

The CORE-MATH acos is correctly rounded

Paul Zimmermann

April 20, 2026

Abstract—We prove that the CORE-MATH binary64 arc-cosine function is correctly rounded.

Index Terms—IEEE 754, binary64 format, correct rounding.

We consider the CORE-MATH `acos.c` file from commit 7457819. The function `acos(x)` is only defined for $|x| \leq 1$ (for $|x| > 1$ it returns NaN), and it decreases from π to 0 on $[-1, 1]$ (see Figure 1). We consider all four rounding modes

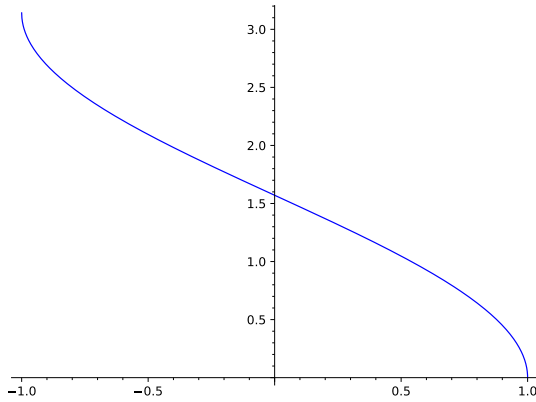


Fig. 1. Graph of $\text{acos}(x)$ on $[-1, 1]$.

from the C23 standard: to nearest-even, toward zero, toward $+\infty$, and toward $-\infty$.

I. BRANCH $|x| < 2^{-15}$

In this case we use Algorithm `FastPathTiny`, which approximates $\text{acos}(x)$ by $p - (x + x^3/6)$, where $p = p_h + p_\ell$ is a double-double approximation of $\pi/2$. The `fasttwosum` routine

Algorithm 1 (FastPathTiny)

Input: x a binary64 number, $|x| \leq 2^{-15}$

Output: the correct rounding $\circ(\text{acos}(x))$ or FAIL

- 1: $p_h = 0x1.921fb54442d18p+0$
 - 2: $p_\ell = 0x1.1a62633145c07p-54$
 - 3: $c = -0x1.555555555555555p-3 \triangleright$ approximates $-1/6$
 - 4: $t = \circ(x^2)$, $u = \circ(c \cdot x)$, $v = \circ(t \cdot u)$
 - 5: $h, w = \text{fasttwosum}(p_h, -x)$
 - 6: $\ell = \circ(v + \circ(w + p_\ell))$, $\epsilon = 0x1.34p-79$
 - 7: $\ell_b = h + \circ(\ell - \epsilon)$, $u_b = h + \circ(\ell + \epsilon)$
 - 8: if $\circ(\ell_b) = \circ(u_b)$ return $\circ(\ell_b)$ else return FAIL
-

is defined in Algorithm 2.

Lemma 1: For x a binary64 number, $|x| < 2^{-15}$, and any rounding mode, if Algorithm `FastPathTiny` does not return FAIL, then it yields the correct rounding of $\text{acos}(x)$.

Algorithm 2 (fasttwosum)

Input: a, b binary64 numbers

Output: h, ℓ such that $h + \ell$ approximates $a + b$

- 1: $h \leftarrow \circ(a + b)$
 - 2: $t \leftarrow \circ(h - a)$
 - 3: $\ell \leftarrow \circ(b - t)$
-

Proof: Sollya [2] proves that the absolute error is bounded by $2^{-78.736}$:

```
prec=64;
ph=0x1.921fb54442d18p+0;
pl=0x1.1a62633145c07p-54;
c=0x1.555555555555555p-3;
q=x+c*x^3;
y=ph+pl-q;
d=[-2^-15], 2^-15];
dirtyinfnorm(y-acos(x), d);
1.9852337356788177956e-24
```

We thus have:

$$|[p - (x + cx^3)] - \text{acos}(x)| \leq 2^{-78.736}. \quad (1)$$

We then use the following Gappa script to bound the rounding error in Algorithm `FastPathTiny`, with `RND` substituted by all rounding modes:

```
@rnd = float<ieee_64, RND>;
ph = 0x1.921fb54442d18p+0;
pl = 0x1.1a62633145c07p-54;
c = -0x1.555555555555555p-3;
t = x*x;
tr rnd= x * x;
u = c*x;
ur rnd=c*x;
v = t*u;
vr rnd= tr*ur;
y = -x;
# h, w = fasttwosum(ph, y)
h rnd= ph + y;
w = rnd(-(h-(ph+y)));
l1 rnd= w+pl;
l rnd= vr+l1;
```

```
{ x in [-1b-15, 1b-15] -> x // ph in [-1, 1]
  /\ |(h+1)-((ph+pl)+(y+v))| in ? }
```

```
(h+w) - (ph+y) -> w - (-(h-(ph+y)));
```

$$(h+1) - ((ph+pl) + (y+v)) \rightarrow ((h+w) - (ph+y)) + ((l1 - (w+pl)) + (1 - (vr+l1)) + (vr-v));$$

This script shows that the absolute error is bounded by $2^{-98.0891}$:

$$|h + \ell - [p - (x + cx^3)]| \leq 2^{-98.0891}. \quad (2)$$

Combining Eq. (1) with Eq. (2), we get:

$$|h + \ell - \text{acos}(x)| < 2^{-78.735}.$$

With the goal `|l| in ?`, this script also proves $|\ell| \leq 2^{-47.5012}$. Thus the rounding error in $\circ(\ell \pm \epsilon)$ for ℓ_b and u_b at line 7 of Algorithm `FastPathTiny` is bounded by $\text{ulp}(2^{-47.5}) = 2^{-100}$. Thus since $\epsilon > 2^{-78.735} + 2^{-100}$, we have $\ell_b < \text{acos}(x) < u_b$. By monotonicity of rounding, if $\circ(\ell_b) = \circ(u_b)$, this concludes the proof. ■

REMARK: in the actual code, to avoid a spurious underflow in the computation of v , we set $v = 0$ when $|x| \leq x_0 := 0 \times 1 . \text{cb3b3869747f4p-55}$, since then $\text{acos}(x)$ rounds to the same value as $\text{acos}(0)$. The proof remains correct since this adds an absolute error less than $|cx^3| < 2^{-165}$. Moreover, for $|x| \leq x_0$, we always have $\ell_b = u_b$, thus the accurate path is never called (Lemma 2).

Lemma 2: For $|x| \leq x_0$, and any C23 rounding mode, Algorithm `FastPathTiny` never returns FAIL.

Proof: Since $x_0 < \text{ulp}(p_h)$, at line 5, h is either p_h , $p_h + \text{ulp}(p_h)$ or $p_h - \text{ulp}(p_h)$, depending on the rounding mode and the sign of x . We know from [Theorem 9-11] [3] that $h + w$ either equals $p_h - x$, or equals its rounding in double precision $\circ_{106}(p_h - x)$. Thus $h + w$ is decreasing with x , and since for a given rounding mode and sign of x , h is fixed, then w is decreasing with x . For each rounding mode, it thus suffices to check the property is satisfied for $x = \pm x_0$. And indeed, for $x = \pm x_0$, we have $\circ(\ell_b) = \circ(u_b)$ for all rounding modes. ■

A. Accurate path

This branch uses the same routine for the accurate path as in §II, which is proven correct in §IV.

II. BRANCH $2^{-15} \leq |x| < 2^{-4}$

For $|x| < 0.5$, the code computes the integer $j = \lfloor 2^7 \circ(x^2) \rfloor$, and a polynomial approximation of $\text{asin}(x)$ indexed by j is used:

$$q(x) = x \cdot (c_0 + c_1 + t(c_2 + c_3t + \dots + c_7t^5)),$$

where $t = \text{fma}(x, x, -2^{-7}j)$, and then $\text{acos}(x)$ is approximated by:

$$p_h + p_\ell - q(x), \quad (3)$$

where $p_h + p_\ell$ is a double-double approximation of $\pi/2$. The integer j goes from 0 for $2^{-15} \leq |x| \leq 2^{-4}$ to 32 for $|x|$ near 0.5. We only study the case $|x| < 2^{-4}$ here. For $2^{-4} \leq |x| < 0.5$, the correctness will be proven by exhaustive tests in Section III.

In the case $|x| < 2^{-4}$, we have $\circ(x^2) < 2^{-8}$, thus $j = 0$ (it is easy to check that for $x = \text{nextdown}(2^{-4})$,

$\text{RU}(x^2) < 2^{-8}$). In this $j = 0$ case, Algorithm `FastPathSmall` is used, with c_0, \dots, c_7 be the coefficients for index $j = 0$, and t simplifies to $\circ(x^2)$.

Algorithm 3 (FastPathSmall)

Input: x a binary64 number, $2^{-15} \leq |x| < 2^{-4}$

Output: the correct rounding $\circ(\text{acos}(x))$ or FAIL

- 1: $p_h = 0 \times 1 . 921 \text{fb}54442 \text{d}18 \text{p}+0$
 - 2: $p_\ell = 0 \times 1 . 1 \text{a}62633145 \text{c}07 \text{p}-54$
 - 3: $t = \circ(x^2)$, $z = -x$, $z_\ell = 0$
 - 4: $t_2 = \circ(t \cdot t)$
 - 5: $g_2 = \circ(c_2 + \circ(tc_3))$, $h_4 = \circ(c_4 + \circ(tc_5))$
 - 6: $h_6 = \circ(c_6 + \circ(tc_7))$, $g_4 = \circ(h_4 + \circ(t_2h_6))$
 - 7: $e_2 = \circ(g_2 + \circ(t_2g_4))$, $d = \circ(te_2)$
 - 8: $f_h = c_0$, $f_\ell = \circ(c_1 + d)$
 - 9: $f_h, f_\ell = \text{muldd}(z, z_\ell, f_h, f_\ell)$
 - 10: $h, \ell = \text{fastsum}(p_h, p_\ell, f_h, f_\ell)$
 - 11: $\epsilon = \circ(\circ(zt), 0 \times 1 . 81 \text{p}-52)$
 - 12: $\ell_b = h + \circ(\ell - \epsilon)$, $u_b = h + \circ(\ell + \epsilon)$
 - 13: if $\circ(\ell_b) = \circ(u_b)$ return $\circ(\ell_b)$ else return FAIL
-

REMARK: in line 11 of `FastPathSmall`, ϵ is negative for $x > 0$. In this case, ℓ_b is an upper bound and u_b a lower bound, but the rounding test remains correct. In the following we assume ϵ represents $|\epsilon|$.

Algorithm `fasttwosum` is the same as in §I, and `fastsum` and `muldd` are given in Algorithms 4 and 5.

Algorithm 4 (fastsum)

Input: $a := a_h + a_\ell, b := b_h + b_\ell$ double-double numbers

Output: h, ℓ such that $h + \ell$ approximates $a + b$

- 1: $h, s \leftarrow \text{fasttwosum}(a_h, b_h)$
 - 2: $u \leftarrow \circ(a_\ell + b_\ell)$
 - 3: $\ell \leftarrow \circ(u + s)$
-

Algorithm 5 (muldd)

Input: x_h, x_ℓ, c_h, c_ℓ binary64 numbers

Output: h, ℓ approximating $(x_h + x_\ell)(c_h + c_\ell)$

- 1: $h \leftarrow \circ(x_h c_h)$
 - 2: $\ell_1 \leftarrow \text{fma}(x_h, c_h, -h)$
 - 3: $\ell_2 \leftarrow \circ(x_h c_\ell)$
 - 4: $\ell_3 \leftarrow \circ(x_\ell c_h)$
 - 5: $\ell \leftarrow \circ(\circ(\ell_2 + \ell_3) + \ell_1)$
-

Lemma 3: For x a binary64 number, $2^{-15} \leq |x| < 2^{-4}$, and any rounding mode, if Algorithm `FastPathSmall` does not return FAIL, then it yields the correct rounding of $\text{acos}(x)$.

Proof: In $[2^{-15}, 2^{-4}]$, we can prove with Sollya, for $y = p_h + p_\ell - q(x)$ (see Figure 2):

$$10219 \cdot 2^{-70} \leq (y - \text{acos}(x))/x^3 \leq 10985 \cdot 2^{-70}. \quad (4)$$

We then use the following Gappa script, with Gappa option `-Echange-threshold=0`, with `RND` substituted by all rounding modes, `EPS` substituted by `0 \times 1 . 81 \text{p}-52`, `XMIN`

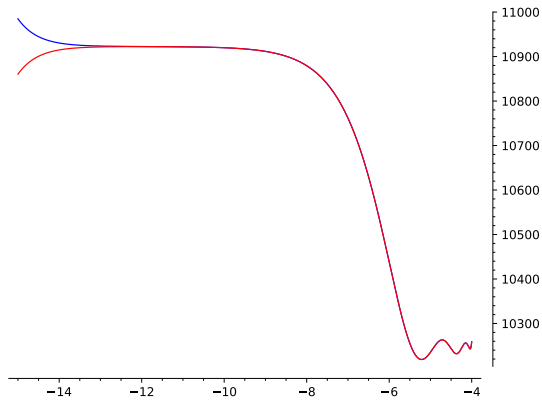


Fig. 2. Graph of $((p_h + p_l - q(x)) - \text{acos}(x))/x^3$ scaled by 2^{70} , for $2^{-15} \leq |x| \leq 2^{-4}$, where the x -axis gives $\log_2 |x|$. The blue curve is for $x > 0$, the red curve for $x < 0$.

and XMAX substituted by 2^{-15} and 2^{-4} for positive x , by -2^{-4} and -2^{-15} for negative x , and BND substituted by the upper bound $10985 \cdot 2^{-70}$:

```
@rnd = float<ieee_64,RND>;
ph = 0x1.921fb54442d18p+0;
pl = 0x1.1a62633145c07p-54;
p = ph+pl; # approximates pi/2
x = rnd(x_);
x2 = x*x;
t rnd= x*x;
z = -x;
c0 = 1;
c1 = 0;
c2 = 0x1.55555555555555p-3;
c3 = 0x1.33333333333333e4p-4;
c4 = 0x1.6db6db6d31f82p-5;
c5 = 0x1.f1c71f6889397p-6;
c6 = 0x1.6e874b7045b46p-6;
c7 = 0x1.1f753132271e2p-6;
x4 = x2*x2;
t2 rnd= t*t;
g2 = c2+x2*c3;
g2r rnd= c2+t*c3;
h4 = c4+x2*c5;
h4r rnd= c4+t*c5;
h6 = c6+x2*c7;
h6r rnd= c6+t*c7;
g4 = h4+x4*h6;
g4r rnd= h4r+t2*h6r;
e2 = g2+x4*g4;
e2r rnd= g2r+t2*g4r;
d = x2*e2;
dr rnd= t*e2r;
fl = c1+d;
flr rnd= c1+dr;
# fh1, fl1 = muldd(z, 0, c0, flr)
q = z*(c0+fl);
```

```
fh1 rnd= z*c0;
l1 = -(fh1-(z*c0));
l2 rnd= z*flr;
# l3=0
fl1 rnd= l2+l1;
# fh1+fl1 approximates q
# h, l = fastsum(ph, pl, fh1, fl1)
# h, s = fasttwosum(ph, fh1)
h rnd= ph+fh1;
s = rnd(-(h-(ph+fh1)));
u rnd= pl+fl1;
l rnd= u+s;
eps rnd= |z*t|*EPS;
lb = h + rnd(l - eps);
# h+l approximates ph+pl+q
y = ph+pl+q;
x3 = x2*x;
```

```
{ x in [XMIN,XMAX] -> fh1 // ph in [-1,1]
  /\ (lb-y)+x3*BND <= 0 }
```

```
fh1+l1 -> z*c0;
(fh1+fl1)-q -> ((fh1+l1) - z*c0)
  + (fl1 - (l2+l1)) + (l2 - z*flr) + z*(flr-fl);
(h+s) - (ph+fh1) -> s - (-(h-(ph+fh1)));
(h+l)-(p+(fh1+fl1)) -> ((h+s) - (ph+fh1))
  + (u - (pl+fl1)) + (l - (u+s));
```

This script proves:

$$\frac{\ell_b - y}{x^3} + 10985 \cdot 2^{-70} \leq 0.$$

Combined with Eq. (4), this proves $\ell_b \leq \text{acos}(x)$.

A similar Gappa script which differs only by the following lines proves $\text{acos}(x) \leq u_b$:

```
ub = h + rnd(l + eps);
{ x in [XMIN,XMAX] -> fh1 // ph in [-1,1]
  /\ (ub-y)+x3*BND >= 0 }
```

We thus have $\ell_b \leq \text{acos}(x) \leq u_b$, whence since rounding is monotonic $\circ(\ell_b) \circ(\text{acos}(x)) \leq \circ(u_b)$, if $\circ(\ell_b) = \circ(u_b)$, they equal the correct rounding of $\text{acos}(x)$. ■

If we reduce the error bound $0x1.81p-52$ to $0x1.d3p-53$, the fast path fails for $x = 0x1.7cb54339263fbbp-12$ without FMA. On the other hand, the largest error bound (as a 9-bit number) we found that defeats the rounding test is $0x1.00p-52$ with $x = 0x1.d12b3716d66e4p-7$. The smallest error bound we could prove with Gappa for $2^{-15} \leq |x| < 2^{-4}$ is $0x1.5fp-52$, but it took about 4 hours on our computer, and this branch of the code shares the same error bound than $2^{-4} \leq |x| < 0.5$, for which $0x1.81p-52$ was proven to be enough by exhaustive tests (§III). It might be that an error bound smaller than $0x1.81p-52$ works too, while defeating the condition $\ell_b \leq \text{acos}(x) \leq u_b$ of the rounding test. Indeed, it can be that cases where this condition is not satisfied are

far from rounding boundaries, and thus $\circ(\ell_b)$ is still the correct rounding. However, finding the smallest “working” error bound is only possible with exhaustive search, which would be quite expensive.

A. Accurate path

The accurate path is analyzed in §IV.

III. CASE $2^{-4} \leq |x| \leq 1$

For $2^{-4} \leq |x| \leq 1$, we performed an exhaustive search for all rounding modes, with and without FMA. This exhaustive search revealed that a previous error bound was too small.

IV. ACCURATE PATH

In this section we prove the correctness of the accurate path (routine `as_acos_refine` in the code) for $|x| < 0.5$. In theory, we only have to prove it for $|x| < 2^{-4}$ since for $2^{-4} \leq |x| < 1$ we performed an exhaustive search (§III), which checks the accurate path is correct when it is called. Moreover we know from Lemma 2 that the accurate path is never called for $|x| \leq x_0 := 0 \times 1. \text{cb3b3869747f4p-55}$, thus we assume $x_0 < |x| < 0.5$ in this section.

The accurate path takes as input the binary64 number x and the approximation ϕ of $\text{acos}(x)$ computed by `FastPathSmall` ($\phi = \ell_b$). From ϕ , an integer $j \in [0, 32]$ is deduced by $j \approx \lfloor 64 \cdot |\phi - \pi/2|/\pi \rfloor$. For our domain of interest $|x| < 0.5$, we have $0 \leq j \leq 11$. Since $\phi \leq \pi/2$ for $x > 0$ and $\phi \geq \pi/2$ for $x < 0$, we have:

$$\text{acos}(x) = \frac{\pi}{2} \mp j \frac{\pi}{64} - \delta, \quad (5)$$

where \mp is a minus sign for $x > 0$, a plus sign for $x < 0$, and $|\delta|$ is bounded by $\approx \pi/128$. Indeed, due to approximations and rounding errors, $|\delta|$ might slightly exceed $\pi/128$. For a given j , sign of x , and rounding mode, $|\delta|$ is maximal at the boundaries of the range of x with given sign yielding this value of j . We found that at all these boundaries, $|\delta| < 0.0246$. Eq. (5) translates to:

$$\delta = \frac{\pi}{2} \mp j \frac{\pi}{64} - \text{acos}(x), \quad (6)$$

and to:

$$\sin(\delta) = \cos(j \frac{\pi}{64})x \mp \sin(j \frac{\pi}{64})\sqrt{1-x^2}. \quad (7)$$

Double-double approximations of $\cos(j \frac{\pi}{64})$ and $\sin(j \frac{\pi}{64})$ are tabulated, which yield double-double approximations of $\sqrt{1-x^2}$ and $\sin(\delta)$. Using a polynomial approximation of $\text{asin}(t)$, we deduce a double-double approximation of δ . Then it suffices to plug this approximation into Eq. (5). Five hard-to-round cases are not rounded correctly by this algorithm, and need to be dealt with separately.

For simplicity, we omit roundings in lines 20 and 26 of Algorithm `AccuratePath`. At lines 15-16, we split $S_h d_{ch}$ into $S_c + d_{S_c}$ (up to rounding), where S_c is a multiple of $\text{ulp}(M) = 2^{-56}$, and $|d_{S_c}| < 2^{-56}$. The same holds for $C_h d_{sh}$ which is split into $C_s + d_{C_s}$ at lines 17-18. This ensures that the subtraction $C_s - S_c$ at line 19 is exact.

Algorithm 6 (`AccuratePath`)

Input: x, ϕ binary64 numbers, with $\phi \approx \text{acos}(x)$

Output: $p_h + p_\ell$ approximating $\text{acos}(x)$

- 1: $s_2 = \circ(x^2)$, $d_2 = \text{fma}(x, x, -s_2)$
 - 2: $c_{2h}, c_{2\ell} = \text{fasttwosum}(1, -s_2)$
 - 3: $c_{2\ell} = \circ(c_{2\ell} - d_2)$
 - 4: $c_{2h}, c_{2\ell} = \text{fasttwosum}(c_{2h}, c_{2\ell})$
 - 5: $c_h = \circ(\sqrt{c_{2h}})$
 - 6: $c_\ell = \circ(\circ(c_{2\ell} - \text{fma}(c_h, c_h, -c_{2h})) \cdot \circ(0.5/c_h))$
 - 7: $P = 0 \times 1.921 \text{fb}54442 \text{d}18 \text{p}+0$ ▷ $P \approx \pi/2$
 - 8: $I = 0 \times 1.45 \text{f}306 \text{dc}9 \text{c}883 \text{p}+4$ ▷ $I \approx 64/\pi$
 - 9: $j = \lfloor \circ(|\phi - P|) \cdot I \rfloor$
 - 10: $C_h + C_\ell \approx \cos(j\pi/64)$ ▷ tabulated
 - 11: $S_h + S_\ell \approx \sin(j\pi/64)$ ▷ tabulated
 - 12: $d_{sh} = \circ(|x| - S_h)$, $d_{s\ell} = -S_\ell$
 - 13: $d_{ch} = \circ(c_h - C_h)$, $d_{c\ell} = \circ(c_\ell - C_\ell)$
 - 14: $M = 0 \times 1.8 \text{p}-4$
 - 15: $S_c = \circ(\text{fma}(S_h, d_{ch}, M) - M)$
 - 16: $d_{S_c} = \text{fma}(S_h, d_{ch}, -S_c)$
 - 17: $C_s = \circ(\text{fma}(C_h, d_{sh}, M) - M)$
 - 18: $d_{C_s} = \text{fma}(C_h, d_{sh}, -C_s)$
 - 19: $v = \circ(C_s - S_c)$
 - 20: $d_v = (C_h d_{s\ell} + C_\ell d_{sh}) - (S_h d_{c\ell} + S_\ell d_{ch}) - (d_{S_c} - d_{C_s})$
 - 21: $v, d_v = \text{fasttwosum}(v, d_v)$
 - 22: $j' = 32 - j$ if $x > 0$, $j' = 32 + j$ if $x < 0$
 - 23: $v_2, d_{v2} = \text{muldd}(v, d_v, v, d_v)$
 - 24: $v = -\text{sign}(x)v$, $d_v = -\text{sign}(x)d_v$
 - 25: let $c_0 t + c_1 t^3 + \dots + c_7 t^{15}$ approximating $\text{asin}(x)$ where $c_i = c_{ih} + c_{i\ell}$ for $i \leq 4$ ▷ tabulated
 - 26: $r = v_2(c_5 + v_2(c_6 + v_2 c_7))$
 - 27: $f_h, f_\ell = \text{polydd}(v_2, d_{v2}, 5, r)$
 - 28: $f_h, f_\ell = \text{muldd}(v, d_v, f_h, f_\ell)$
 - 29: $p_h = \circ(0 \times 1.921 \text{fb}54442 \text{dp}-5 \cdot j')$
 - 30: $p_\ell = \circ(0 \times 1.8469898 \text{cc}518 \text{p}-53 \cdot j')$
 - 31: $p_s = \circ(-0 \times 1. \text{fc}8 \text{f}8 \text{cbb}5 \text{bf}6 \text{cp}-102 \cdot j')$
 - 32: $p_\ell, p_s = \text{Sum}(f_h, f_\ell, p_\ell, p_s)$
 - 33: $p_h, p_\ell = \text{fasttwosum}(p_h, p_\ell)$
 - 34: $p_\ell = \circ(p_\ell + p_s)$
-

Algorithm `TwoSum` is Algorithm 4.4 from [4]: From Theorem 4.1 from [1], Algorithm `TwoSum` satisfies, if all the rounding are faithful:

$$|s + t - (a + b)| \leq 2^{-52} \text{ulp}(s). \quad (8)$$

We split the analysis of Algorithm `AccuratePath` in three parts: in §IV-A we analyze the approximation of $\sqrt{1-x^2}$, in §IV-B we analyze the approximation of $\sin(\delta)$, from which in §IV-C we deduce an approximation of δ and thus of $\text{acos}(x)$ using Eq. (5).

A. Approximation of $\sqrt{1-x^2}$

The double-double approximation $c_h + c_\ell$ of $\sqrt{1-x^2}$ is computed at lines 1-6.

Algorithm 7 (polydd)

Input: $x = x_h + x_\ell$ double-double, n integer, a double r **Output:** $h + \ell$ approximating $c_0 + c_1 \cdot x^2 + \dots + (c_{n-1} + r)x^{2n-2}$

- 1: $h, t \leftarrow \text{fasttwosum}(c_{n-1}, h, r)$
 - 2: $\ell \leftarrow \circ(t + c_{n-1}, \ell)$
 - 3: for i from $n - 2$ downto 0
 - 4: $h, \ell \leftarrow \text{muldd}(x_h, x_\ell, h, \ell)$
 - 5: $h, e \leftarrow \text{fasttwosum}(c_{i,h}, h)$
 - 6: $\ell \leftarrow \circ(\ell + c_{i,\ell}) + e$
-

Algorithm 8 (Sum)

Input: $x_h + x_\ell$ and $c_h + c_\ell$ two double-double pairs**Output:** $s_h + s_\ell$ approximating $(x_h + x_\ell) + (c_h + c_\ell)$

- 1: $s_h, \ell \leftarrow \text{TwoSum}(x_h, c_h)$
 - 2: $s_\ell \leftarrow \circ(\ell + c_\ell) + \ell$
-

Lemma 4: For $|t| \leq 2^{-7}$:

$$|\sqrt{1+t} - (1 + \frac{t}{2})| < 0.126t^2.$$

Proof: We can see graphically that the function $f(t) = (\sqrt{1+t} - (1 + t/2))/t^2$ is negative and increasing on $[-2^{-7}, 2^{-7}]$, thus reaches its maximum in absolute value at $t = -2^{-7}$, for which it is about 0.125491. ■

Lemma 5: If $|\ell/h| \leq 2^{-7}$, then:

$$\left| \sqrt{h} + \frac{\ell}{2\sqrt{h}} - \sqrt{h+\ell} \right| < 0.126 \frac{\ell^2}{h^2} \sqrt{h}.$$

Proof: Apply Lemma 4 with $t = \ell/h$. ■

We apply Lemma 5 with $h = c_h^2$ and $\ell = c_{2h} + c_{2\ell} - h$, with $c_{2h}, c_{2\ell}$ the values computed at line 4. Our Gappa script proves $|\ell/h| \leq 2^{-50.415}$, so the condition of the lemma is fulfilled, thus we deduce from Lemma 5 and the fact that $c_h \leq 1$:

$$\left| c_h + \frac{\ell}{2c_h} - \sqrt{c_{2h} + c_{2\ell}} \right| < 2^{-103.818}. \quad (9)$$

At line 6, the lower term $c_{2\ell}$ is added to $-\text{fma}(c_h, c_h, -c_{2h})$ which takes into account the rounding error in $c_h = \circ(\sqrt{c_{2h}})$. The Gappa script proves:

$$\left| c_\ell - \frac{\ell}{2c_h} \right| < 2^{-102.677}. \quad (10)$$

Combining Eq. (9) and Eq. (10) we obtain for the values of $c_{2h}, c_{2\ell}$ at line 4, and those of c_h, c_ℓ at lines 5-6:

$$|c_h + c_\ell - \sqrt{c_{2h} + c_{2\ell}}| < 2^{-102.137}. \quad (11)$$

Now we have to bound $|c_h + c_\ell - \sqrt{1-x^2}|$:

$$\begin{aligned} \left| c_h + c_\ell - \sqrt{1-x^2} \right| &\leq \left| c_h + c_\ell - \sqrt{c_{2h} + c_{2\ell}} \right| \\ &\quad + \left| \sqrt{c_{2h} + c_{2\ell}} - \sqrt{1-x^2} \right|. \end{aligned} \quad (12)$$

Gappa proves that after line 4, $|c_{2h} + c_{2\ell} - (1-x^2)| \leq 2^{-104}$. Let $\alpha = \sqrt{c_{2h} + c_{2\ell}}$ and $\beta = \sqrt{1-x^2}$. We thus have $|\alpha^2 - \beta^2| \leq 2^{-104}$, thus $|\alpha - \beta| \leq 2^{-104}/(\alpha + \beta)$. Gappa proves

Algorithm 9 (TwoSum)

Input: a and b binary64 numbers**Output:** $s + t$ approximating $a + b$

- 1: $s \leftarrow \circ(a + b)$
 - 2: $a' \leftarrow \circ(s - b)$
 - 3: $b' \leftarrow \circ(s - a')$
 - 4: $\delta_a \leftarrow \circ(a - a')$
 - 5: $\delta_b \leftarrow \circ(b - b')$
 - 6: $t \leftarrow \circ(\delta_a + \delta_b)$
-

$c_{2h} + c_{2\ell} \geq 0.75$, and since $|x| \leq 0.5$, $1 - x^2 \geq 0.75$ too, thus $\alpha, \beta \geq \sqrt{0.75} \geq 0.866$, and $|\alpha - \beta| \leq 2^{-104}/(2 \cdot 0.866) < 2^{-104.792}$. Plugging into Eq. (12) and using Eq. (11), we obtain:

$$\left| c_h + c_\ell - \sqrt{1-x^2} \right| < 2^{-101.924}. \quad (13)$$

B. Approximation of $\sin(\delta)$

A double-double approximation $v + d_v$ of $\sin(\delta)$ is computed at lines 7-21. As seen above, for $|x| < 0.5$, the integer j at line 9 satisfies $0 \leq j \leq 11$. The analysis depends on the precomputed constants C_h, C_ℓ, S_h, S_ℓ at lines 10 and 11. We analyze separately the case $j = 0$ since Gappa does not like variables which are zero.

Case $j = 0$. This case may happen when $x_0 < |x| \leq 0 \times 1.92155f7a366b7p-6$. We then have $C_h = 1$, and $C_\ell = S_h = S_\ell = 0$. It follows $d_{sh} = |x|$, $d_{s\ell} = 0$, Gappa proves $0.999 < c_h < 1.001$ thus $d_{ch} = c_h - C_h$ is exact, $d_{c\ell} = c_\ell$, $S_c = d_{S_c} = 0$, $C_s = \circ(\circ(|x| + M) - M)$ is the upper part of $|x|$ (rounded to a multiple of 2^{-56}), d_{C_s} is the lower part of $|x|$ such that $|x| = C_s + d_{C_s}$, $v = C_s$, $d_v = d_{C_s}$, thus $v + d_v = |x|$. We thus have exactly $\sin|\delta| = v + d_v$ in Eq. (7).

Let $V = (C_h + C_\ell)|x| - (S_h + S_\ell)(c_h + c_\ell)$, and $\delta' = \text{sign}(x)\delta$. For $1 \leq j \leq 11$, Gappa finds the bound (obtained for $j = 11$):

$$|v + d_v - V| \leq 2^{-103.064}.$$

The approximations $C_h + C_\ell$ of $\cos(j\pi/64)$ and $S_h + S_\ell$ of $\sin(j\pi/64)$ have an absolute error bounded by $2^{-108.574}$. It follows:

$$\begin{aligned} |V - \sin(\delta')| &\leq |C_h + C_\ell - \cos(j\pi/64)| \cdot |x| \\ &\quad + |S_h + S_\ell| \cdot |c_h + c_\ell - \sqrt{1-x^2}| \\ &\quad + |S_h + S_\ell - \sin(j\pi/64)| \cdot \sqrt{1-x^2}. \end{aligned}$$

Since $|x| < 0.5$, the first term above is bounded by $2^{-109.574}$. Since $|S_h + S_\ell| \leq 1$, and $|c_h + c_\ell - \sqrt{1-x^2}| < 2^{-101.924}$ by Eq. (13), their product is bounded by $2^{-101.924}$. The last term is bounded by $2^{-108.574}$, thus:

$$|V - \sin(\delta')| \leq 2^{-109.574} + 2^{-101.924} + 2^{-108.574} < 2^{-101.902}.$$

Combined with the above bound on $|v + d_v - V|$, this yields:

$$|v + d_v - \sin(\delta')| < 2^{-103.064} + 2^{-101.902} < 2^{-101.369}. \quad (14)$$

C. Approximation of δ

Lines 22-28 compute a double-double approximation $f_h + f_\ell$ of δ , from the approximation $v + d_v$ of $\sin(\delta')$. First, since $|\delta| < 0.0246$ and with Eq.(14), we have $|v + d_v| < 0.0247$, thus after line 21, we have $|v| < 0.025$ and $|d_v| \leq 2^{-58}$.

Let us denote $q(t) = c_0 t + \dots + c_7 t^{15}$ the polynomial approximation of $\text{asin}(t)$ used in Algorithm AccuratePath. After the muldd call at line 28, Gappa proves for $t = v + d_v$:

$$|f_h + f_\ell - q(t)| < 2^{-106.759}.$$

We now have:

$$|f_h + f_\ell - \delta| \leq |f_h + f_\ell - q(t)| + |q(t) - \text{asin}(t)| + |\text{asin}(t) - \delta|.$$

The first term in the right-hand side is bounded by $2^{-106.759}$ from just above. Sollya bounds the second term over $[-0.025, 0.025]$ by $|q(t) - \text{asin}(t)| < 2^{-108.871}$. Since $|t - \sin(\delta)| < 2^{-101.369}$, the third term is bounded by $2^{-101.369}$ multiplied by a bound for the derivative of $\text{asin}(t)$ over $[-0.025, 0.025]$. This derivative is maximal at $t = \pm 0.025$, where it is less than 1.0004, this yields a bound of $2^{-101.369} \cdot 1.0004 < 2^{-101.368}$ for the third term. Combining, we get:

$$|f_h + f_\ell - \delta| \leq 2^{-106.759} + 2^{-108.871} + 2^{-101.368} < 2^{-101.326}.$$

At lines 29-31, we multiply the integer j' , which lies in $[0, 64]$ (remember $0 \leq j \leq 32$) by a triple-word approximation of $\pi/64$, with error less than 2^{-155} . Since $j \leq 11$ we have $21 \leq j' \leq 43$. The constants involved in p_h and p_ℓ are such that the multiplication by j' is exact. Since $|j'| \leq 64$, we have $|p_s| < 2^{-95}$, thus the rounding error on p_s is bounded by $\text{ulp}(2^{-96}) = 2^{-148}$. The error from the approximation of $\pi/64$ multiplied by $j' \leq 2^6$ is bounded by 2^{-149} , thus after line 31:

$$|p_h + p_\ell + p_s - j' \frac{\pi}{4}| < 2^{-148} + 2^{-149} < 2^{-147.415}.$$

Gappa proves that the rounding error of calling Algorithm Sum is bounded by $2^{-105.752}$, thus after line 32:

$$\begin{aligned} & |p_h + p_\ell + p_s - \text{acos}(x)| \\ & < 2^{-101.326} + 2^{-147.415} + 2^{-105.752} < 2^{-101.260}. \end{aligned}$$

On line 33, Gappa proves that the precondition $|p_\ell| \leq |p_h|$ of `fasttwosum` is satisfied, and the rounding error due to `fasttwosum` is bounded by 2^{-104} . Finally, on line 34, Gappa proves the rounding error is bounded by 2^{-103} . We thus have for the final values p_h and p_ℓ :

$$|p_h + p_\ell - \text{acos}(x)| \leq 2^{-101.260} + 2^{-104} + 2^{-103} < 2^{-100.724}. \quad (15)$$

Theorem 1: For any binary64 number x , $|x| < 0.5$, such that $\text{acos}(x)$ has less than 46 identical bits after the round bit, if $p_h + p_\ell$ is the double-word approximation of $\text{acos}(x)$ computed by Algorithm AccuratePath, then $\circ(p_h + p_\ell)$ is the correct rounding of $\text{acos}(x)$.

Proof: Let $y = \text{acos}(x)$. If y has less than 46 identical bits after the round bit, then y is at distance at least $2^{-47} \text{ulp}(y)$ from a rounding boundary z : $|y - z| \geq 2^{-47} \text{ulp}(y)$. Since

$\text{ulp}(y) > 2^{-53} |y|$, this yields $|y - z| \geq 2^{-100} |y| \geq 2^{-100}$, since for $|x| < 0.5$, $\text{acos}(x) \geq 1$. Thus Eq. (15) shows that $|p_h + p_\ell - y| < |y - z|$, i.e., $p_h + p_\ell$ is closer from $\text{acos}(x)$ than any rounding boundary. It follows $\circ(p_h + p_\ell) = \circ(\text{acos}(x))$. ■

V. HARD-TO-ROUND CASES

During this proof, we realized that hard-to-round cases for $x < 0$ were missing in the CORE-MATH test suite. Indeed, since $\text{acos}(-x)$ differs from $\pm \text{acos}(x)$, it does not suffice to compute hard-to-round cases for $x > 0$.

REFERENCES

- [1] BOLDO, S., GRAILLAT, S., AND MULLER, J. On the robustness of the `2sum` and `fast2sum` algorithms. *ACM Trans. Math. Softw.* 44, 1 (2017), 4:1–4:14.
- [2] CHEVILLARD, S., JOLDEŞ, M., AND LAUTER, C. Sollya: An environment for the development of numerical codes. In *Mathematical Software - ICMS 2010* (Heidelberg, Germany, September 2010), K. Fukuda, J. van der Hoeven, M. Joswig, and N. Takayama, Eds., vol. 6327 of *Lecture Notes in Computer Science*, Springer, pp. 28–31.
- [3] JEANNEROD, C.-P., AND ZIMMERMANN, P. FastTwoSum revisited. In *32nd IEEE Symposium on Computer Arithmetic, ARITH 2025, El Paso, TX, USA, May 4-7 (2025)*.
- [4] MULLER, J.-M., BRUNIE, N., DE DINECHIN, F., JEANNEROD, C.-P., JOLDES, M., LEFÈVRE, V., MELQUIOND, G., REVOL, N., AND TORRES, S. *Handbook of Floating-point Arithmetic (2nd edition)*. Birkhäuser Basel, July 2018.