



HAL
open science

Correct Rounding in Double Extended Precision

Sélène Corbineau, Paul Zimmermann

► **To cite this version:**

Sélène Corbineau, Paul Zimmermann. Correct Rounding in Double Extended Precision. 2025. hal-04861251

HAL Id: hal-04861251

<https://inria.hal.science/hal-04861251v1>

Preprint submitted on 2 Jan 2025

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

Correct Rounding in Double Extended Precision

S el ene Corbineau
D epartement d'Informatique de l'ENS
 cole Normale Sup erieure-PSL
F-75005 Paris, France
selene.corbineau@ens.fr

Paul Zimmermann
Universit e de Lorraine
CNRS, Inria, LORIA
F-54000 Nancy, France
paul.zimmermann@inria.fr

Abstract—The double extended precision format is an 80-bit floating-point format introduced in the 80x87 series of floating-point processors by Intel. Since the introduction of vector instructions in the x86 processors, its use has fallen due to speed concerns. We implement the first correctly-rounded routines for double extended precision. These implementations use modern microprocessor features and double-double arithmetic, avoiding x87-specific features, and achieve up to 2x speedup over state-of-the-art implementations which are not correctly rounded. This demonstrates that double extended precision could be viable as a large computational format.

Index Terms—IEEE 754, floating-point, correct rounding, double extended format, efficiency.

I. INTRODUCTION

The double extended precision format was introduced in Intel's x87 series of floating-point coprocessors. Featuring excess precision and range compared to the double precision format (now binary64), it allowed double precision intermediate computations to be carried out without excessive accumulation of errors. Except for some minor encoding quirks, it can be thought of as the natural extension of the IEEE 754 standard's binary floating-point formats to 80 bits. Most C compilers for the x86 and x86-64 architectures assign double extended precision to the `long double` type, the most precise floating-point type described by the C standards. Intel later introduced vector instructions in the x86 architecture to assist multimedia and signal processing applications. These SIMD extensions offer at most double precision, but are faster and easier to compile for. They proved popular and usage of double extended precision dwindled.

Current `long double` mathematical functions implementations perform all their calculations in the slow native double extended format. Most are a thin wrapper around x87 complex mathematical instructions such as `FPSIN` or `FPEXP2` which are microcoded, unmaintained, and yield processor-dependent results. These instructions are slow even compared to naive computation in the double extended format and their accuracy has been the object of errata.

A solution against the reproducibility problem is to enforce usage of correctly-rounded routines. However, prior to this work, the major correctly-rounded libraries (CRLIBM, MathLib, LLVM-libm, RLIBM, CORE-MATH) did not offer any extended double precision support. This is partly because correctly-rounded routines must compute their result to a higher intermediate precision. Except for the single precision

format, where using double precision for internal computations is quite efficient, this is usually achieved by double-word arithmetic. This technique gives best performance when fused multiply-add (FMA) is available in hardware; however since this is not the case for extended precision on current processors, it would yield unacceptably slow routines.

Leveraging the fact that double-double arithmetic is more precise than the double extended format, we implement correctly-rounded routines using double-double arithmetic in the critical path. This nearly eliminates x87 usage, allowing a significant performance boost and could be ported to allow double extended precision computation on non-x86 architectures.

Another difficulty implementing correctly-rounded routines is that certifying them traditionally involves looking for the hardest-to-round cases. This is hard for double extended precision routines because of the large input space. In this work, we show that this search can be replaced by a proof that all cases are sufficiently easy. This allows us to use more restrictive parameters in our search tools and speeds up correctness checking.

The contributions of this article are the following: after some useful routines in §II, we present in §III routines to convert a double extended number to a double-double representation, to convert a double-double representation to a 128-bit floating-point number, and to round such a 128-bit number to a double extended number. These routines are building blocks independent from the function to be implemented. Then in §IV we present a case study for the exponential function, correctly rounded in double extended precision. Finally, we present experimental results and conclude in §V.

II. BACKGROUND

A. Correct rounding

For a mathematical function f and a floating-point number x , the correct rounding is the floating-point number y which is closest to $f(x)$ according to the given rounding mode. Correct rounding is a way to ensure portability of programs and reproducibility of results. For basic arithmetic operations, IEEE 754 requires correct rounding. For mathematical functions, many people believe that correct rounding is expensive, but Ziv demonstrated already in 1991—with the MathLib library—that it can be computed efficiently [9]. For more details on correct rounding, we refer the reader to [1].

B. Double-double arithmetic

A “double-double” number is an evaluated sum $h + \ell$ of two double-precision numbers. If $h = \circ(h + \ell)$ with the rounding mode $\circ(\cdot)$, we say the double-double number is “normalized”. Normalizing double-double numbers after each operation simplifies the error analysis, but is expensive: we usually try to avoid it, still trying to keep a rigorous bound on the overlap between h and ℓ . For more details on double-double arithmetic, see [8].

We recall here three main algorithms on double-double arithmetics: Algorithm ExactMul (resp. FastTwoSum) produces a double-double number for the product (resp. sum) of two double-precision numbers (we use here the formulation from [5]), and Algorithm DMul approximates the product of two double-double numbers. In the absence of underflow and overflow, ExactMul is always exact, whatever the rounding mode. In ExactMul, we denote by $\circ(ab - h)$ the correct

Algorithm 1 (ExactMul)

Input: $a, b \in \mathbb{F}$

Output: h, ℓ such that $h + \ell = ab$

- 1: $h \leftarrow \circ(ab)$
 - 2: $\ell \leftarrow \circ(ab - h)$
-

rounding of $ab - h$, which can be performed with a fused multiply-add.

Algorithm 2 (DMul)

Input: $a_h, a_\ell, b_h, b_\ell \in \mathbb{F}$

Output: x_h, x_ℓ approximating $(a_h + a_\ell)(b_h + b_\ell)$

- 1: $(x_h, s) \leftarrow \text{ExactMul}(a_h, b_h)$
 - 2: $t \leftarrow \circ(s + a_\ell b_h)$
 - 3: $x_\ell \leftarrow \circ(t + a_h b_\ell)$
-

We will make extensive use of the following lemma.

Lemma 1: Using the notations of Algorithm 2, let $A_h, A_\ell, B_h, B_\ell \in \mathbb{R}_+$ be bounds on $|a_h|, |a_\ell|, |b_h|, |b_\ell|$ respectively. Let RU be the rounding toward infinity. Let us call

$$T = \text{RU}(\text{ulp}(\text{RU}(A_h B_h)) + A_\ell B_h), X_\ell = \text{RU}(T + A_h B_\ell).$$

Then, assuming no underflow or overflow, the error of DMUL is at most $\text{ulp}(T) + \text{ulp}(X_\ell) + A_\ell B_\ell$, and X_ℓ is such that $|x_\ell| \leq X_\ell$.

Furthermore, we also have $|x_h| \leq X_h = \text{RU}(A_h B_h)$.

Proof: We have $|x_h| \leq \text{RU}(A_h B_h)$ given EXACTMUL. Then we see that $s < \text{ulp}(x_h)$, so that $s < \text{ulp}(\text{RU}(A_h B_h))$. Given line 2 we get that $|t| \leq T$. The rounding error on line 2 is therefore at most $\text{ulp}(T)$. In the same vein, we get that $|x_\ell| \leq \text{RU}(T + A_h B_\ell) = X_\ell$ so that the rounding error on line 3 is at most $\text{ulp}(X_\ell)$. Since we neglected the term $a_\ell b_\ell$ from the whole product, we get the claimed error bound. ■

Still in the absence of underflow and overflow, FastTwoSum is exact for rounding to nearest, but might not be for directed rounding modes (see [2] for more details).

Algorithm 3 (FastTwoSum)

Input: $a, b \in \mathbb{F}$ with $a = 0$ or $|a| \geq |b|$

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

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

C. The double extended format

The double extended format is an 80-bit format, with a 1-bit sign s , a 15-bit exponent e , and a 64-bit significand m . Apart from special values (NaN, Inf) and subnormal numbers, the corresponding value is $(-1)^s \cdot 2^{e-16383} \cdot m$, where m is interpreted as a number in $[1, 2)$, and always has its most significant bit set (contrary to the binary32 and binary64 formats, the leading bit is explicit). This format has 11 bits more precision than binary64, and a larger exponent range, since it can represent numbers as large as about 2^{16384} , and as small as 2^{-16445} . This format usually corresponds to the C type `long double` on x86 and x86-64 processors. In the rest of this article, since we target such processors, both “double extended” and “long double” refer to that format.

D. Notations

We denote \mathbb{F}_p the set of p -bit floating-point numbers, and DD the set of double-double numbers. For $x \in \mathbb{F}_p$, we denote by $\text{ulp}(x)$ the unit-in-last place of x . We also denote $\text{ulp}_q(x)$ the unit-in-last place of x , considered as a q -bit number (with unbounded exponent).

III. EVALUATION STRATEGY

Since our aim is to use only double-double arithmetic in the critical path, we present here routines to convert between the extended double and double-double representations. These routines are independent from the function to be evaluated.

A. Conversion from double extended to double-double

The double-double representation has 11 bits for encoding the exponent, which does not allow representing all double extended values. Very large and very small inputs thus have to be handled separately. However, this is usually the first reduction step of most algorithms computing mathematical functions anyway. Depending on the function, we discarded inputs or scaled them directly using the exponent as stored on the stack prior to the routine call. This avoids an x87 `FXSCALE` instruction, replacing it by a few integer instructions.

Assuming we reduced to a long double input in the double precision exponent range, we use the following algorithm:

Algorithm 4 Long double splitting (LDSPLIT)

Input: $x \in \mathbb{F}_{80}$

Output: $(a, b) \in \mathbb{F}_{64}^2$ such that $|b| < \text{ulp}(a)$ and $a + b = x$

- $a \leftarrow \circ_{64}(x)$
 - $b \leftarrow \circ_{64}(x - \mathbb{F}_{80}(a))$
-

Lemma 2: If $|x| \leq 2^{1024}(1 - 2^{-53})$ and $|x| \geq 2^{-1011}$ then Algorithm 4 is correct.

Proof: The condition $|x| \leq 2^{1024}(1 - 2^{-53})$ ensures there is no overflow computing a . Then we know that $|x - a| < \text{ulp}(a)$. Let m be a 's exponent and n be x 's exponent. Then $m = n$ or $m = n + 1$. This implies that $\text{ulp}(a) = 2^{64-53+m-n} \text{ulp}(x) \leq 2^{12} \text{ulp}(x)$. Moreover, since a and x are integer multiples of $\text{ulp}(x)$, so is $x - a$. Therefore $x - a = k \cdot \text{ulp}(x)$ for some $k \in [-2^{12}, 2^{12}]$. Since $|x| \geq 2^{-1011}$, $\text{ulp}(x) \geq 2^{-1074}$. Together with $|x - a| \leq |x|$ this ensures that computing b is exact, thus the correctness of the algorithm. ■

Except for a final FLD due to System V's ABI, this algorithm constitutes our only use of x87 instructions in the critical path. It compiles to a FLDT, FSTL, FSUBL, FSTPL dependency chain. In our implementations, integer reduction computations using x directly from the stack partially hide the latency.

If x87 support is not available, a variant of Algorithm 4 could be emulated in software using bit-truncation of the significand.

B. Reconstruction

Near the end of the critical path, we get a good double-double approximation of $f(x)$ as $a + b$ with $|b| < \text{ulp}(a)$. We convert (a, b) to an intermediate representation s whose significand has 128 bits (with no implicit bit), then round s to a long double, checking whether our error bounds ensure this result is correct in the process. To convert (a, b) to s , we use Algorithm 5.

Algorithm 5 Computing a 128-bit approximation s

Input: $a, b \in \mathbb{F}_{64}$, $a = (-1)^{a_s}(2^{53} + a_m)2^{a_e}$ with $a_s \in \{0, 1\}$ and $0 \leq a_m < 2^{52}$, similarly for b

Output: $s = (a_s, s_e, s_m) \in \mathbb{F}_{128}$ representing $(-1)^{a_s}2^{s_e-127} \cdot s_m$ and approximating $a + b$

- 1: $s_e \leftarrow a_e$
 - 2: $s_m \leftarrow 2^{127} + 2^{64+11}a_m + (-1)^{b_s-a_s}(2^{63} + 2^{11}b_m)2^{64+b_e-a_e}$
 - 3: **if** $s_m < 2^{127}$ **then**
 - 4: $s_e \leftarrow s_e - 1$
 - 5: $s_m \leftarrow 2s_m$
 - return** $s = (a_s, s_e, s_m)$
-

In line 2, the result is rounded to a 128-bit integer by truncation (if $-63 \leq t := 64 + b_e - a_e < 0$, we shift $2^{63} + 2^{11}b_m$ by $-t$ bits to the right, and if $t \leq -64$, the contribution of b is simply neglected). Notice that since $|b| < \text{ulp}(a)$, we do not have to handle overflows.

Lemma 3: Assume a and b are normal double precision numbers such that $|b| < \text{ulp}(a)$. Then Algorithm 5 returns s such that $s = a + b + \epsilon$ with $|\epsilon| \leq 2^{-126}|s|$.

Proof: Assume a and b have same sign ($b_s = a_s$). Then, since $|b| < \text{ulp}(a)$, by summing the (scaled) significands of a and b with infinite precision, the most significant bit set stays that of a . This remains true in our 128-bit format because we truncate the expansion. Then, the branch cannot be taken and

there cannot be an overflow computing s_m . The correctness is then straightforward.

Assume a and b have different signs ($b_s \neq a_s$). Since $|b| < \text{ulp}(a) \leq |a|/2$, $s_m \geq 2^{-1}(2^{127} + 2^{64+11}a_m) \geq 2^{-126}$. If the branch is not taken, the correctness is obvious. If the branch is taken, our previous inequality shows that we only need to shift s_e by 1 to obtain a normalized 128-bit representation. The operation in the branch does not alter $s_m \cdot 2^{s_e-127}$ which proves correctness.

The truncated part of b (if any) is less than 1, and is multiplied by 2 when the branch is taken, thus contributes to less than 2 compared to $s_m \geq 2^{127}$ at the end, which proves the bound on ϵ . ■

Note that the precondition $|b| < \text{ulp}(a)$ can be ensured with a prior FastTwoSum. If we relax the precondition $|b| < \text{ulp}(a)$ to say $|b| < |a|/2$, we have to handle potential overflows in line 2. Computing whether to take the (rarely taken) branch ended up being slower than the version presented here.

Once we obtain s , we round it to a long double according to the current rounding mode. We also need to get the distance to the nearest rounding boundary. We do the following, assuming s is not in the subnormal range:

Algorithm 6 Rounding a 128-bit number s to long double

- 1: write $s_m = m_h 2^{64} + m_\ell$ with $0 \leq m_h, m_\ell < 2^{64}$
 - 2: **if** $r = \text{NEAREST}$ **then**
 - 3: $m_\ell \leftarrow m_\ell + 2^{63} \pmod{2^{64}}$
 - 4: $a \leftarrow m_\ell \text{ cmod } 2^{64}$ $\triangleright -2^{63} \leq a < 2^{63}$
 - 5: **if** $C_r(m_h, m_\ell, s_s)$ **then**
 - 6: $m_h \leftarrow m_h + 1$
 - 7: **if** $m_h = 2^{64}$ **then**
 - 8: $s_e \leftarrow s_e + 1$
 - 9: $m_h \leftarrow 2^{63}$
 - 10: $a \leftarrow a/2$ \triangleright rounded towards zero
 - 11: $y \leftarrow \mathbb{F}_{64}(s_s, s_e, m_h)$
 - 12: **if** $s_e \geq 16384$ **then**
 - 13: **return** $((-1)^{s_s} \infty, a)$
 - 14: **return** (y, a)
-

Here C_r is a predicate deciding whether (m_h, m_ℓ, s_s) has to be rounded away from zero, and line 4 reinterprets the bits of m_ℓ .

Lemma 4: Assuming $s_e \geq -16383$ at the beginning of Algorithm 6, let (y, a) be its output. Let also $u = a \cdot 2^{-64} \text{ulp}_{64}(y)$. Then y is the rounding of s in the current rounding mode and u is the signed distance between s and the nearest rounding boundary, rounded toward zero.

Proof: The rounding procedure is mostly straightforward, given that we exclude the denormal range. It remains to show that u is what is claimed. First, assume that there is no overflow in line 6. We show that u is exactly the signed distance between s and the nearest rounding boundary.

Assume r is a directed rounding mode. Then when s takes all possible values between two adjacent 64-bit precision numbers z^- and z^+ , m_h is fixed and m_ℓ linearly increases

between 0 and $2^{64} - 1$. When considered as a signed number, a goes between 0 and $2^{63} - 1$ when $z^- \leq s < (z^+ + z^-)/2$ and between -2^{63} and -1 when $(z^+ + z^-)/2 \leq s < z^+$. It is then easy to check that u is appropriately scaled such that it is as claimed.

If r is rounding to nearest, take z^- and z^+ as before. Then, when $z^- \leq s < (z^- + z^+)/2$ we see that a goes from -2^{63} to 0. Similarly, when $(z^- + z^+)/2 < s < z^+$ then a goes from 1 to $2^{63} - 1$. A similar argument to above allows us to finish the proof when there is no overflow.

The same argument can be adapted to the case where there is an overflow in line 6: the division by 2 of a explains the rounding toward zero in the theorem. ■

Algorithm 6 can be augmented to support denormal output, by shifting the significand of s to the left as needed.

The reconstruction phase (mainly Algorithm 6) is a bottleneck for two reasons:

- The branch of line 5 of Algorithm 6 is not easily predictable. To switch C_r and determine whether $r = \text{NEAREST}$, we use floating-point operations instead of directly checking the relevant CSRs.
- The System V ABI imposes that long double return values are passed in a x87 register. Since our computations are done on the stack, this imposes an `FLDT` instruction which is the hottest instruction of our routines.

IV. CASE STUDY: `EXPL`

We implemented several functions with correct rounding using the above mentioned techniques. We illustrate here the usage of our techniques on the example of the `expl` function (exponential in double extended format).

We use a truncated Ziv iteration method, with two precision levels: a fast path using double-double arithmetic and an accurate path using a 192-bit software ad-hoc format using 64-bit integer arithmetic. The hard-to-round cases were searched using `BaCSel`.

A. Fast path

The overall evaluation strategy is that of Markstein [7]. The input is first filtered: inputs with $|x| \geq 2^{14}$ round either to infinity, to the largest long double number, to 0, or to the smallest positive subnormal number depending on rounding mode and sign. By directly checking the exponent's representation, this test also catches infinities and NaNs in the input and correctly deals with them. Similarly, inputs with $|x| < 2^{-64}$ round trivially.

After this filter, all inputs easily fit in the double exponent range. We use Algorithm 4 to split the input x into (x_h, x_ℓ) exactly. All further computations proceed with double-double arithmetic using the following pseudocode.

Here, `SPLIT` is a subroutine such that `SPLIT(x')` returns $(n, f, y) \in \mathbb{Z}^2 \times \text{DD}$ satisfying $0 \leq f < 2^{20}$ and such that $n + f/2^{20} + y$ approximates x' and $|y| \leq 2^{-19.999}$. Subroutine `LOOKUP` computes an approximation of $2^{f/2^{20}}$. Subroutine `POLY` is a straightforward polynomial approximation of 2^y , knowing $|y| \leq 2^{-19.999}$. Finally, `RECONSTRUCT` is a small

Algorithm 7 Computing `expl`

Input: $x \in \mathbb{F}_{64}, 2^{-64} \leq |x| < 2^{14}$
 $x' \leftarrow \text{REDUCTION}(x)$
 $(n, f, y) \leftarrow \text{SPLIT}(x')$
 $z \leftarrow \text{LOOKUP}(f)$
 $t \leftarrow \text{POLY}(y)$
 $s \leftarrow \text{DMUL}(t, z)$
return `RECONSTRUCT`(n, s)

modification of the reconstruction phase (§III-B) to correctly round $2^n \cdot s$ to a long double. It also implements the fast path validity check knowing Algorithm 7's error bounds. Note that `LOOKUP` and `POLY` are independent and can be efficiently pipelined together.

a) *Subroutine REDUCTION:* We use the following algorithm, where constants K_h, K_ℓ are binary64 numbers.

Algorithm 8 Divide by $\log 2$ (`REDUCTION`)

Input: $x \in \mathbb{F}_{64}$ with $|x| < 2^{14}$
Output: $x' = (x'_h, x'_\ell)$ approximating $x/\log 2$
1: $(x_h, x_\ell) \leftarrow \text{LDSPLIT}(x)$
2: $K_h = 0 \times 1.71547652\text{b}82\text{fep}+0$
3: $K_\ell = 0 \times 1.777\text{d}0\text{ffda}0\text{d}24\text{p}-56$
4: $(x'_h, x'_\ell) \leftarrow \text{DMUL}(K_h, K_\ell, x_h, x_\ell)$

Lemma 5: The double-double number $x' = x'_h + x'_\ell$ returned by Algorithm 8 satisfies $|x'_h| < 2^{14.7}$, $|x'_\ell| < 2^{-37.141}$ and

$$|x' - x/\log 2| < 2^{-88.951}.$$

Moreover, if the value x'_h computed by Algorithm 8 satisfies $|x'_h| < 2^{-20}$, then $|x'_\ell| < 2^{-71.63}$ and

$$|x' - x/\log 2| < 2^{-123.3}.$$

Proof: From Algorithm 4 we know that $|x_h| \leq 2^{14}$ and $|x_\ell| < \text{ulp}(x) \leq 2^{-39}$. The constants K_h and K_ℓ satisfy $|K_h + K_\ell - 1/\log 2| < 2^{-109.5}$. Directly applying Lemma 1 (and adapting the notations), we get

$$\begin{cases} T = \text{RU}(\text{ulp}(K_h 2^{14}) + K_\ell 2^{14}) \leq 2^{-37.873} \\ X'_\ell \leq \text{RU}(T + K_h \cdot 2^{-39}) \leq 2^{-37.141} \\ |x'_h| \leq \text{RU}(K_h 2^{14}) < 2^{14.7}. \end{cases}$$

Since $\text{ulp}(2^{-38}) = 2^{-90}$, the error of the `DMUL` call is at most $2 \cdot 2^{-90} + K_\ell \cdot 2^{-39}$. The error due to $K_h + K_\ell$ not being $1/\log 2$ is at most $|K_h + K_\ell - 1/\log 2| \cdot |x| \leq 2^{14-109.5}$. The total error of the algorithm is thus at most

$$2^{-89} + K_\ell \cdot 2^{-39} + 2^{-95.5} \leq 2^{-88.951}.$$

If $|x'_h| < 2^{-20}$, necessarily $|x_h| \leq t_0$ where $t_0 = 0 \times 1.62\text{e}42\text{fefa}39\text{efp}-21$, because $x'_h = \text{o}(K_h x_h)$, and if we multiply K_h by $t_0 + \text{ulp}(t_0)$ with rounding towards zero, we get 2^{-20} . This implies $|x_\ell| < \text{ulp}(t_0) = 2^{-73}$. Using Lemma 1 with the improved bounds, we get

$$\begin{cases} T = \text{RU}(\text{ulp}(\text{RU}(K_h t_0)) + K_\ell t_0) \leq 2^{-72.827} \\ X'_\ell \leq \text{RU}(T + K_h \cdot 2^{-73}) \leq 2^{-71.638}. \end{cases}$$

This allows us to show that the total error is at most $2^{-123.3}$ in the same way as above. ■

b) *Subroutine SPLIT*: We use a variant of the absolute splitting of [6] to truncate x' to a fixed precision 2^{-20} . We use the “magic” double precision constant $C = 3 \cdot 2^{31}$.

Algorithm 9 Splitting to 2^{-20} precision

Input : $x' = (x'_h, x'_\ell)$ with $|x'_h| < 2^{14.7}$ and $|x'_\ell| < 2^{-37.141}$

Output : (n, f, y) such that $n + f/2^{20} + y$ approximates x'

```

1: if  $|x'_h| < 2^{-20}$  then
2:   return  $(0, 0, x')$ 
3: else
4:    $S \leftarrow \circ(C + x'_h)$ 
5:   Extract  $(n, f)$  from  $S$  bitwise
6:    $r \leftarrow \circ(x'_h - \circ(S - C))$ 
7:    $(y_h, y_\ell) = \text{FASTTWO\SUM}(r, x'_\ell)$ 
8:   return  $(n, f, (y_h, y_\ell))$ 
9: 
```

In line 5, f is formed as an unsigned integer by the 20 low bits of S and n is extracted from the 16 following upper bits interpreted as a two’s complement integer.

Lemma 6: Algorithm 9 returns (n, f, y) such that

$$x' = n + \frac{f}{2^{20}} + y + \epsilon$$

with n, f integers, $0 < f < 2^{20}$, $|y|, |y_h| < 2^{-19.999}$, $|y_\ell| \leq 2^{-71.63}$ and $|\epsilon| < 2^{-87.922}$.

Proof: If $|x'_h| < 2^{-20}$, we have $y_h = x'_h$ and $y_\ell = x'_\ell$, and the statement follows from Lemma 5, with $\epsilon = 0$. Assume now $|x'_h| \geq 2^{-20}$, thus we use the else branch of the algorithm. Since $|x'_h| < 2^{14.7}$, we have $2^{32} \leq C + x'_h < 3 \cdot 2^{31} + 2^{14.7}$, thus S lies in the binade $[2^{32}, 2^{33})$, with $\text{ulp}(S) = 2^{-20}$. We thus have $S = C + x'_h + \tau$ with $|\tau| < 2^{-20}$, the sign of τ depending on the rounding mode. In line 5, we read the bits of $x'_h + \tau$ as a two’s complement integer, where $|x'_h + \tau| < 2^{15}$: $S = C + n + f \cdot 2^{-20}$. By Sterbenz’s lemma, the subtraction $S - C$ is exact, thus $r = \circ(x'_h - (S - C))$. Since $|x'_h| \geq 2^{-20}$, $\text{ulp}(x'_h) \geq 2^{-72}$, thus $-\tau = x'_h - (S - C)$ is an integer multiple of 2^{-72} , with $|\tau| < 2^{-20}$, thus $x'_h - (S - C)$ is exact too, and $r = -\tau = x'_h - (n + f \cdot 2^{-20})$. Since the extraction of (n, f) is exact, and likewise for the computation of r , the only rounding error might occur in the FASTTWO\SUM call, which we denote by ϵ : $r + x'_\ell = y_h + y_\ell + \epsilon$. It yields $x'_h + x'_\ell = n + f \cdot 2^{-20} + y + \epsilon$.

Now let us bound $|\epsilon|$ and $|y_\ell|$. If $|r| \geq |x'_\ell|$, the precondition of FASTTWO\SUM is fulfilled, and using Theorem 1 of [2], the FASTTWO\SUM error is bounded by $2^{-105}|y_h|$. Since $y_h = \circ(r + x'_\ell)$ with $|r| = |\tau| < 2^{-20}$ and $|x'_\ell| < 2^{-37.141}$, this yields $|y_h| < 2^{-19.9999}$ and $|\epsilon| < 2^{-124.999}$. If e is the rounding error in $y_h = \circ(r + x'_\ell)$, i.e., $y_h = r + x'_\ell - e$, we know (see for example the proof of Theorem 1 in [2]) that $y_\ell = \circ(e)$; since $|e| < \text{ulp}(y_h) \leq \text{ulp}(2^{-19.9999}) = 2^{-72}$, it yields $|y_\ell| \leq 2^{-72}$.

If $|r| < |x'_\ell|$, the precondition of FASTTWO\SUM is not fulfilled. However, it is known (see for example Theorem 4.3 from [8]) that FASTTWO\SUM behaves normally when $|r|$

and $|x'_\ell|$ lie in the same binade. Thus since $|x'_\ell| < 2^{-37.141}$, the precondition of FASTTWO\SUM is not fulfilled only when $|r| < 2^{-38}$. In that case Theorem 2 of [2] yields $|\epsilon| < 3 \cdot 2^{-53}|y_h|$. It gives $|y_h| \leq \circ(2^{-38} + 2^{-37.141}) < 2^{-36.507}$, and $|\epsilon| < 3 \cdot 2^{-89.507} < 2^{-87.922}$. Then since $y_h = r + x'_\ell + \alpha$ with $|\alpha| < \text{ulp}(y_h) \leq 2^{-89}$, and $y_h + y_\ell + \epsilon = r + x'_\ell$, we get $y_\ell + \epsilon = -\alpha$, thus $|y_\ell| < |\epsilon| + |\alpha| < 2^{-87.922} + 2^{-89} < 2^{-72}$.

Since $|y_h| < 2^{-19.9999}$ and $|y_\ell| < 2^{-72}$, we have $|y| < 2^{-19.9999} + 2^{-72} < 2^{-19.999}$. ■

Now if we use the output x' of Algorithm 8 in Algorithm 9, combining the bounds of Lemma 5 and Lemma 6, we get when $|x'_h| \geq 2^{-20}$:

$$\left| n + \frac{f}{2^{20}} + y - \frac{x}{\log 2} \right| < 2^{-87.346},$$

and when $|x'_h| < 2^{-20}$ since $\epsilon = 0$ in Lemma 6:

$$\left| n + \frac{f}{2^{20}} + y - \frac{x}{\log 2} \right| < 2^{-123.3}.$$

By picking alternate values of C , Algorithm 9 can be adapted to different scales. It does however break with overflows and underflows.

c) *Subroutine LOOKUP*: Tabulating the whole 2^{20} possible values for $2^{f/2^{20}}$ would be prohibitively expensive. We split f in base 32 as $f = \sum_{i=0}^3 f_i \cdot 2^{5i}$ with $0 \leq f_i < 32$ so that

$$2^{f/2^{20}} = \prod_{i=0}^3 2^{f_i/2^{20-5i}}.$$

For $j = 0..3$, let t_j be a table of double-double values approximating $2^{x/2^{20-5j}}$ to nearest for $x = 0..31$. We use the following algorithm:

Algorithm 10 LOOKUP

Input : integer $f \in [0, 2^{20} - 1]$

Output : $z = (z_h, z_\ell)$ approximating $2^{f/2^{20}}$

```

Split  $f = \sum_{i=0}^3 f_i 2^{5i}$  bitwise
 $a_i = (a_{ih}, a_{i\ell}) \leftarrow t_i[f_i]$  for  $i = 0..3$ 
 $b = (b_h, b_\ell) \leftarrow \text{DMUL}(a_0, a_1)$ 
 $b' = (b'_h, b'_\ell) \leftarrow \text{DMUL}(a_2, a_3)$ 
 $z \leftarrow \text{DMUL}(b, b')$ 

```

Since each double-double needs 16 bytes of memory, the total lookup table size is only $16 \times 32 \times 4$ bytes, or 2kB. Moreover the tables are easily cache-aligned. We precomputed the tables t_j using SageMath.

Lemma 7: Subroutine LOOKUP’s result $z = (z_h, z_\ell)$ is such that

$$\left| z - 2^{f/2^{20}} \right| \leq 2^{f/2^{20}} \cdot 2^{-100.540}.$$

Moreover $|z_h| < 2$ and $|z_\ell| < 2^{-49.271}$.

Proof: Since the total table size is 2^{20} , we can do an exhaustive check against MPFR for each rounding mode. Then, the maximal error can be shown to be between $v = 0 \times 1.5\text{fff}27\text{a94p-101}$ (for rounding toward zero or downward) and $v + 2^{-130}$. This exhaustive check also provides the maximum values of $|z_h|$ and $|z_\ell|$. ■

$$\begin{aligned}
Q &= 1 \\
&+ (0x1.62e42fefa39efp-1 \\
&+ 0x1.abc9e3b369936p-56) \cdot y \\
&+ 0x1.ebfbdf82c696p-3 \cdot y^2 \\
&+ 0x1.c6b08d704a1cdp-5 \cdot y^3.
\end{aligned}$$

Fig. 1. Degree-3 polynomial approximating 2^y for $|y| < 2^{-19.999}$.

d) *Subroutine POLY*: We implement POLY using the degree-3 polynomial of Figure 1. This polynomial was found using Sollya and perfectly accurate evaluation of it would yield an absolute error at most $2^{-89.218}$ for $|y| < 2^{-19.999}$. We evaluate it using the following scheme:

$$Q(y) \approx 1 + y \cdot q_1 + y_h^2(q_2 + y_h \cdot q_3),$$

where q_1 decomposes into $q_{1h} + q_{1\ell}$ as double-double. The terms of order 2 and above are only computed in double-precision. Using an FMA and an independent multiply, this is particularly efficient.

Algorithm 11 POLY

Input: $y = y_h + y_\ell$, $|y|, |y_h| < 2^{-19.999}$, $|y_\ell| < 2^{-71.63}$

Output: $t = t_h + t_\ell$ approximating 2^y

- 1: $w \leftarrow \circ(y_h^2)$
 - 2: $u \leftarrow \circ(q_3 y_h + q_2)$
 - 3: $v \leftarrow \circ(wu)$
 - 4: $d_h, d_\ell \leftarrow \text{DMUL}(q_{1h}, q_{1\ell}, y_h, y_\ell)$
 - 5: $f_h, f_\ell \leftarrow \text{FASTTWO\SUM}(1, v)$
 - 6: $t_h, z \leftarrow \text{FASTTWO\SUM}(f_h, d_h)$
 - 7: $t_\ell \leftarrow \circ(z + \circ(f_\ell + d_\ell))$
-

Lemma 8: Let $y = (y_h, y_\ell)$ be a double-double value such that $|y|, |y_h| \leq 2^{-19.999}$ and $|y_\ell| \leq 2^{-71.63}$. Then POLY's result t is such that $|t_h| < 1.01$, $|t_\ell| < 2^{-50.948}$, and

$$|t - 2^y| \leq 2^y \cdot 2^{-89.0001}.$$

Proof: Since $|y_h| < 2^{-19.999}$, we have $|w| = \circ(y_h^2) < 2^{-39.997}$, and the rounding error on w is bounded by $\text{ulp}(2^{-39.997}) = 2^{-92}$, and $|w - y^2| \leq 2^{-92} + 2|y_h y_\ell| + y_\ell^2 \leq 2^{-90.157}$.

After $u = \circ(q_3 y_h + q_2)$, since $|q_3| < 1$ and $|q_2| < 0.2403$, we have $|q_3 y_h + q_2| < 2^{-19.999} + 0.2403 < 0.241$. Thus $|u| < 0.242$ and the rounding error on u is bounded by $\text{ulp}(0.242) = 2^{-55}$. Since we neglected $q_3 y_\ell$ which is bounded in absolute value by $2^{-71.63}$, we have $|u - (q_3 y + q_2)| < 2^{-55} + 2^{-71.63} < 2^{-54.999}$.

Now $v = \circ(wu)$ approximates $q_2 y^2 + q_3 y^3$ with $|wu| \leq 2^{-39.997} \cdot 0.242$, thus $|v| < 2^{-42.043}$, and the rounding error

on v is bounded by $\text{ulp}(2^{-42.043}) = 2^{-95}$:

$$\begin{aligned}
|v - (q_2 y^2 + q_3 y^3)| &\leq 2^{-95} + |wu - y^2(q_2 + q_3 y)| \\
&\leq 2^{-95} + |w - y^2|u + y^2|u - (q_2 + q_3 y)| \\
&< 2^{-95} + 2^{-90.157} \cdot 0.242 + 2^{-39.998} \cdot 2^{-54.999} \\
&< 2^{-91.838}.
\end{aligned}$$

After $d_h, d_\ell \leftarrow \text{DMUL}(q_{1h}, q_{1\ell}, y_h, y_\ell)$, $d_h + d_\ell$ approximates $q_1 y$. Knowing that $|y_h| < 2^{-19.999}$ and $|y_\ell| < 2^{-71.63}$, Lemma 1 gives

$$\begin{cases} T = \text{RU}(\text{ulp}(\text{RU}(q_{1h} 2^{-19.999})) + q_{1\ell} 2^{-19.999}) \leq 2^{-72.72}, \\ D_\ell = \text{RU}(T + q_{1h} 2^{-71.63}) \leq 2^{-71.41}, \\ D_h = \text{RU}(q_{1h} 2^{-19.999}) \leq 2^{-20.52}. \end{cases}$$

The total error of the DMUL call is therefore at most

$$|d_h + d_\ell - q_1 y| < 2^{-125} + 2^{-124} + 2^{-55.25} \cdot 2^{-71.63} < 2^{-123.289}.$$

After $f_h, f_\ell \leftarrow \text{FASTTWO\SUM}(1, v)$, $f_h + f_\ell$ approximates $1 + q_2 y^2 + q_3 y^3$. Since $|v| < 2^{-42.043}$, the FastTwoSum condition is fulfilled, and using Theorem 1 of [2] the rounding error is bounded by $2^{-105}|f_h| < 2^{-104.999}$. Also, since $|f_h| < 2$ it follows $|f_\ell| < \text{ulp}(1) = 2^{-52}$:

$$|f_h + f_\ell - (1 + q_2 y^2 + q_3 y^3)| < 2^{-104.999} + 2^{-91.838} < 2^{-91.837}.$$

After $t_h, z \leftarrow \text{FASTTWO\SUM}(f_h, d_h)$, $t_h + z + f_\ell + d_\ell$ approximates $Q(y)$. Since $|f_h| \geq \circ(1 - 2^{-42.043}) \geq 1/2$ and $|d_h| < 2^{-20.5}$, the FastTwoSum condition is fulfilled. Furthermore $|f_h + d_h| < \circ(1 + 2^{-42.043}) + 2^{-20.5} < 1.01$, and by the same argument as above, the FastTwoSum rounding error is bounded by $2^{-105} \cdot 1.01 < 2^{-104.985}$, and $|z| < \text{ulp}(1.01) = 2^{-52}$:

$$\begin{aligned}
|t_h + z + f_\ell + d_\ell - Q(y)| &< 2^{-104.985} + 2^{-91.837} + 2^{-123.289} < 2^{-91.836}.
\end{aligned}$$

Finally after $t_\ell \leftarrow \circ(z + \circ(f_\ell + d_\ell))$, $t_h + t_\ell$ approximates $Q(y)$. Let $\sigma = \circ(f_\ell + d_\ell)$. Since $|f_\ell| < 2^{-52}$ and $|d_\ell| < 2^{-71.4}$, it follows $|\sigma| < 2^{-51.9}$ and the rounding error on σ is bounded by $\text{ulp}(2^{-51.9}) = 2^{-104}$. Since $|z| < 2^{-52}$, we have $|z + \sigma| < 2^{-52} + 2^{-51.9} < 2^{-50.949}$, thus $|t_\ell| < 2^{-50.948}$ and the rounding error on t_ℓ is bounded by $\text{ulp}(2^{-50.948}) = 2^{-103}$:

$$|t_\ell - (z + f_\ell + d_\ell)| < 2^{-104} + 2^{-103} < 2^{-102.415}.$$

Summarizing we get:

$$|t_h + t_\ell - Q(y)| < 2^{-91.836} + 2^{-102.415} < 2^{-91.835}.$$

Since $Q(y)$ approximates 2^y for $|y| < 2^{-19.999}$ with absolute error bounded by $2^{-89.218}$, and using $2^y > 0.999999$:

$$|t_h + t_\ell - 2^y| < 2^{-91.835} + 2^{-89.218} < 2^{-89.0001} \cdot 2^y.$$

We can now state the main theorem: ■

Theorem 1: Given $x \in \mathbb{F}_{64}$, $2^{-64} \leq |x| < 2^{14}$, the integer n and the double-double number s computed by Algorithm 7 satisfy:

$$e^x = 2^n \cdot s \cdot (1 + \mu)$$

with $|\mu| < 2^{-87.193}$, $0.999999 < s_h < 2.03$ and $|s_\ell| < 2^{-48.2}$.

Proof: From Lemma 5, the double-double number $x' = x'_h + x'_\ell$ computed by Algorithm 8 satisfies $|x'_h| < 2^{14.7}$, $|x'_\ell| < 2^{-37.141}$ and $x' = x/\log 2 + \tau$ with $|\tau| < 2^{-88.951}$. Thus Lemma 6 applies and (n, f, y) satisfy $x' = n + f/2^{20} + y + \epsilon$ with $0 < f < 2^{20}$, $|y|, |y_h| < 2^{-19.999}$, $|y_\ell| \leq 2^{-71.63}$ and $|\epsilon| < 2^{-87.922}$. It follows:

$$e^x = 2^n \cdot 2^{f/2^{20}} \cdot 2^y \cdot 2^{\epsilon - \tau}.$$

Now from Lemma 7, the result z from subroutine LOOKUP satisfies $z = 2^{f/2^{20}}(1 + \eta)$ with $|\eta| < 2^{-100.540}$, and by Lemma 8, the double-double value t computed by subroutine POLY satisfies $t = 2^y(1 + \alpha)$ with $|\alpha| < 2^{-89.0001}$. thus

$$e^x = 2^n \cdot z \cdot t \cdot \frac{2^{\epsilon - \tau}}{(1 + \eta)(1 + \alpha)}.$$

Using Lemma 1 on the DMUL(t, z) call using the postconditions of Lemma 7 and 8 and renaming A and B the quantities from the lemma to avoid conflicts, we have

$$\begin{cases} A = \text{RU}(2.02 \cdot 2^{-52} + 2 \cdot 2^{-50.948}) \leq 2^{-49.3} \\ B \leq \text{RU}(A + 1.01 \cdot 2^{-49.271}) \leq 2^{-48.2}, \end{cases}$$

so that the last DMUL's error is bounded by $2^{-102} + 2^{-101} + 2^{-50.948} \cdot 2^{-49.271} < 2^{-99.313}$.

Since $|z| \geq 1$ (by looking at the table values), and $|t| \geq 0.999999$ (since t approximates 2^y with $y \geq -2^{-19.999}$), we deduce the relative error of DMUL is bounded by $2^{-99.313}/0.999999 < 2^{-99.312}$:

$$s = tz(1 + \beta) \quad \text{with } |\beta| < 2^{-99.312}.$$

Summarizing, we get:

$$e^x = 2^n \cdot s \cdot \frac{2^{\epsilon - \tau}}{(1 + \eta)(1 + \alpha)(1 + \beta)},$$

which implies $e^x = 2^n \cdot s \cdot (1 + \mu)$ with $|\mu| < 2^{-87.193}$. ■

B. Accurate path

The accurate path follows the same computation scheme as the fast path, leveraging the extra software precision to attain a relative error bound of $2^{-167.006}$.

To reduce total lookup table size we reuse the fast path's lookup tables using the "accurate tables" trick from [7]. Let $1 \leq k \leq 4$ and $0 \leq t < 32$. Assume the fast path's lookup tables approximated $2^{t/2^{5k}}$ as x . We then know that

$$\left| x - 2^{t/2^{5k}} \right| \leq 2^{-107}.$$

As a consequence, there must be some small s such that

$$x = 2^s \cdot 2^{t/2^{5k}}.$$

By storing a 62-bit approximation of s , we can make the lookup phase output an approximation pair (s', y) with $y \approx 2^{f/2^{20} + s'}$ and s' small. By adding s' to the reduced value, this allows us to get a lookup accuracy of $107 + 62 = 169$ bits using only 4 tables of 32 words of 64 bits, which fit in 1kB.

C. Computing hard-to-round cases with BaCSeL

We have computed hard-to-round inputs using BaCSeL. BaCSeL [4] is a software tool to compute hard-to-round cases of a mathematical function. BaCSeL takes as input a mathematical function f , a range $[x_0, x_1)$, an input precision n , an output precision n' , an integer m , and outputs all n -bit floating-point numbers x in $[x_0, x_1)$ such that:

$$|f(x) - \circ(f(x))| < 2^{-m} \text{ulp}(f(x)),$$

where $\circ(\cdot)$ is any rounding mode. In other words, $f(x)$ is at distance 2^{-m} ulp from a rounding boundary, or there are at least $m - 1$ identical bits after the round bit. The complexity of the underlying algorithm depends on several factors, but the larger is m , the faster is the algorithm.

Let x_2 be the smallest long double value such that $\exp x > 2^{16384}$ and x_1 the largest such that $\exp x_1 < 2^{-16382}$, which is the smallest positive normal number. Notice that for $|x| < 2^{-65}$, $\exp x$ rounds to the same value as $1 + x$, thus we only have to search for hard-to-round inputs in the ranges $[x_1, -2^{-65})$ and $[2^{-65}, x_2)$. For $|x|$ large, the efficiency of the SLZ algorithm used by BaCSeL decreases, because the difference between $\exp x$ and $\exp y$ for two adjacent numbers x and y becomes large with respect to $\text{ulp}(\exp x)$. For this reason, we further subdivided the search into sub-ranges:

- for $-2^4 < x \leq -2^{-65}$ and $2^{-65} \leq x < 0x1.484p+9$, we search worst cases with at least 54 identical bits after the round bit. We found 158,662 such worst cases. The largest number of identical bits after the rounding bit is 126 for $x = 0x1.ffffffffffffffffffffep-64$. Apart from such inputs very close to zero and with a special form due to the Taylor expansion of \exp , the largest number of identical bits after the rounding bit is 75 for $x = -0x1.625ac7bfa54aba72p-14$, and then 71 for $x = 0x1.fe3ab8dba4df80d6p+0$;
- for $x_1 \leq x \leq -2^4$ and $0x1.484p+9 \leq x < x_2$, we search worst cases with at least 101 identical bits after the round bit. We found no such worst case.

a) *Subnormal output:* Let x_0 be the smallest long double such that $\exp x_0 \geq 2^{-16445}$. For $x_0 \leq x < x_1$, $\exp x$ lies in the subnormal range $[2^{-16445}, 2^{-16382})$. For each exponent e , $-16444 \leq e \leq -16382$, we compute the largest interval (ℓ, h) such that $2^{e-1} \leq \exp(\ell) \leq \exp(h) < 2^e$, and search with BaCSeL inputs in this interval with at least 101 identical bits after the round bit (which has weight 2^{-16446} for all intervals). We found no such input.

b) *Correctness of the accurate path:*

Lemma 9: Assume a floating-point format with a precision of p bits, and an accurate path with relative error less than 2^k . Then if x is such that $f(x)$ is at distance at least 2^{-m} ulp from a floating-point number, or from the middle of two adjacent

	expl	pow1
our code	46.2	161.7
Intel Math Library (2025.0.0)	63.0	287.6
GNU Libc 2.40	124.1	747.2
Openlibm 0.8.3	149.5	626.2
Musl 1.2.5	112.2	535.5

Fig. 2. Reciprocal throughput (in cycles) of some double extended precision functions on a Intel Xeon Silver 4214 with GCC 14.2.0.

floating-point numbers, then rounding the approximation y of the accurate path delivers the correct rounding of $f(x)$ as long as:

$$m \leq -k - p - 1 + \log_2(1 - 2^k). \quad (1)$$

Moreover this remains true in the subnormal range.

Proof: Let x as in the statement of the lemma, and z a $(p+1)$ -bit floating-point number. We then have:

$$|f(x) - z| \geq 2^{-m} \text{ulp}_p(z) > 2^{-m-p}|z|.$$

The approximation y of $f(x)$ computed by the accurate path satisfies $|f(x) - y| < \varepsilon$ with $\varepsilon = 2^k|y|$. Assume the confidence interval $[y - \varepsilon, y + \varepsilon]$ crosses a rounding boundary z . We then have:

$$|z - y| \leq 2^k|y|, |f(x) - y| < 2^k|y|, |f(x) - z| > 2^{-m-p}|z|.$$

By the triangular inequality, it follows $2^{k+1}|y| > 2^{-m-p}|z| \geq 2^{-m-p}(1 - 2^k)|y|$. This contradicts the hypothesis on m , thus the confidence interval $[y - \varepsilon, y + \varepsilon]$ cannot contain any rounding boundary, which proves that any number in this interval rounds to the same value. Therefore rounding y yields the correct rounding of $f(x)$.

In the subnormal range, the output “local” precision p' satisfies $p' \leq p$: if Eq. (1) holds for p , it also holds for p' . ■

In the `expl` case, the relative error of the accurate path is less than 2^k with $k = -167.006$, we have $p = 64$, thus from Lemma 9 we find that $m \leq 102$ suffices.

This corresponds to `-m 102` in BaCSeL syntax, which means searching all inputs such that $|\text{exp } x - z| < 2^{-102} \text{ulp}_{64}(z)$, where z is a 65-bit floating-point number (exact double extended or midpoint).

The correctness of the accurate path follows from the following remarks:

- if $f(x)$ has at least 101 identical bits after the round bit, x will be found by our search with BaCSeL. It thus suffices to check that $f(x)$ is correctly rounded by the accurate path, and otherwise treat x as exceptional case;
- otherwise, $f(x)$ has less than 101 identical bits after the round bit, and by Lemma 9 it is correctly rounded.

V. RESULTS AND CONCLUSION

Using the previously-described methods, we implemented correctly-rounded versions of `expl`, `log21` and `pow1`. Using

double-double arithmetic and Newton’s method, we implemented `cbrtl` and `rsqrtl`. This is particularly effective because the double precision version of these functions gives a nice starting point for Newton’s method.

We certified the correctness of our implementations using BaCSeL, and tested their performance with the CORE-MATH benchmark suite. We consistently obtained better performance than the other mathematical libraries with double extended support (see Figure 2). We attribute this to our use of double-double arithmetic. Other implementations extensively use x87 features, which in modern processors are badly pipelined and slow [3].

The GNU Libc and the Intel Math Library use complex x87 features such as `FSQRT` which are not completely specified. As such, their results depend on the specific processors on which they are run. Our code on the other hand only uses FMA operations, which are completely specified. Because the correctness of our algorithms is based on the double precision FMA, we expect them to be portable across architectures, operating systems and processor manufacturers.

REFERENCES

- [1] BRISEBARRE, N., HANROT, G., MULLER, J.-M., AND ZIMMERMANN, P. Correctly-rounded evaluation of a function: why, how, and at what cost? <https://hal.science/hal-04474530>, 2024.
- [2] CORBINEAU, S., AND ZIMMERMANN, P. Note on FastTwoSum with Directed Roundings. <https://inria.hal.science/hal-03798376>, 2024.
- [3] FOG, A. *The microarchitecture of Intel, AMD and VIA CPUs: An optimization guide for assembly programmers and compiler makers*. 2024.
- [4] HANROT, G., LEFÈVRE, V., STEHLÉ, D., AND ZIMMERMANN, P. The BaCSeL software tool. <https://gitlab.inria.fr/zimmerma/bacsel>.
- [5] HUBRECHT, T., JEANNEROD, C.-P., AND ZIMMERMANN, P. Towards a correctly-rounded and fast power function in binary64 arithmetic. In *2023 IEEE 30th Symposium on Computer Arithmetic (ARITH 2023)* (Portland, Oregon (USA), United States, 2023), vol. 2023 IEEE 30th Symposium on Computer Arithmetic (ARITH).
- [6] JEANNEROD, C.-P., MULLER, J.-M., AND ZIMMERMAN, P. On various ways to split a floating-point number. In *2023 IEEE 30th Symposium on Computer Arithmetic (ARITH 2023)* (Amherst, Massachusetts (USA), United State, 2018), pp. 53–60.
- [7] MARKSTEIN, P. *IA-64 and Elementary Functions: Speed and Precision*. Hewlett-Packard professional books. Prentice Hall PTR, 2000.
- [8] 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, 2018.
- [9] ZIV, A. Fast evaluation of elementary mathematical functions with correctly rounded last bit. *ACM Trans. Math. Softw.* 17, 3 (1991), 410–423.