correct) implementation provided by the CORE-MATH project [SZG22] [ZSG+24]. More precisely, the algorithm proceeds through the following steps:

- perform argument reduction using Cody and Waite's algorithm in double precision (see [Mul+10, p. 379]);
- if argument reduction loses too many bits (*i.e.*, the argument is close to a multiple of  $\frac{\pi}{2}$ ), fall back to `cr_sin` or `cr_cos`;
- otherwise, uses accurate tables and polynomial approximations to compute `Sin` or `Cos` with extra accuracy;
- if the result fails Muller's rounding test ([Mul+10, pp. 397–400]), fall back to `cr_sin` or `cr_cos`;
- otherwise return the rounded result of the preceding computation.

## Notation and Accuracy Model

In this document we assume a base-2 floating-point number system with  $M$  significand bits<sup>1</sup> similar to the IEEE formats. We define a real function  $m$  and an integer function  $e$  denoting the *significand* and *exponent* of a real number, respectively:

$$x = \pm m(x) \times 2^{e(x)} \quad \text{with} \quad 2^{M-1} \leq m(x) \leq 2^M - 1$$

Note that this representation is unique. Furthermore, if  $x$  is a floating-point number,  $m(x)$  is an integer.

The *unit of the last place*<sup>2</sup> of  $x$  is defined as:

$$u(x) := 2^{e(x)}$$

In particular,  $u(1) = 2^{1-M}$  and:

$$\frac{|x|}{2^M} < \frac{|x|}{2^M - 1} \leq u(x) \leq \frac{|x|}{2^{M-1}} \quad (1)$$

We ignore the exponent bias, overflow and underflow as they play no role in this discussion.

---

<sup>1</sup>In binary64,  $M = 53$ .

<sup>2</sup>Intuitively, this is the ULP “above” for powers of 2.

Finally, for error analysis we use the accuracy model of [Higo2], equation (2.4): everywhere they appear, the quantities  $\delta_i$  represent a roundoff factor such that  $|\delta_i| < u = u(\frac{1}{2}) = 2^{-M}$  (see pages 37 and 38). We also use  $\theta_n$  and  $\gamma_n$  with the same meaning as in [Higo2], lemma 3.1.

## Approximation of $\frac{\pi}{2}$

To perform argument reduction, we need to build approximations of  $\frac{\pi}{2}$  with extra accuracy and analyse the circumstances under which they may be used and the errors that they entail on the reduced argument.

Let  $z \geq 0$ . We start by defining the truncation function  $\text{Tr}(\kappa, z)$  which clears the last  $\kappa$  bits of the significand of  $z$ :

$$\text{Tr}(\kappa, z) := \lfloor 2^{-\kappa} m(z) \rfloor 2^\kappa u(z)$$

We have:

$$z - \text{Tr}(\kappa, z) = (2^{-\kappa} m(z) - \lfloor 2^{-\kappa} m(z) \rfloor) 2^\kappa u(z)$$

The definition of the floor function implies that the quantity in parentheses is in  $[0, 1[$  and therefore:

$$0 \leq z - \text{Tr}(\kappa, z) < 2^\kappa u(z)$$

Furthermore if the bits that are being truncated start with exactly  $k$  zeros we have the stricter inequality:

$$2^{\kappa'-1} u(z) \leq z - \text{Tr}(\kappa, z) < 2^{\kappa'} u(z) \quad \text{with} \quad \kappa' = \kappa - k \quad (2)$$

This leads to the following upper bound for the unit of the last place of the truncation error:

$$u(z - \text{Tr}(\kappa, z)) < 2^{\kappa'-M+1} u(z)$$

which can be made more precise by noting that the function  $u$  is always a power of 2:

$$u(z - \text{Tr}(\kappa, z)) = 2^{\kappa'-M} u(z) \quad (3)$$

### Two-Term Approximation

In this scheme we approximate  $\frac{\pi}{2}$  as the sum of two floating-point numbers:

$$\frac{\pi}{2} \simeq C_1 + \delta C_1$$

which are defined as:

$$\begin{cases} C_1 & := \text{Tr}(\kappa_1, \frac{\pi}{2}) \\ \delta C_1 & := \left\lfloor \frac{\pi}{2} - C_1 \right\rfloor \end{cases}$$

Equation (2) applied to the definition of  $C_1$  yields:

$$2^{\kappa'_1-1} u\left(\frac{\pi}{2}\right) \leq \frac{\pi}{2} - C_1 < 2^{\kappa'_1} u\left(\frac{\pi}{2}\right)$$

where  $\kappa'_1 \leq \kappa_1$  accounts for any leading zeroes in the bits of  $\frac{\pi}{2}$  that are being truncated. Accordingly equation (3) yields, for the unit of the last place:

$$u\left(\frac{\pi}{2} - C_1\right) = 2^{\kappa'_1-M} u\left(\frac{\pi}{2}\right)$$

Noting that the absolute error on the rounding that appears in the definition of  $\delta C_1$  is bounded by  $\frac{1}{2} u\left(\frac{\pi}{2} - C_1\right)$ , we obtain the absolute error on the two-term approximation<sup>3</sup>:

$$\left| \frac{\pi}{2} - C_1 - \delta C_1 \right| \leq \frac{1}{2} u\left(\frac{\pi}{2} - C_1\right) = 2^{\kappa'_1-M-1} u\left(\frac{\pi}{2}\right) \quad (4)$$

<sup>3</sup>With the chosen value  $\kappa_1 = 8$ , the two sides of this inequality turn out to be  $2^{-103.217}$  and  $2^{-101}$ . The difference comes from the fact that  $\left| \frac{\pi}{2} - C_1 - \delta C_1 \right|$  has three zeroes after the last bit of its mantissa, but the theoretical computation above assumes the worst, *i.e.*, a run of ones.

From this we derive the following upper bound for  $\delta C_1$ :

$$\begin{aligned} |\delta C_1| &< \frac{\pi}{2} - C_1 + \frac{1}{2} u\left(\frac{\pi}{2} - C_1\right) \\ &< 2^{\kappa'_1} u\left(\frac{\pi}{2}\right) + 2^{\kappa'_1 - M - 1} u\left(\frac{\pi}{2}\right) = 2^{\kappa'_1} (1 + 2^{-M-1}) u\left(\frac{\pi}{2}\right) \end{aligned} \quad (5)$$

This scheme gives a representation with a significand that has effectively  $2M - \kappa'_1$  bits and is such that multiplying  $C_1$  by an integer less than or equal to  $2^{\kappa'_1}$  is exact.

### Three-Term Approximation

In this scheme we approximate  $\frac{\pi}{2}$  as the sum of three floating-point numbers:

$$\frac{\pi}{2} \simeq C_2 + C'_2 + \delta C_2$$

which are defined as:

$$\begin{cases} C_2 &:= \text{Tr}\left(\kappa_2, \frac{\pi}{2}\right) \\ C'_2 &:= \text{Tr}\left(\kappa'_2, \frac{\pi}{2} - C_2\right) \\ \delta C_2 &:= \left\llbracket \frac{\pi}{2} - C_2 - C'_2 \right\rrbracket \end{cases}$$

Equation (2) applied to the definition of  $C_2$  yields:

$$2^{\kappa'_2 - 1} u\left(\frac{\pi}{2}\right) \leq \frac{\pi}{2} - C_2 < 2^{\kappa'_2} u\left(\frac{\pi}{2}\right) \quad (6)$$

where  $\kappa'_2 \leq \kappa_2$  accounts for any leading zeroes in the bits of  $\frac{\pi}{2}$  that are being truncated. Accordingly equation (3) yields, for the unit of the last place:

$$u\left(\frac{\pi}{2} - C_2\right) = 2^{\kappa'_2 - M} u\left(\frac{\pi}{2}\right)$$

Similarly, equation (2) applied to the definition of  $C'_2$  yields:

$$\begin{aligned} 2^{\kappa''_2 - 1} u\left(\frac{\pi}{2} - C_2\right) &\leq \frac{\pi}{2} - C_2 - C'_2 < 2^{\kappa''_2} u\left(\frac{\pi}{2} - C_2\right) \\ 2^{\kappa'_2 + \kappa''_2 - M - 1} u\left(\frac{\pi}{2}\right) &\leq < 2^{\kappa'_2 + \kappa''_2 - M} u\left(\frac{\pi}{2}\right) \end{aligned}$$

where  $\kappa''_2 \leq \kappa_2$  accounts for any leading zeroes in the bits of  $\frac{\pi}{2} - C_2$  that are being truncated. Note that normalization of the significand of  $\frac{\pi}{2} - C_2$  effectively drops the zeroes at positions  $\kappa_2$  to  $\kappa'_2$  and therefore the computation of  $C'_2$  applies to a significand aligned on position  $\kappa'_2$ .

It is straightforward to transform these inequalities using (6) to obtain bounds on  $C'_2$ :

$$2^{\kappa'_2} \left(\frac{1}{2} - 2^{\kappa''_2 - M}\right) u\left(\frac{\pi}{2}\right) < C'_2 < 2^{\kappa'_2} (1 - 2^{\kappa''_2 - M - 1}) u\left(\frac{\pi}{2}\right)$$

Equation (3) applied to the definition of  $C'_2$  yields, for the unit of the last place:

$$\begin{aligned} u\left(\frac{\pi}{2} - C_2 - C'_2\right) &= 2^{\kappa''_2 - M} u\left(\frac{\pi}{2} - C_2\right) \\ &= 2^{\kappa'_2 + \kappa''_2 - 2M} u\left(\frac{\pi}{2}\right) \end{aligned}$$

Noting that the absolute error on the rounding that appears in the definition of  $\delta C_2$  is bounded by  $\frac{1}{2} u\left(\frac{\pi}{2} - C_2 - C'_2\right)$ , we obtain the absolute error on the three-term approximation:

$$\left| \frac{\pi}{2} - C_2 - C'_2 - \delta C_2 \right| \leq \frac{1}{2} u\left(\frac{\pi}{2} - C_2 - C'_2\right) = 2^{\kappa'_2 + \kappa''_2 - 2M - 1} u\left(\frac{\pi}{2}\right) \quad (7)$$

and the following upper bound for  $\delta C_2$ :

$$|\delta C_2| < 2^{\kappa'_2 + \kappa''_2 - M} (1 + 2^{-M-1}) u\left(\frac{\pi}{2}\right) \quad (8)$$

This scheme gives a representation with a significand that has effectively  $3M - \kappa'_2 - \kappa''_2$  bits and is such that multiplying  $C_2$  and  $C'_2$  by an integer less than or equal to  $2^{\kappa'_2}$  is exact.

## Argument Reduction

Given an argument  $x$ , the purpose of argument reduction is to compute a pair of floating-point numbers  $(\tilde{x}, \delta\tilde{x})$  such that:

$$\begin{cases} \tilde{x} + \delta\tilde{x} \cong x \pmod{\frac{\pi}{2}} \\ \tilde{x} \text{ is approximately in } \left[-\frac{\pi}{4}, \frac{\pi}{4}\right] \\ |\delta\tilde{x}| \leq \frac{1}{2} u(\tilde{x}) \end{cases}$$

### Argument Reduction for Small Angles

If  $|x| < \left\lceil\frac{\pi}{4}\right\rceil$  then  $\tilde{x} = x$  and  $\delta\tilde{x} = 0$ .

### Argument Reduction Using the Two-Term Approximation

If  $|x| \leq 2^{\kappa_1} \left\lceil\frac{\pi}{2}\right\rceil$  we compute:

$$\begin{cases} n &= \left\lceil\left\lceil x \left\lceil\frac{2}{\pi}\right\rceil\right\rceil\right\rceil \\ y &= x - n C_1 \\ \delta y &= \llbracket n \delta C_1 \rrbracket \\ (\tilde{x}, \delta\tilde{x}) &= \text{TwoDifference}(y, \delta y) \end{cases}$$

The first thing to note is that  $|n| \leq 2^{\kappa_1}$ . We have:

$$|x| \leq 2^{\kappa_1} \left\lceil\frac{\pi}{2}\right\rceil = 2^{\kappa_1} \frac{\pi}{2} (1 + \delta_1)$$

and:

$$\left\lceil\left\lceil x \left\lceil\frac{2}{\pi}\right\rceil\right\rceil\right\rceil = x \frac{2}{\pi} (1 + \delta_2)(1 + \delta_3) \quad (9)$$

from which we deduce the upper bound:

$$\begin{aligned} |n| &\leq \left\lceil 2^{\kappa_1} \frac{\pi}{2} (1 + \delta_1) \frac{2}{\pi} (1 + \delta_2)(1 + \delta_3) \right\rceil \\ &\leq \lceil 2^{\kappa_1} (1 + \gamma_3) \rceil \end{aligned}$$

If  $2^{\kappa_1} \gamma_3$  is small enough (less than 1/2), the rounding cannot cause  $n$  to exceed  $2^{\kappa_1}$ . In practice we choose a relatively small value for  $\kappa_1$ , so this condition is met.

Now if  $x$  is close to an odd multiple of  $\frac{\pi}{4}$  it is possible for misrounding to happen. There are two kinds of misrounding, with different bounds.

A misrounding of the first kind occurs if, assuming  $n > 0$ :

$$x < \left(n - \frac{1}{2}\right) \frac{\pi}{2} \quad \text{and} \quad \left\lceil\left\lceil x \left\lceil\frac{2}{\pi}\right\rceil\right\rceil\right\rceil > n - \frac{1}{2}$$

Using equation (9) we find that this misrounding is only possible if:

$$x > \frac{\pi}{2} \left(n - \frac{1}{2}\right) \frac{1}{(1 + \delta_2)(1 + \delta_3)} \geq \frac{\pi}{2} \left(n - \frac{1}{2}\right) \frac{1}{1 + \gamma_2}$$

In which case the computation of  $n$  results in:

$$n \frac{\pi}{2} - x < \frac{\pi}{4} \left(1 + \frac{\gamma_2}{1 + \gamma_2} (2n - 1)\right)$$

In this case, misrounding causes the absolute value of the reduced angle to increase and it may thus exceed  $\frac{\pi}{4}$  by as much as:

$$\frac{\pi}{4} \frac{\gamma_2}{1 + \gamma_2} (2^{\kappa_1+1} - 1)$$

A misrounding of the second kind occurs if, assuming  $n \geq 0$ :

$$x > \left(n + \frac{1}{2}\right) \frac{\pi}{2} \quad \text{and} \quad \left\lceil \left\lceil x \left\lceil \frac{2}{\pi} \right\rceil \right\rceil \right\rceil < n + \frac{1}{2}$$

A derivation similar to the one above gives the following condition for this misrounding to be possible. Using equation (9):

$$x < \frac{\pi}{2} \left(n + \frac{1}{2}\right) \frac{1}{(1 + \delta_2)(1 + \delta_3)} \leq \frac{\pi}{2} \left(n + \frac{1}{2}\right) (1 + \gamma_2)$$

we derive the bound:

$$x - n \frac{\pi}{2} < \frac{\pi}{4} (1 + \gamma_2 (2n + 1))$$

In this case, misrounding causes the absolute value of the reduced angle to decrease by as much as:

$$\frac{\pi}{4} \gamma_2 (2^{\kappa_1+1} + 1) \quad (10)$$

This is however not a concern for the accurate tables as it cannot cause the reduced angle to become negative.

Using the bound on  $|n|$  and the fact that  $C_1$  has  $\kappa_1$  trailing zeroes, we see that the product  $n C_1$  is exact. The subtraction  $x - n C_1$  is exact by Sterbenz's Lemma. Finally, the last step performs an exact addition<sup>4</sup> using algorithm 4 of [HLBo8].

To compute the overall error on argument reduction<sup>5</sup>, first remember that, from equation (4), we have:

$$C_1 + \delta C_1 = \frac{\pi}{2} + \zeta \quad \text{with} \quad |\zeta| \leq 2^{\kappa'_1 - M - 1} u\left(\frac{\pi}{2}\right)$$

The error computation proceeds as follows:

$$\begin{aligned} y - \delta y &= x - n C_1 - n \delta C_1 (1 + \delta_4) \\ &= x - n(C_1 + \delta C_1) - n \delta C_1 \delta_4 \\ &= x - n \frac{\pi}{2} - n(\zeta + \delta C_1 \delta_4) \end{aligned}$$

from which we deduce an upper bound on the absolute error of the reduction:

$$\begin{aligned} \left| y - \delta y - \left(x - n \frac{\pi}{2}\right) \right| &\leq 2^{\kappa_1} 2^{\kappa'_1} (2^{-M-1} + 2^{-M} + 2^{-2M-1}) u\left(\frac{\pi}{2}\right) \\ &= 2^{\kappa_1 + \kappa'_1 - M} \left(\frac{3}{2} + 2^{-M-1}\right) u\left(\frac{\pi}{2}\right) \\ &< 2^{\kappa_1 + \kappa'_1 - M + 1} u\left(\frac{\pi}{2}\right) \end{aligned}$$

where we have used the upper bound for  $\delta C_1$  given by equation (5).

The exact TwoDifference yields a pair such that  $|\delta \tilde{x}| \leq \frac{u(\tilde{x})}{2} \leq 2^{-M} |\tilde{x}|$ . Furthermore, misrounding of the first kind and the above error on the reduction may combine to cause  $|\tilde{x}|$  to move above  $\frac{\pi}{4}$  by as much as:

$$\frac{\pi}{4} \frac{\gamma_2}{1 + \gamma_2} (2^{\kappa_1+1} - 1) + 2^{\kappa_1 + \kappa'_1 - M + 1} u\left(\frac{\pi}{2}\right)$$

---

<sup>4</sup>The more efficient QuickTwoDifference is not usable here. First, note that  $|y|$  is equal to  $u(x)$  if we take  $x$  to be the successor or the predecessor of  $n C_1$  for any  $n$ . Ignoring rounding errors we have:

$$|\delta y| \geq n 2^{\kappa'_1 - 1} u\left(\frac{\pi}{2}\right) \geq 2^{\kappa'_1 + M - 2} u\left(\frac{\pi}{2}\right) u(n)$$

where we used the bound given by equation (1). Now the computation of  $n$  can result in a value that is either in the same binade or in the binade below that of  $x$ . Therefore  $u(n) \geq \frac{1}{2} u(x)$  and the above inequality becomes:

$$|\delta y| \geq 2^{\kappa'_1 + M - 3} u\left(\frac{\pi}{2}\right) u(x)$$

plugging  $u\left(\frac{\pi}{2}\right) = 2^{1-M}$  we find:

$$|\delta y| \geq 2^{\kappa'_1 - 2} u(x)$$

Therefore, as long as  $\kappa'_1 > 2$ , there exist arguments  $x$  for which  $|\delta y| > |y|$ .

<sup>5</sup>Note that this error analysis is correct even in the face of misrounding.

The accurate tables must be constructed so that the last interval covers angles misrounded in that manner<sup>6</sup>.

In the computation of the trigonometric functions, we need  $\tilde{x} + \delta\tilde{x}$  to provide enough accuracy that the final result is correctly rounded most of the time. The above error bound shows that, if  $\tilde{x}$  is very small (*i.e.*, if  $x$  is very close to a multiple of  $\frac{\pi}{2}$ ), the two-term approximation may not provide enough correct bits. Formally, say that we want to have  $M + \kappa_3$  correct bits in the mantissa of  $\tilde{x} + \delta\tilde{x}$ . The error must be less than  $2^{-\kappa_3}$  half-units of the last place of the result:

$$2^{\kappa_1 + \kappa'_1 - M + 1} u\left(\frac{\pi}{2}\right) \leq 2^{-\kappa_3 - 1} u(\tilde{x}) \leq 2^{-\kappa_3 - M} |\tilde{x}| \quad (11)$$

which leads to the following condition on the reduced angle:

$$|\tilde{x}| \geq 2^{\kappa_1 + \kappa'_1 + \kappa_3 + 1} u\left(\frac{\pi}{2}\right) = 2^{\kappa_1 + \kappa'_1 + \kappa_3 - M + 2}$$

The rest of the implementation assumes that  $\kappa_3 = 18$  to achieve correct rounding with high probability. If we choose  $\kappa_1 = 8$  we find that  $\kappa'_1 = 5$  (because there are three consecutive zeroes at this location in the significand of  $\frac{\pi}{2}$ ) and the desired accuracy is obtained as long as  $|\tilde{x}| \geq 2^{-20} \simeq 9.5 \times 10^{-7}$ .

### Argument Reduction Using the Three-Term Approximation

If  $|x| \leq 2^{\kappa_2} \left\lfloor \frac{\pi}{2} \right\rfloor$  we compute:

$$\begin{cases} n &= \left\lfloor x \left\lfloor \frac{2}{\pi} \right\rfloor \right\rfloor \\ y &= x - n C_2 \\ y' &= n C'_2 \\ \delta y &= \llbracket n \delta C_2 \rrbracket \\ (z, \delta z) &= \text{QuickTwoSum}(y', \delta y) \\ (\tilde{x}, \delta\tilde{x}) &= \text{DWPlusFP}((-z, -\delta z), y) \end{cases}$$

The products  $n C_2$  and  $n C'_2$  are exact thanks to the  $\kappa_2$  trailing zeroes of  $C_2$  and  $C'_2$ . The subtraction  $x - n C_2$  is exact by Sterbenz's Lemma. QuickTwoSum performs an exact addition using algorithm 3 of [HLBo8]; it is usable in this case because clearly  $|\delta y| < |y'|$ . DWPlusFP is algorithm 5 from [MR22], which implements precise (but not exact) double-precision arithmetic.

It is straightforward to show, like we did in the preceding section, that:

$$|n| \leq \lceil 2^{\kappa_2} (1 + \gamma_3) \rceil$$

and therefore that  $|n| \leq 2^{\kappa_2}$  as long as  $2^{\kappa_2} \gamma_3 < 1/2$ . Similarly, the misrounding bound (10) is applicable with  $\kappa_2$  replacing  $\kappa_1$ .

To compute the overall error on argument reduction, first remember that, from equation (7), we have:

$$C_2 + C'_2 + \delta C_2 = \frac{\pi}{2} + \zeta_1 \quad \text{with} \quad |\zeta_1| \leq 2^{\kappa'_2 + \kappa''_2 - 2M - 1} u\left(\frac{\pi}{2}\right)$$

Let  $\zeta_2$  be the relative error introduced by DWPlusFP. Theorem 2.2 of [MR22] indicates that  $|\zeta_2| < 2^{1-2M}$ . The error computation proceeds as follows:

$$\begin{aligned} y - y' - \delta y &= (x - n C_2 - n C'_2 - n \delta C_2 (1 + \delta_4)) (1 + \zeta_2) \\ &= \left( x - n \frac{\pi}{2} - n (\zeta_1 + \delta C_2 \delta_4) \right) (1 + \zeta_2) \\ &= x - n \frac{\pi}{2} - n (\zeta_1 + \delta C_2 \delta_4) (1 + \zeta_2) + \left( x - n \frac{\pi}{2} \right) \zeta_2 \end{aligned}$$

<sup>6</sup>In practice this is not a stringent constraint because the distance between accurate table entries is much larger than this quantity.

from which we deduce an upper bound on the absolute error of the reduction, noting that  $\left|x - n\frac{\pi}{2}\right| \leq \frac{\pi}{4}(1 + \gamma_2(2^{\kappa_2+1} + 1))$  as per (10):

NONONO

$$\begin{aligned}
& \left|y - y' - \delta y - \left(x - n\frac{\pi}{2}\right)\right| \\
& \leq 2^{\kappa_2+\kappa'_2+\kappa''_2}(2^{-2M-1} + 2^{-2M} + 2^{-3M-1})(1 + 2^{1-2M})u\left(\frac{\pi}{2}\right) + 2^{1-2M}\frac{\pi}{4}(1 + \gamma_2(2^{\kappa_2+1} + 1)) \\
& < 2^{\kappa_2+\kappa'_2+\kappa''_2-2M}\left(\frac{3}{2} + 2^{-M-1}\right)(1 + 2^{1-2M})u\left(\frac{\pi}{2}\right) + 2^{-2M}\frac{\pi}{2}\left(1 + 3 \times 2^{\kappa_2}u\left(\frac{\pi}{2}\right)\right) \\
& < 2^{\kappa_2-2M}(2^{\kappa'_2+\kappa''_2+1} + 5)u\left(\frac{\pi}{2}\right) + 2^{-2M}\frac{\pi}{2}
\end{aligned}$$

where the second inequality uses  $\gamma_2(2^{\kappa_2+1} + 1) < 3u\ 2^{\kappa_2+1}$ .

A sufficient condition for the reduction to guarantee  $\kappa_3$  extra bits of accuracy is for this error to be less than  $2^{-\kappa_3-1}u(\tilde{x})$  which itself is less than  $2^{-\kappa_3-M}|\tilde{x}|$ . Therefore we want:

$$\begin{aligned}
|\tilde{x}| & \geq 2^{\kappa_3-M}\left(2^{\kappa_2}(2^{\kappa'_2+\kappa''_2+1} + 5)u\left(\frac{\pi}{2}\right) + \frac{\pi}{2}\right) \\
& = 2^{\kappa_3-M}\left(2^{\kappa_2-M+1}(2^{\kappa'_2+\kappa''_2+1} + 5) + \frac{\pi}{2}\right)
\end{aligned}$$

and it is therefore sufficient to have:

$$|\tilde{x}| \geq 2^{\kappa_3-M}(2^{\kappa_2+\kappa'_2+\kappa''_2-M+2} + 2)$$

If we choose  $\kappa_3 = 18$  as above, and  $\kappa_2 = 18$  we find that  $\kappa'_2 = 14$  and  $\kappa''_2 = 15$ . Therefore, the desired accuracy is obtained as long as  $|\tilde{x}| \geq 33 \times 2^{-39} \simeq 6.0 \times 10^{-11}$ .

## Fallback

If any of the conditions above is not met, we fall back on the CORE-MATH implementation.

## Accurate Tables and Their Generation

First, a remark on the notation. This section (as well as the code in `accurate_table_generator*.hpp`) follows the notation of [SZ05]. In particular,  $N$  is related to the number of bits in the mantissa of the floating-point type, and  $M$  is related to the number of zeroes or ones after the mantissa of the accurate values. Specifically, since we use `binary64` and want 18 zeroes or ones after the mantissa (to minimize the probability of expensive fall backs to the CORE-MATH implementation), we have  $N = 2^{53}$  and  $M = 2^{18}$ .

The accurate tables are made of triples  $(x_k, s_k, c_k)$  where  $x_k$  is chosen such that  $\sin x_k$  and  $\cos x_k$  are very close to machine numbers, *i.e.*, have a string of zeroes or ones after the last bit of their mantissa. The tabulated value  $s_k$  (respectively,  $c_k$ ) is then obtained by rounding  $\sin x_k$  (respectively,  $\cos x_k$ ) to the nearest machine number. See Overview of the Algorithms for Sin and Cos below for the details of how these tables are used for implementing sin and cos.

The argument range (which is approximately  $\left[0, \frac{\pi}{4}\right]$ ) is split into intervals of length  $2\Delta$  centered on  $2k\Delta$ . The interval<sup>7</sup> containing  $x_k$  is  $\mathcal{J}_k := [(2k-1)\Delta, (2k+1)\Delta]$  (for  $k = 0$  the interval is  $\mathcal{J}_0 := [0, \Delta]$  and  $x_0 = 0$ ). We want to select  $x_k$  as close as possible to the midpoint of  $\mathcal{J}_k$ , because we will need to approximate sin and cos over the interval  $[-h_{\max}, h_{\max}]$  with:

$$h_{\max} := \max_k (x_k - (2k-1)\Delta, (2k+1)\Delta - x_k)$$

<sup>7</sup>Obviously the intervals  $\mathcal{J}_k$  cannot all contain their bounds. Because of the way they are computed, the odd multiples of  $\Delta$  which separate the intervals are rounded to the nearest even  $k$ . Therefore, there is an alternation of open intervals (for  $k$  odd) and closed intervals (for  $k$  even). We do not have a convenient notation for this, and anyway this detail is mostly irrelevant.

and we want this interval to be as small as possible to improve the accuracy of the approximation.

In practice we choose  $\Delta = 2^{-10}$  and the accurate tables construction yields  $h_{\max} < \Delta + 2^{-17.834}$ . The largest perturbation for  $x_k$  is reached when  $k = 1$ , with 35.629 bits, but that is because, for  $k = 1$ , we only search below  $2\Delta$  (see Sin and Cos Around Table Entries). Other than this special value, the perturbation ranges from 23.755 bits to 35.331 bits. This is consistent with the random model of [SZo5] but note that some intervals are 3000 times “harder” than others.

The naïve table construction method described in [Gal86] (random sampling over  $J_k$ ) has two problems. First, it doesn’t guarantee that  $x_k$  is close to the midpoint of  $J_k$  and can in fact, in the worst case, result in  $x_k$  being very close to one of the bounds of the interval, requiring the construction of a minimax polynomial over an interval which is (nearly) two times too big. Second, for each triple of the accurate table it has to inspect, on average, around  $M^2 = 2^{36}$  values. On a modern processor, testing one candidate takes approximately 50  $\mu$ s, which amounts to about 900 hours for each triple of the table, and therefore 43 years for the 402 triples that we need.

Consequently, we use the method described in [SZo5] to build our accurate tables. It is much faster than the method in [Gal86], but must still be implemented with great care to achieve a satisfactory speed. In the rest of this section, we describe the algorithms and optimizations that we used to generate our tables efficiently. On a modern processor (AMD Ryzen Threadripper PRO 5965WX) the entire table generation takes about 17 seconds.

## The Stehlé-Zimmermann Algorithm

For convenience, we copy algorithm 1 from [SZo5].

---

### Algorithm 1: SimultaneousBadCaseSearch.

---

**Input:** Two functions  $F_1$  and  $F_2$ , and two positive integers  $M, T$ .

**Output:** All  $t \in [-T, T]$  such that  $|F_i(t) \bmod 1| < \frac{1}{M}$  for  $i \in \{1, 2\}$ .

1. Let  $P_1(t), P_2(t)$  be the degree-2 Taylor expansions of  $F_1(t), F_2(t)$ .
  2. Compute  $\epsilon$  such that  $|P_i(t) - F_i(t)| < \epsilon$  for  $|t| \leq T$  and  $i \in \{1, 2\}$ .
  3. Compute  $M' = \left\lfloor \frac{1/2}{1/M + \epsilon} \right\rfloor$ ,  $C = 3M'$ , and<sup>\*</sup>  $\tilde{P}_i(\tau) = [C \cdot P_i(T\tau)]$  for  $i \in \{1, 2\}$ .
  4. Let  $e_1 = 1, e_2 = \tau, e_3 = \tau^2, e_4 = v, e_5 = \phi$ .
  5. Let  $g_1 = C, g_2 = CT \cdot \tau, g_3 = \tilde{P}_1(\tau) + 3v, g_4 = \tilde{P}_2(\tau) + 3\phi$ .
  6. Create the  $4 \times 5$  integral matrix  $L$  where  $L_{k,l}$  is the coefficient of the monomial  $e_l$  in  $g_k$ .
  7.  $V \leftarrow \text{LatticeReduce}(L)$ .
  8. Let  $\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3$  be the three shortest vectors of  $V$ , and  $Q_i(\tau)$  the associated polynomials.
  9. **if** there exist  $i \in \{1, 2, 3\}$  such that  $\|\mathbf{v}_i\|_1 \geq C$  **then return FAIL**.
  10. Let  $Q(\tau)$  be a linear combination of the  $Q_i$ ’s which is independent of  $v$  and  $\phi$ . We have  $\deg Q \leq 1$ .
  11. Let  $q(t) = Q\left(\frac{t}{T}\right)$ .  
**foreach**  $t_0$  **in**  $\text{IntegerRoots}(q, [-T, T])$  **do** **if**  $|F_i(t_0) \bmod 1| < \frac{1}{M}$  **for all**  $i \in \{1, 2\}$  **then return**  $t_0$ .
- 

When algorithm 1 returns a value at step 11, that value can be used to build an accurate table entry. When it does not return a value at step 11, the interval  $[-T, T]$  is guaranteed not to contain a point suitable for building an accurate table entry. When it fails at step 9, the algorithm must be retried over  $[-T/2, T/2]$  until a conclusive answer is found at step 11. Note that, despite the phrase “All  $t$ ” in the description of the result, algorithm 1 returns at most one value (because  $\deg Q \leq 1$ ).

---

<sup>\*</sup>The notation  $[C \cdot P_i(T\tau)]$  means that we round to the nearest each coefficient of  $C \cdot P_i(T\tau)$ . This gives an element of  $\mathbb{Z}[\tau]$ . (This footnote is missing from [SZo5] but present in the preprint [SZo4].)



## The Complete Search

In order to construct the accurate tables we need to build on algorithm 1 to perform a search over a arbitrary interval. In this context, if  $T$  is small enough, we are better off doing an exhaustive search than going through the fairly expensive lattice reduction. For an exhaustive search, we let the caller specify whether the search should proceed by increasing or decreasing values to make sure that we find, to the extent possible, the solution closest to the centre of  $\mathcal{J}_k$ . Algorithm 2 spells out the processing of one interval specifically for the sin and cos functions.

---

**Algorithm 2:** IntervalSimultaneousBadCaseSearch.

---

**Input:** The interval to search  $\left[x_{\text{mid}} - \frac{T}{N}, x_{\text{mid}} + \frac{T}{N}\right]$  defined by its midpoint  $x_{\text{mid}}$  and its scaled radius  $T$ , and a positive integer  $M$ .  
**Output:** A  $t \in [-T, T]$  such that  $x_{\text{mid}} + \frac{t}{N}$  is an accurate value for sin and cos; furthermore, if the search direction is UP,  $t$  is as close as possible to  $-T$  and the search direction is DOWN,  $t$  is as close as possible to  $T$ .

1. Let  $F_1 = t \mapsto N \sin\left(x_{\text{mid}} + \frac{t}{N}\right)$  and  $F_2 = t \mapsto N \cos\left(x_{\text{mid}} + \frac{t}{N}\right)$ .
2. **if**  $T < \text{some threshold}$  **then**
3.     **switch** *search direction* **do**
4.         **case** *UP* **do**
5.             **foreach**  $t_0 \in [-T, T]$  *by ascending values* **do**
6.                 **if**  $|F_i(t_0) \bmod 1| < \frac{1}{M}$  *for all*  $i \in \{1, 2\}$  **then**
7.                     **return**  $t_0$ .
8.                 **end**
9.             **end**
10.         **case** *DOWN* **do**
11.             **foreach**  $t_0 \in [-T, T]$  *by descending values* **do**
12.                 **if**  $|F_i(t_0) \bmod 1| < \frac{1}{M}$  *for all*  $i \in \{1, 2\}$  **then**
13.                     **return**  $t_0$ .
14.                 **end**
15.             **end**
16.     **end**
17. **else**
18.     **return** SimultaneousBadCaseSearch( $F_1, F_2, M, T$ ).
19. **end**

---

Algorithm 2 must be integrated into a retry loop that searches progressively farther above and below  $2k\Delta$  until a solution is found. This is what algorithm 3 is doing: it first searches an interval of radius  $T_0$  above  $2k\Delta$ , and retries with intervals progressively smaller if algorithm 2 returns a failure. When algorithm 2 returns an empty solution set, the last interval searched is marked as “covered” and the algorithm proceeds to perform a similar search below  $2k\Delta$ . If that search does not find a solution, a new interval of radius  $T_0$  above the covered interval is tried, and so on repeatedly. The algorithm returns when a solution is found. That solution is nearly optimal in terms of its distance to  $2k\Delta$  with the caveat that, since we search above  $2k\Delta$  before searching below, we could return a solution above when there is a slightly closer solution below.

The value of  $T_0$  is chosen based on the Coppersmith bound established in [SZ05]:  $T^3 = \mathcal{O}(M \cdot N)$ . Evidently this bound does not unambiguously determine  $T_0$ , but in practice  $T_0 = \sqrt[3]{MN}$  yields good results: a value too large would cause too many splits of the search interval before finding a solution or (more likely) rejecting the interval; a value too small would require to process too many small intervals.

## Representation of Functions and Polynomials

Our implementation is based on the Boost multiprecision library. The functions sin and cos take `cpp_rational` values as arguments and return `cpp_bin_float_50`

---

**Algorithm 3:** FullSimultaneousBadCaseSearch.

---

**Input:**  $\Delta$ , the half-size of intervals, and  $k$ , the index of the interval to search.

**Output:** An accurate value sufficiently close to  $2k\Delta$ .

```

1.  $T_0 \leftarrow \sqrt[3]{MN}$ 
2.  $\mathcal{J}_{\text{covered}} \leftarrow [2k\Delta, 2k\Delta]$ 
3. loop
4.    $T \leftarrow T_0$ 
5.   loop
6.     Let  $\mathcal{J}_{\text{to\_cover}}$  be the interval of radius  $T$  immediately above  $\mathcal{J}_{\text{covered}}$ .
7.      $t \leftarrow \text{IntervalSimultaneousBadCaseSearch}(\mathcal{J}_{\text{to\_cover}}, UP, M)$ .
8.     if  $t$  is a solution then
9.       return  $\text{Midpoint}(\mathcal{J}_{\text{to\_cover}}) + \frac{t}{N}$ .
10.    else if  $t = FAIL$  then
11.       $T \leftarrow T/2$ .
12.    else
13.       $\mathcal{J}_{\text{covered}} = \mathcal{J}_{\text{covered}} \cup \mathcal{J}_{\text{to\_cover}}$ 
14.      break
15.    end
16.  end
17.   $T \leftarrow T_0$ 
18.  loop
19.    Let  $\mathcal{J}_{\text{to\_cover}}$  be the interval of radius  $T$  immediately below  $\mathcal{J}_{\text{covered}}$ .
20.     $t \leftarrow \text{IntervalSimultaneousBadCaseSearch}(\mathcal{J}_{\text{to\_cover}}, DOWN, M)$ .
21.    if  $t$  is a solution then
22.      return  $\text{Midpoint}(\mathcal{J}_{\text{to\_cover}}) + \frac{t}{N}$ .
23.    else if  $t = FAIL$  then
24.       $T \leftarrow T/2$ .
25.    else
26.       $\mathcal{J}_{\text{covered}} = \mathcal{J}_{\text{covered}} \cup \mathcal{J}_{\text{to\_cover}}$ 
27.      break
28.    end
29.  end
30. end

```

---

values (50 decimal digits should be plenty for our purposes since we are looking for values that are accurate to approximately 21 digits; at any rate, while the construction of the accurate values is expensive, checking that they are correct is straightforward so we are not concerned about accuracy issues in this context). The polynomial approximations take `cpp_rational` values as arguments, have coefficients that are `cpp_rational` and return `cpp_rational`.

Step 2 of algorithm 1 requires that we compute an upper bound for  $|P_i(t) - F_i(t)|$  over an interval. To do this efficiently (and in particular, without doing multiple evaluation of the functions and the polynomials), we approximate  $P_i(t) - F_i(t)$  by the degree-3 term of the Taylor expansion of  $F_i(t)$  and find the extremum of this term. While ignoring the terms of degree 4 and higher introduces a small error, this error is of no importance as  $\epsilon$  is typically much smaller than  $1/M$  and does not affect the value of  $M'$  (see step 3).

It would be possible (and correct) to compute the polynomials  $P_i$  using a Taylor approximation around  $2k\Delta$  and to translate them just like we translate the functions in algorithm 2:  $t \mapsto NP_i(x_{\text{mid}} + \frac{t}{N})$ . In practice however this results in polynomial coefficients with very large numerators and denominators, greatly affecting the performance of algorithm 1. It is much more efficient to recompute the Taylor expansion around  $x_{\text{mid}}$  as this produces coefficients with smaller numerators and denominators for the same accuracy. (The coefficients of the Taylor expansion involve evaluations of  $\sin$  and  $\cos$ , so the two methods do not generally result in the exact same polynomial.)

## Lattice Reduction

Even though the matrix  $L$  is small, a fair amount of time is spent in the lattice reduction at step 7 of algorithm 1. The classical algorithm for lattice reduction is the Lenstra-Lenstra-Lovász algorithm, described for instance in [HPS14, p. 444] but it turns out to be fairly costly in our case because the matrix elements are of type `cpp_int` and operations on this type are expensive. Instead we use the algorithm described in [NS09], which is quadratic and speeds up the search significantly.

## Speculative Execution

The natural way to parallelize the search on a multicore processor is to run multiple instance of algorithm 3 for distinct intervals  $\mathcal{I}_k$  each in its own thread. Unfortunately this does not work very well because the searches have a “long tail”: for most intervals, a solution is found relatively quickly, but for some intervals the search has to explore values far from  $2k\Delta$ . So after a short time only a few “expensive” intervals are left to process, effectively doing sequential searches while most of the cores are idle. This does not properly take advantage of the available parallelism.

It is much more effective to parallelize using speculative execution as follows. Each of the intervals of radius  $T_0$  processed in steps 5 and 17 of algorithm 3 is a *slice* and slices can be searched independently for a given value of  $k$ . As soon as a core is available, it picks a slice to process, which is guaranteed not to have been processed already, but which may not be adjacent to the interval  $\mathcal{I}_{\text{covered}}$  (it is centered at  $2k\Delta + (2j + 1)\frac{T_0}{N}$  where  $j$  is the *slice index*). The execution is speculative in the sense that the search for slice index  $j$  may not be necessary if a solution exists, say, for index  $j - 1$ . However, most slices do not contain a solution, so the amount of wasted work is minimal. Note that, for speculative execution to be deterministic, it is important, if a solution is found in slice  $j$ , to wait for the searches with indices preceding  $j$  to complete before returning, as one of these slices may contain a solution closer to  $2k\Delta$ .

## Summary

The following table summarizes the impact of each of the optimizations described above. The “Baseline” row is the time with all optimizations enabled. Each of the other rows gives the timing with exactly one optimization disabled:

Method	Elapsed Time	Ratio
Baseline	16.6 s	
No speculative execution	107.2 s	6.5×
Reduction using Lenstra-Lenstra-Lovász	237.1 s	14.3×
Polynomials centred on $2k\Delta$	around 363 000 s	22 000×

## Overview of the Algorithms for Sin and Cos

This section presents an overview of the algorithms used to compute sin and cos. They take as input the result of argument reduction,  $(\tilde{x}, \delta\tilde{x})$  and produce a pair  $(y, \delta y)$  which is passed to the rounding test described in [Mul+10, p. 397] to decide whether  $\llbracket y + \delta y \rrbracket$  is the correctly-rounded result (this is the case with high probability) or whether we need to fall back to the CORE-MATH implementation.

The lower part of the reduced angle,  $\delta\tilde{x}$ , must be used with care: on the one hand, it would not make sense to explicitly compute expressions like  $\llbracket \tilde{x} + \delta\tilde{x} \rrbracket$  as the second term vanishes before the first one; on the other hand, we must avoid the computation of expressions involving  $\delta\tilde{x}$  that have no bearing on the final result (see [Mul+10, pp. 402–404] for a discussion of this issue). We are going to spell out the places where  $\delta\tilde{x}$  is actually used and we will prove later that neglecting it in other places has no effect.

### Sin Near Zero

For the sin function near zero the accurate tables method is not usable because the correction term is not small compared to the tabulated value of the function (which can be arbitrarily close to zero). Instead we use a polynomial approximation that minimizes the relative error on the result. Since  $\sin t$  is an odd function and since its dominant term is  $t$ , we are using an approximation of the form:

$$\sin t \simeq t + t^3 p_{s0}(t^2)$$

and we compute:

$$\sin x \simeq \sin(\tilde{x} + \delta\tilde{x}) \simeq \tilde{x} + \tilde{x}^3 p_{s0}(\tilde{x}^2) + \delta\tilde{x}$$

### Sin and Cos Around Table Entries

Let  $(x_k, s_k, c_k)$  be an accurate table entry. The implementation of sin and cos starts by choosing, from the argument  $\tilde{x}$  (which in the case of sin, is not close to zero), a  $k$  such that  $\tilde{x} \in J_k$  and computing  $h = \tilde{x} - x_k$ . We then use the addition formulæ to separate out the terms in  $x_k$ :

$$\begin{aligned}
 \sin x &\simeq \sin(\tilde{x} + \delta\tilde{x}) \\
 &= \sin(x_k + h + \delta\tilde{x}) \\
 &= \sin x_k \cos(h + \delta\tilde{x}) + \cos x_k \sin(h + \delta\tilde{x}) \\
 \cos x &\simeq \cos(\tilde{x} + \delta\tilde{x}) \\
 &= \cos(x_k + h + \delta\tilde{x}) \\
 &= \cos x_k \cos(h + \delta\tilde{x}) - \sin x_k \sin(h + \delta\tilde{x})
 \end{aligned}$$

The quantities  $\sin x_k$  and  $\cos x_k$  are extremely close to  $s_k$  and  $c_k$  (this is the core of the accurate table method). The terms in  $h + \delta\tilde{x}$  may be approximated by polynomials chosen to respect the parity of sin and cos and their values at 0:

$$\begin{aligned}
 \sin t &\simeq t + t^3 p_s(t^2) \\
 \cos t &\simeq 1 + t^2 p_c(t^2)
 \end{aligned}$$

After neglecting the terms in  $\delta\tilde{x}$  that do not contribute to the final result we obtain<sup>9</sup>:

$$\begin{aligned}\sin x &\approx s_k(1 + h(h + 2\delta\tilde{x})p_c(h^2)) + c_k(h + \delta\tilde{x} + h^3p_s(h^2)) \\ &= (s_k + c_k h) + s_k h(h + 2\delta\tilde{x})p_c(h^2) + c_k h^3p_s(h^2) + c_k \delta\tilde{x} \\ \cos x &\approx c_k(1 + h(h + 2\delta\tilde{x})p_c(h^2)) - s_k(h + \delta\tilde{x} + h^3p_s(h^2)) \\ &= (c_k - s_k h) + c_k h(h + 2\delta\tilde{x})p_c(h^2) - s_k h^3p_s(h^2) - s_k \delta\tilde{x}\end{aligned}$$

where  $h(h + 2\delta\tilde{x})$  is an approximation of  $(h + \delta\tilde{x})^2$ . For accuracy reasons, the leading term of these formulæ must be computed with extra accuracy.

We are now going to look at the techniques used to obtain polynomial approximation and we will later analyze the errors committed by these formulæ.

## Polynomial Approximations

The *Mathematica* function `GeneralMiniMaxApproximation` produces a minimax polynomial  $p$  such that  $p(q(t))$  approximates a function  $f(t)$  by minimizing the  $L^\infty$  norm of the *residual* function defined as:

$$\mathcal{E}(t) := \frac{f(t) - p(q(t))}{g(t)}$$

The residual is equioscillatory and bounded. The choice of  $f$  and  $q$  is dictated by the properties desired for the polynomial approximation: parity, value at zero, etc. The error function  $g$ , however, can be freely chosen to achieve the desired error bound on the approximation (e.g., minimizing the relative or absolute error).

### Sin Near Zero

Near zero we are looking for an approximation that minimizes the relative error of the result over the interval  $[0, \Delta]$ , where  $\Delta$  is chosen so that  $\Delta^2 \ll 1$ .

We are therefore calling `GeneralMiniMaxApproximation` with:

$$\begin{cases} q(t) &:= t^2 \\ f(t) &:= \frac{\sin t - t}{t^3} \\ g(t) &:= \frac{\sin t}{t^3} \end{cases}$$

The degree of  $p_{s_0}$  is chosen so that the approximation error is less than  $u(\Delta^2)$ .

In practice we choose  $\Delta = 2^{-10}$ , and compute a degree-1 polynomial with a residual smaller than  $2^{-85.560}$  (before rounding the coefficients to machine numbers).

### Sin and Cos Around Table Entries

Around table entries, it would be natural to approximate  $\cos$  by a polynomial minimizing the absolute or relative error. Unfortunately this creates singularities in the error analysis. To see why, consider the term  $h(h + 2\delta\tilde{x})p_c(h^2)$  which appears in the computations above, and expand it by making the error function  $g$  and the residual  $\mathcal{E}$  explicit:

$$\begin{aligned}h(h + 2\delta\tilde{x})p_c(h^2) &= h(h + 2\delta\tilde{x})\left(\frac{\cos h - 1}{h^2} + g(h)\mathcal{E}(h)\right) \\ &= \cos h - 1 + 2\delta\tilde{x}\frac{\cos h - 1}{h} + 2h\delta\tilde{x}g(h)\mathcal{E}(h) + h^2g(h)\mathcal{E}(h)\end{aligned}$$

When  $h \rightarrow 0$ , the dominant term is  $2h\delta\tilde{x}g(h)\mathcal{E}(h)$ . Since  $\mathcal{E}(h)$  is bounded, this term is of order  $\mathcal{O}(hg(h))$ . To obtain an absolute or relative error on  $\cos$  we would need

---

<sup>9</sup>The alert reader will note that the formulæ for  $\sin$  and  $\cos$  are very similar, which creates an opportunity for implementing a function that computes both together with significant performance savings.

to choose  $g(h) = 1/h^2$  or  $g(h) = \cos h/h^2$ , respectively. However, these choices yield a dominant term of order  $\mathcal{O}(1/h)$  which diverges when  $h \rightarrow 0$ . The root cause is that we are computing an approximation based on  $h$  alone, but we use it in combination with the term in  $\delta\tilde{x}$  which was not part of the optimization.

To avoid this issue, we choose  $g(h) = (\cos h - 1)/h^2$  which minimizes the relative error on  $p_c(h)$  and leads to a straightforward error analysis. Therefore, for  $\cos$  we call `GeneralMiniMaxApproximation` with:

$$\begin{cases} q(t) &:= t^2 \\ f(t) &:= \frac{\cos t - 1}{t^2} \\ g(t) &:= \frac{\cos t - 1}{t^2} \end{cases}$$

For  $\sin$  there is no such issue, therefore we simply minimize the relative error on  $\sin t$ :

$$\begin{cases} q(t) &:= t^2 \\ f(t) &:= \frac{\sin t - t}{t^3} \\ g(t) &:= \frac{\sin t}{t^3} \end{cases}$$

The minimax computation results in a 1-degree polynomial  $p_s$  with with a residual smaller than  $2^{-85.534}$  and a 1-degree polynomial  $p_c$  with a residual smaller than  $2^{-51.466}$  (before rounding the coefficients to machine numbers).

## Error Analysis

We use  $\llbracket expr \rrbracket$  to denote evaluation where appropriate rounding happens on each literal or operation of the expression  $expr$ . This notation is used for error intervals computed using the function `IEEEEvaluateWithRelativeError` in file `mathematica/ieee754_floating_point_evaluation.wl`. In particular,  $\llbracket \cdot \rrbracket$  takes into account that the evaluation of the polynomials  $p_{s0}$ ,  $p_s$ , and  $p_c$  has two sources of errors: the rounding of the coefficients to machine numbers, and the error due to the floating-point operations; this may lead to asymmetrical error intervals.

TODO(phl): Try to tune the high-degree coefficient after rounding the low-degree one.

We assume that  $\tilde{x}$  is positive ( $\delta\tilde{x}$  may be positive or negative) and that the rounding direction is `roundTiesToEven`. We *do not* assume that the machine has an FMA instruction in our error analysis: even though we actually use this instruction when available (for performance) our rounding bounds are valid even in the absence of an FMA<sup>10</sup>.

## Reduced Angle

Argument reduction took the input angle  $x$  and produced a pair  $(\tilde{x}, \delta\tilde{x})$  approximating the angle reduced modulo  $\frac{\pi}{2}$ . That approximation is correct to  $M + \kappa_3$  bits. We therefore have, from equation (11):

$$x \equiv \tilde{x} + \delta\tilde{x} + \zeta_0 \tilde{x} \pmod{\frac{\pi}{4}} \quad |\zeta_0| \leq 2^{-\kappa_3 - M}$$

For simplicity of the error analysis we can assume  $0 \leq x \leq \frac{\pi}{4}$ : even though we do not do argument reduction when  $x$  is less than  $\frac{\pi}{4}$ , the error analysis would work exactly the same if  $x$  was simply provided as a high-accuracy value (e.g., a double-double) and reduced using the techniques above.

<sup>10</sup>This means that the accuracy computed below is possibly pessimistic by a few hundredth of bits, which has no effect on performance.

## Rounding Test

The rounding test described in [Mul+10, pp. 397–400] is done by comparing  $y$  and  $\llbracket y + \llbracket \delta y e \rrbracket \rrbracket$  for equality, where  $e > 1$  is computed based on the bound on the relative error of  $y + \delta y$  with respect to  $\sin x$  or  $\cos x$ . In our case, when an FMA is available, we want to compute the second part as  $\llbracket y + \delta y e \rrbracket$ . We must analyse what this implies for the computation of  $e$ .

The proof of the rounding test is rather convoluted, but for our purpose the important part is the implication<sup>11</sup>, when  $y$  is not a power of 2 or  $\delta y \geq 0$ :

$$y = \llbracket y + \llbracket \delta y e \rrbracket \rrbracket \Rightarrow \llbracket \delta y e \rrbracket \leq \frac{u(y)}{2} \Rightarrow \delta y e \left(1 - u\left(\frac{1}{2}\right)\right) \leq \frac{u(y)}{2}$$

When an FMA is used the smaller term has to be less than half a ULP of the larger term, otherwise the result would round away from  $y$ . This gives the tighter bound:

$$y = \llbracket y + \delta y e \rrbracket \Rightarrow \delta y e \leq \frac{u(y)}{2}$$

Similarly when  $y$  is a power of 2 and  $\delta y < 0$  the proof depends on the implications:

$$y = \llbracket y + \llbracket \delta y e \rrbracket \rrbracket \Rightarrow \llbracket \delta y e \rrbracket \leq \frac{u(y)}{4} \Rightarrow -\delta y e \left(1 - u\left(\frac{1}{2}\right)\right) \leq \frac{u(y)}{4}$$

When an FMA is used the smaller term, which is negative, has to be less than half a ULP *below* the larger term otherwise the result would round below  $y$ . Because the ULP below is half the ULP above, this gives the tighter bound:

$$y = \llbracket y + \delta y e \rrbracket \Rightarrow -\delta y e \leq \frac{u(y)}{4}$$

Putting these tighter bounds together we obtain:

$$k := \lfloor -\log_2 \bar{\epsilon}_1 - M \rfloor$$

$$e_{\text{FMA}} := 1 + \frac{1 + 2^{M+1} \bar{\epsilon}_1}{1 - \bar{\epsilon}_1 - 2^{-k+1}}$$

In the absence of FMA, [Mul+10] gives the same value for  $k$  and for  $e$ :

$$e_{\text{NoFMA}} := (1 - 2^{-M})^{-1} \left(1 + \frac{1 + 2^{M+1} \bar{\epsilon}_1}{1 - \bar{\epsilon}_1 - 2^{-k+1}}\right)$$

We see immediately that  $e_{\text{FMA}} < e_{\text{NoFMA}}$ , and it is therefore correct to perform the rounding test using  $e_{\text{NoFMA}}$ , even if an FMA is used to compute  $\llbracket y + \delta y e \rrbracket$ .

Obviously, if the above expression does not produce a machine number, it must be rounded towards  $+\infty$ . In each of the error analyses below we document the concrete factor  $\llbracket e_{\text{NoFMA}} \rrbracket_+$  to use for the rounding test.

## Sin Near Zero

If  $|\tilde{x}| \leq \Delta = 2^{-10}$  the steps of the computation are as follows:

$$\begin{cases} t_1 &:= \llbracket p_{s0}(\tilde{x}^2) \rrbracket \dots \\ t_2 &:= \llbracket \llbracket \tilde{x}^2 \rrbracket \tilde{x} \rrbracket \\ t_3 &:= \llbracket \llbracket t_1 t_2 \rrbracket + \delta \tilde{x} \rrbracket \\ y &:= \tilde{x} \\ \delta y &:= t_3 \end{cases}$$

<sup>11</sup>We adapt the proof to use our notation.

We gave above the upper bound of the residual of the minimax approximation. It may be rewritten as:

$$p_{s0}(t^2) = \frac{(1 + \zeta_1) \sin t - t}{t^3} \quad |\zeta_1| < 2^{-85.560}$$

for  $|t| \leq \Delta$ .

The errors committed at each step are as follows:

$$\begin{cases} t_1 = p_{s0}(\tilde{x}^2)(1 + \zeta_2) & \zeta_2 \in ]-2^{-52.415}, 2^{-53.999}[ \\ = \frac{(1 + \zeta_1) \sin \tilde{x} - x}{\tilde{x}^3} (1 + \zeta_2) \\ t_2 = \tilde{x}^3(1 + \delta_1)(1 + \delta_2) \\ t_3 = (t_1 t_2(1 + \delta_3) + \delta \tilde{x})(1 + \delta_4) \end{cases}$$

The relative error of the entire computation is:

$$r(\tilde{x}, \delta \tilde{x}) := \frac{y + \delta y}{\sin(\tilde{x} + \delta \tilde{x} + \zeta_0 \tilde{x})} - 1$$

TODO(phl): See if this can be proven rigourously.

The function  $r$  can be computed using *Mathematica* interval arithmetic over the triangular domain  $|\tilde{x}| \leq \Delta$ ,  $|\delta \tilde{x}| \leq 2^{-M} |\tilde{x}|$ . Plotting shows that this function reaches its extrema at the corners of the domain, which is logical because we expect the minimax polynomial to reach its largest errors on the bounds of the optimization interval. In practice we find that  $|r(\tilde{x}, \delta \tilde{x})| < 2^{-70.583}$  (as expected, this is a bit worse than the error originating from the angle reduction). The rounding test must be done with  $e = 0 \times 1.0000' AAD0' 391A' Dp0$ .

TODO(phl):  $\kappa_3 = 19$  would gain 0.7 bits. Worthwhile?

As mentioned above (see Overview of the Algorithms for Sin and Cos), in these steps we do not compute the terms  $\tilde{x}^n \delta \tilde{x}$  for  $n > 0$ . The largest such term is for  $n = 2$  and the relative error that it induces is:

$$-\frac{1}{2} \frac{\tilde{x}^2 \delta \tilde{x}}{\sin(\tilde{x} + \delta \tilde{x})}$$

This function reaches its extrema on the corners of the triangular domain and is found to be smaller than  $2^{-73.999}$ , so its contribution to the overall relative error would be very small.

## Sin and Cos Around Table Entries

If  $\tilde{x} \in \mathcal{I}_k$  we first compute  $h = \tilde{x} - x_k$ . It is essential that this subtraction be exact, and therefore that the conditions of Sterbenz's lemma be met: we must have  $\tilde{x}/2 \leq x_k \leq 2\tilde{x}$ . Based on the range of  $\tilde{x}$ , a sufficient condition for this is:

$$\frac{2k+1}{2} \Delta \leq x_k \leq 2(2k-1)\Delta$$

Because  $x_k$  is close to  $2k\Delta$  the left part of this condition is trivially met for all  $k > 0$ . However the second part is only trivially met if  $k > 1$ . If  $k = 1$  it becomes  $x_1 \leq 2\Delta$ . In other words, when building the accurate tables we must only look for  $x_1$  below  $2\Delta$ . This is done by skipping the loop at step 5 of algorithm 3.

For sin we do not use the interval  $\mathcal{I}_0$  (see Sin Near Zero, above) but for cos we do. That interval is special because  $x_0 = 0$ : the computation of  $h$  is therefore trivially correct.

Interestingly,  $x_1 < 2\Delta$  is also the necessary condition to compute  $h$   $c_k + s_k$  exactly using an FMA as explained in section 2.1 of [SZ05]. The subtraction in the second step of their algorithm must be exact, therefore the following Sterbenz inequalities must hold<sup>12</sup>:

$$\frac{s_k}{2} \leq h \quad c_k + s_k \leq 2 s_k$$

<sup>12</sup>Note that  $\mathcal{I}_0$  is uninteresting because  $c_0 = 1$  and  $s_0 = 0$  so the computation is trivially exact.



This may be rewritten as:

$$-\frac{1}{2} \tan x_k \leq h \leq \tan x_k$$

Let  $x_k = 2k\Delta + \epsilon_k$  for some small  $\epsilon_k$ . By the definition of  $h$  we have:

$$-\Delta - \epsilon_k \leq h \leq \Delta - \epsilon_k$$

Now we know that  $x_k < \tan x_k$  and since  $\Delta - \epsilon_k < 2k\Delta + \epsilon_k$  for all  $k > 0$  the right side of the Sterbenz condition is always true. The left side is more interesting though as:

$$-\frac{1}{2} \tan x_k < -\frac{x_k}{2} = -k\Delta - \frac{\epsilon_k}{2}$$

For  $k > 1$ , this quantity is clearly smaller than  $-\Delta - \epsilon_k$ , the lower bound of  $h$ . However, when  $k = 1$ ,  $-\Delta - \epsilon_1/2 \leq -\Delta - \epsilon_1$  requires that  $\epsilon_1 \leq 0$ , in other words, that  $x_1$  be below  $2\Delta$ . Note that we do not need to go through a similar proof for the exactness of  $-h s_k + c_k$  because  $h s_k$  is much smaller than  $c_k$  for all  $k$ .

Finally, for the purpose of error analysis, the relative errors of the minimax polynomials may be rewritten as:

$$\begin{aligned} p_s(t^2) &= \frac{\sin t (1 + \zeta_1) - t}{t^3} & |\zeta_1| &< 2^{-85.534} \\ p_c(t^2) &= \frac{\cos t - 1}{t^2} (1 + \zeta_2) & |\zeta_2| &< 2^{-51.466} \end{aligned}$$

## Sin

The first step of the computation is to evaluate  $h c_k + s_k$  exactly<sup>13</sup> using an FMA:

$$(z, \delta z) = t_0 = h c_k + s_k$$

TODO(phl): Examine what happens when we don't have FMA.

where  $z$  and  $\delta z$  have nonoverlapping significands. For the purpose of describing the computation and analysing errors we will write  $\delta z = t_0 \delta_0$  and  $z = t_0(1 - \delta_0)$ , where the two terms are exact.

The remaining steps of the computation are then as follows:

$$\left\{ \begin{array}{ll} t_1 &:= \llbracket p_s(h^2) \rrbracket \dots \\ t_2 &:= \llbracket p_c(h^2) \rrbracket \dots \\ t_3 &:= \llbracket h \llbracket h + (\delta \tilde{x} + \delta \tilde{x}) \rrbracket \rrbracket \\ t_4 &:= \llbracket \llbracket h^2 \rrbracket h \rrbracket \\ t_5 &:= \llbracket \llbracket \llbracket s_k \rrbracket t_3 \rrbracket t_2 \rrbracket \\ t_6 &:= \llbracket \llbracket t_4 t_1 \rrbracket + \delta \tilde{x} \rrbracket \\ t_7 &:= \llbracket \llbracket \llbracket c_k \rrbracket t_6 \rrbracket + t_5 \rrbracket \\ t_8 &:= \llbracket \llbracket t_0 \delta_0 \rrbracket + t_7 \rrbracket \\ y &:= t_0(1 - \delta_0) \\ \delta y &:= t_8 \end{array} \right.$$

where we have made the rounding of the accurate table elements  $s_k$  and  $c_k$  explicit, and used the fact that the computation of  $\delta \tilde{x} + \delta \tilde{x}$  is exact.

<sup>13</sup>See Annex: Proof of the FMA Computation for a proof of the correctness of this computation.

The errors committed at each step are as follows:

$$\begin{cases} t_1 = p_s(h^2)(1 + \zeta_3) & \zeta_3 \in ]-2^{-53.221}, 2^{-52.808}[ \\ = \frac{\sin h(1 + \zeta_1) - h}{h^3}(1 + \zeta_3) \\ t_2 = p_c(h^2)(1 + \zeta_4) & \zeta_4 \in ]-2^{-52.855}, 2^{-53.160}[ \\ = \frac{\cos h - 1}{h^2}(1 + \zeta_2)(1 + \zeta_4) \\ t_3 = h(h + 2\delta\tilde{x})(1 + \delta_1)(1 + \delta_2) \\ t_4 = h^3(1 + \delta_3)(1 + \delta_4) \\ t_5 = \llbracket s_k \rrbracket t_3 t_2 (1 + \delta_5)(1 + \delta_6) \\ t_6 = (t_4 t_1 (1 + \delta_7) + \delta\tilde{x})(1 + \delta_8) \\ t_7 = (\llbracket c_k \rrbracket t_6 (1 + \delta_9) + t_5)(1 + \delta_{10}) \\ t_8 = (t_0 \delta_0 (1 + \delta_{11}) + t_7)(1 + \delta_{12}) \end{cases}$$

The relative error of the entire computation is:

$$r(\tilde{x}, \delta\tilde{x}) := \frac{y + \delta y}{\sin(\tilde{x} + \delta\tilde{x} + \zeta_0 \tilde{x})} - 1$$

Great care is required when evaluating this expression using interval arithmetic because the terms in  $\delta_0$  in  $t_8$  and  $y$  could be considered independent and lead to a useless bound. To avoid this issue  $\delta_0$  must be factored out to appear only once. Specifically the  $\delta_0$  term after factoring is:

$$(h\llbracket c_k \rrbracket + \llbracket s_k \rrbracket)((1 + \delta_{11})(1 + \delta_{12}) - 1)$$

which is pleasantly small, and even though  $\delta_{12}$  also appears in the term that does not involve  $\delta_0$ , the effect of replicating it is acceptable.

The function  $r$  must be evaluated separately for each interval around  $x_k$ , with  $\tilde{x} \in [(2k-1)\Delta, (2k+1)\Delta]$  and  $|\delta\tilde{x}| < |\tilde{x}|/2^M$ . This defines a trapezoidal domain and  $r$  reaches its extrema at the corners of that domain. In practice we find using *Mathematica* interval arithmetic that  $|r(\tilde{x}, \delta\tilde{x})| < 2^{-68.734}$ . The relative error for  $k = 1$  is particularly high and is an outlier. Other than that, the worst relative errors correspond to table entries that have many ones after their accurate zeroes (or many zeroes after their accurate ones). The rounding test must be done with  $e = 0 \times 1.0002'6752'8572'Dp0$ .

As mentioned above (see Overview of the Algorithms for Sin and Cos), in these steps we do not compute the terms  $h^n \delta\tilde{x}$  for  $n > 1$ . The largest such term is for  $n = 2$  and the relative error that it induces is:

$$-\frac{1}{2} \frac{h^2 \delta\tilde{x} \cos x_k}{\sin(\tilde{x} + \delta\tilde{x})}$$

This function reaches its extrema on the corners of the trapezoidal domain and is found to be smaller than  $2^{-73.999}$  for all  $k$ , so its contribution to the overall relative error would be very small. (The term for  $n = 1$ , on the other hand, induces an error of the order of  $2^{-63.348}$  so it must be computed.)

## Cos

TODO(phl): Prove that this works, and examine what happens when we don't have FMA.

The first step of the computation is to evaluate  $-h s_k + c_k$  exactly using an FMA:

$$(z, \delta z) = t_0 = -h s_k + c_k$$

where  $z$  and  $\delta z$  have nonoverlapping significands. For the purpose of describing the computation and analysing errors we will write  $\delta z = t_0 \delta_0$  and  $z = t_0(1 - \delta_0)$ , where the two terms are exact.

The remaining steps of the computation are then as follows:

$$\begin{cases} t_1 &:= \llbracket p_s(h^2) \rrbracket \dots \\ t_2 &:= \llbracket p_c(h^2) \rrbracket \dots \\ t_3 &:= \llbracket h \llbracket h + (\delta\tilde{x} + \delta\tilde{x}) \rrbracket \rrbracket \\ t_4 &:= \llbracket \llbracket h^2 \rrbracket h \rrbracket \\ t_5 &:= \llbracket \llbracket c_k \rrbracket t_3 \rrbracket t_2 \rrbracket \\ t_6 &:= \llbracket \llbracket t_4 t_1 \rrbracket + \delta\tilde{x} \rrbracket \\ t_7 &:= \llbracket -\llbracket s_k \rrbracket t_6 \rrbracket + t_5 \rrbracket \\ t_8 &:= \llbracket \llbracket t_0 \delta_0 \rrbracket + t_7 \rrbracket \\ y &:= t_0(1 - \delta_0) \\ \delta y &:= t_8 \end{cases}$$

where we have made the rounding of the accurate table elements  $s_k$  and  $c_k$  explicit, and used the fact that the computation of  $\delta\tilde{x} + \delta\tilde{x}$  is exact.

The errors committed at each step are as follows:

$$\begin{cases} t_1 = p_s(h^2)(1 + \zeta_3) & \zeta_3 \in ]-2^{-53.221}, 2^{-52.808}[ \\ = \frac{\sin h(1 + \zeta_1) - h}{h^3}(1 + \zeta_3) \\ t_2 = p_c(h^2)(1 + \zeta_4) & \zeta_4 \in ]-2^{-52.855}, 2^{-53.160}[ \\ = \frac{\cos h - 1}{h^2}(1 + \zeta_2)(1 + \zeta_4) \\ t_3 = h(h + 2\delta\tilde{x})(1 + \delta_1)(1 + \delta_2) \\ t_4 = h^3(1 + \delta_3)(1 + \delta_4) \\ t_5 = \llbracket c_k \rrbracket t_3 t_2 (1 + \delta_5)(1 + \delta_6) \\ t_6 = (t_4 t_1 (1 + \delta_7) + \delta\tilde{x})(1 + \delta_8) \\ t_7 = (\llbracket s_k \rrbracket t_6 (1 + \delta_9) + t_5)(1 + \delta_{10}) \\ t_8 = (t_0 \delta_0 (1 + \delta_{11}) + t_7)(1 + \delta_{12}) \end{cases}$$

The relative error of the entire computation is:

$$r(\tilde{x}, \delta\tilde{x}) := \frac{y + \delta y}{\sin(\tilde{x} + \delta\tilde{x} + \zeta_0 \tilde{x})} - 1$$

As explained above (see Sin) the expression for  $r$  must be rewritten so that  $\delta_0$  only appears once. In this case the  $\delta_0$  term is:

$$(h\llbracket s_k \rrbracket - \llbracket c_k \rrbracket)(1 - (1 + \delta_{11})(1 + \delta_{12}))$$

and a *Mathematica* computation yields  $|r(\tilde{x}, \delta\tilde{x})| < 2^{-69.217}$ . Again, the largest errors come from table entries that have many ones after their accurate zeroes (or many zeroes after their accurate ones). The rounding test must be done with  $e = 0 \times 1.0001' B838' 5D8B' 6p0$ .

As mentioned above (see Overview of the Algorithms for Sin and Cos), in these steps we do not compute the terms  $h^n \delta\tilde{x}$  for  $n > 1$ . The largest such term is for  $n = 2$  and the relative error that it induces is:

$$\frac{1}{2} \frac{h^2 \delta\tilde{x} \sin x_k}{\cos(\tilde{x} + \delta\tilde{x})}$$

This function reaches its extrema on the corners of the trapezoidal domain and is found to be smaller than  $2^{-74.348}$  for all  $k$ , so its contribution to the overall relative error would be very small. (The term for  $n = 1$ , on the other hand, induces an error of the order of  $2^{-63.346}$  so it must be computed.)

## Annex: Proof of the FMA Computation

[SZ05] presents, in section 2.1, the following algorithm to compute  $s_k + c_k h$  exactly:

$$\begin{cases} h' &= \llbracket s_k + c_k h \rrbracket \\ t &= h' - s_k \\ l &= \llbracket t - c_k h \rrbracket \end{cases}$$

where the computation of  $t$  is exact if  $h'$  and  $s_k$  verify the conditions of Sterbenz's Lemma. [SZ05] then states without proof that  $h' + l = s_k + c_k h$ .

The first thing to note is that there is a sign error in that last formula. Let's make the absolute errors  $\zeta_i$  of each computation explicit:

$$\begin{cases} h' &= s_k + c_k h + \zeta_1 \\ t &= h' - s_k = c_k h + \zeta_1 \\ l &= \llbracket t - c_k h \rrbracket = \llbracket \zeta_1 \rrbracket = \zeta_1 + \zeta_2 \end{cases}$$

It is clear that the sum  $h' + l = s_k + c_k h + 2\zeta_1 + \zeta_2$  doubles the error  $\zeta_1$  (which is large) while the difference  $h' - l = s_k + c_k h - \zeta_2$  eliminates it and only retains  $\zeta_2$  (which is much smaller).

To prove that  $h' - l = s_k + c_k h$  we must prove that  $\zeta_2 = 0$ , which is equivalent to  $\llbracket \zeta_1 \rrbracket = \zeta_1$ . In other words, if  $\zeta_1$  is a machine number then  $\zeta_2 = 0$  and the equality holds.

The first thing that we show is that the exact value  $s_k + c_k h$  has at most  $2M + 1$  bits in its binary representation. Note that  $h = x - x_k$ , where the subtraction is exact, and therefore there exists  $N \in [0, 2^M[ \cap \mathbb{N}$  such that  $h = N u(x_k)$ . Moreover, because  $x_k \in [0, \frac{\pi}{4}]$  we have  $1/2 < c_k < 1$  and therefore  $u(c_k) = u(1/2) = 2^{-M}$ . Finally,  $x_k/2 < s_k < x_k$ , which implies that  $u(s_k)$  is either  $u(x_k)$  or  $u(x_k)/2$ . We need to distinguish two cases:

– If  $u(s_k) = u(x_k)$  then:

$$\begin{aligned} s_k + c_k h &= (m(s_k) + m(c_k)2^{-M}N) u(x_k) \\ &= (m(s_k)2^M + m(c_k)N) u(x_k)2^{-M} \end{aligned}$$

where the quantity in parentheses is an integer strictly less than  $2^{2M+1}$  by the properties of the significand function  $m$ .

– If  $u(s_k) = u(x_k)/2$  then:

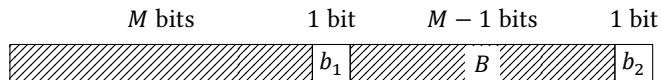
$$\begin{aligned} s_k + c_k h &= \left( \frac{m(s_k)}{2} + m(c_k)2^{-M}N \right) u(x_k) \\ &= (m(s_k)2^{M-1} + m(c_k)N) u(x_k)2^{-M} \end{aligned}$$

where the quantity in parentheses is an integer strictly less than  $3 \cdot 2^{2M-1} < 2^{2M+1}$ .

Taken together these inequalities prove that  $s_k + c_k h$  can be represented using  $2M + 1$  significand bits with an ULP of  $u(x_k)2^{-M}$ .

The computation of  $h'$  is going to round the most significant  $M$  bits of that significand and leave a remainder,  $\zeta_1$ , which has potentially  $M + 1$  bits. We are now going to show that in fact  $\zeta_1$  has only  $M$  bits in its significand.

The following picture represents the significand of  $s_k + c_k h$ , with the least significant bit on the right. The leftmost  $M$  bits are the part that is rounded (either up or down). The rightmost  $M + 1$  bits are the ones that contribute to  $\zeta_1$ . We consider different cases based on the values of bits  $b_1$  and  $b_2$ .



- If  $b_2 = 0$ , the significand of  $s_k + c_k h$  effectively only has  $2M$  bits, and therefore the significand of  $\zeta_1$  trivially has at most  $M$  bits.
- If  $b_2 = 1$ , we have two cases depending on the value of  $b_1$ :
  - If  $b_1 = 0$ , the leading  $M$  bits are rounded *down* and therefore the significand of the positive quantity  $s_k + c_k h - h' = -\zeta_1$  is exactly made of the bits  $(b_1, B, b_2)$ . Since  $b_1 = 0$ , this is only the  $M$  bits  $(B, b_2)$  and therefore  $\zeta_1$  is a machine number (with a “hole” separating it from the significand of  $h'$ ).
  - If  $b_1 = 1$ , the leading  $M$  bits were rounded *up*<sup>14</sup> and therefore the significand of the positive quantity  $h' - s_k + c_k h = \zeta_1$  is obtained by taking the two's complement of the bits  $(b_1, B, b_2)$ . This is exactly  $(0, \neg B, 1)$  where  $\neg$  designates binary negation. Again, this is only  $M$  bits because of the leading zero and  $\zeta_1$  is a machine number.

## References

- [Gal86] S. Gal. “Computing elementary functions: A new approach for achieving high accuracy and good performance”. In: *Accurate Scientific Computations* (Bad Neuenahr, Federal Republic of Germany, Mar. 12–14, 1985). Ed. by W. Miranker and R. Toupin. Lecture Notes in Computer Science 235. Springer Berlin Heidelberg, Sept. 1986, pp. 1–16. DOI: 10.1007/3-540-16798-6\_1.
- [GB91] S. Gal and B. Bachelis. “An Accurate Elementary Mathematical Library for the IEEE Floating Point Standard”. In: *ACM Transactions on Mathematical Software* 17.1 (Mar. 1991), pp. 26–45.
- [Higo2] N. J. Higham. *Accuracy and Stability of Numerical Algorithms*. Society for Industrial and Applied Mathematics, 2002.
- [HLBo8] Y. Hida, X. S. Li, and D. H. Bailey. “Library for Double-Double and Quad-Double Arithmetic”. Preprint at <https://www.davidhbailey.com/dhbpapers/qd.pdf>. May 8, 2008.
- [HPS14] J. Hoffstein, J. Pipher, and J. H. Silverman. *An Introduction to Mathematical Cryptography*. Undergraduate Texts in Mathematics. Springer, 2014.
- [MR22] J.-M. Muller and L. Rideau. “Formalization of Double-Word Arithmetic, and Comments on “Tight and Rigorous Error Bounds for Basic Building Blocks of Double-Word Arithmetic””. In: *ACM Transactions on Mathematical Software* 48.1 (Feb. 2022), pp. 1–24. DOI: 10.1145/3484514.
- [Mul+10] J.-M. Muller, N. Brisebarre, F. De Dinechin, C.-P. Jeannerod, V. Lefèvre, G. Melquiond, N. Revol, D. Stehlé, and S. Torres. *Handbook of Floating-Point Arithmetic*. Birkhäuser, 2010.
- [NS09] P. Q. Nguyen and D. Stehlé. “An LLL Algorithm with Quadratic Complexity”. In: *SIAM Journal on Computing* 39.3 (2009), pp. 874–903. DOI: 10.1137/070705702.
- [SZ04] D. Stehlé and P. Zimmermann. *Gal’s Accurate Tables Method Revisited*. Research Report RR-5359. INRIA, Oct. 2004. eprint: <https://inria.hal.science/inria-00070644v1/file/rr-5359.pdf>.
- [SZ05] D. Stehlé and P. Zimmermann. “Gal’s accurate tables method revisited”. In: *17th IEEE Symposium on Computer Arithmetic (ARITH’05)* (Cape Cod, MA, USA, June 27–29, 2005). Ed. by P. Montuschi and E. Schwarz. IEEE Computer Society, June 2005, pp. 257–264. DOI: 10.1109/ARITH.2005.24.

<sup>14</sup>Because  $b_2 = 1$  the trailing bits are above the midpoint, so the handling of ties is irrelevant.

- [SZG22] A. Sibidanov, P. Zimmermann, and S. Glondou. “The CORE-MATH Project”. In: *2022 IEEE 29th Symposium on Computer Arithmetic (ARITH)*. IEEE, Sept. 2022, pp. 26–34.  
doi: 10.1109/ARITH54963.2022.00014.  
eprint: <https://inria.hal.science/hal-03721525v3/file/core-math-final.pdf>.
- [ZSG+24] P. Zimmermann, A. Sibidanov, S. Glondou, et al. *The CORE-MATH Project*. Software. Apr. 2024.