



**HAL**  
open science

# Correctly Rounded Arbitrary-Precision Floating-Point Summation

Vincent Lefèvre

► **To cite this version:**

Vincent Lefèvre. Correctly Rounded Arbitrary-Precision Floating-Point Summation. 2015. hal-01242127v1

**HAL Id: hal-01242127**

**<https://inria.hal.science/hal-01242127v1>**

Preprint submitted on 11 Dec 2015 (v1), last revised 7 Nov 2016 (v3)

**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.

# Correctly Rounded Arbitrary-Precision Floating-Point Summation

Vincent Lefèvre

INRIA, LIP / CNRS / ENS Lyon / Université de Lyon

Lyon, France

Email: [vincent@vinc17.net](mailto:vincent@vinc17.net)

**Abstract**—We present a fast algorithm together with its low-level implementation of correctly rounded arbitrary-precision floating-point summation. The arithmetic is the one used by the GNU MPFR library: radix 2; no subnormals; each variable (each input and the output) has its own precision. We also describe how the implementation is tested.

**Index Terms**—summation; floating point; arbitrary precision; multiple precision; correct rounding;

## I. INTRODUCTION

In a floating-point system, the summation operation consists in evaluating the sum of several floating-point numbers. The IEEE 754 standard for floating-point arithmetic introduced the sum reduction operation in its 2008 revision [1, Clause 9.4], but provides almost no specification. The IEEE 1788-2015 standard for interval arithmetic goes further by completely specifying this sum operation for IEEE 754 floating-point formats [2, Clause 12.12.12], in particular requiring correct rounding and specifying the sign of an exact zero result (in a way that is incompatible with IEEE 754).

The articles in the literature on floating-point summation mainly focus on IEEE 754 arithmetic and consider the floating-point arithmetic operations (+, −, etc.) as basic blocks; in this context, inspecting bit patterns is generally not interesting. For instance, fast and accurate summation algorithms are presented by Demmel and Hida [3] and by Rump [4]. Correct rounding is not provided. On this subject, the class of algorithms that can provide a correctly rounded sum of  $n \geq 3$  numbers is somewhat limited [5]. In [6], Rump, Ogita and Oishi present correctly rounded summation algorithms. Kulisch proposes a quite different solution: the use of a long accumulator covering the full exponent range (and a bit more to handle intermediate overflows) [7]. A survey of summation methods can be found in [8, Section 6.3].

In IEEE 754, the precision of each floating-point format is fixed. In this paper, we deal with the extension of the summation operation to arbitrary precision in radix 2, where each number has its own precision and results must be correctly rounded, as with the GNU MPFR library<sup>1</sup> [9], where this function is named `mpfr_sum`.

We first give some notation (Section II). In Section III, we present a brief overview of GMP and GNU MPFR. In Section IV, we describe the old `mpfr_sum` implementation and

explain why a new one was needed. In Section V, we give the complete specification of the summation operation in MPFR. In Section VI, we present the completely new algorithm and implementation; since this is a low-level algorithm, the context of MPFR is quite important for the details, but the main ideas could be reused in other contexts. In Section VII, we explain how `mpfr_sum` is tested.

## II. NOTATION

In addition to ISO 80000-2 notation, we will use  $\llbracket$  and  $\rrbracket$  for the bounds of integer intervals, e.g.  $\llbracket 0, 3 \rrbracket = \{0, 1, 2, 3\}$  and  $\llbracket 0, 3 \llbracket = \{0, 1, 2\}$ .

A minor typography difference will be used between variable names (e.g. `minexp`) and their values (e.g. `minexp`).

## III. OVERVIEW OF GMP AND GNU MPFR

GNU MPFR is a free library for efficient arbitrary-precision floating-point computation with well-defined semantics (copying the good ideas from the IEEE 754 standard), in particular correct rounding. It is based on GNU MP (GMP)<sup>2</sup>, which is a free library for arbitrary-precision arithmetic; MPFR mainly uses the low-level GMP layer called “`mpn`”, and we will restrict to it here. As said on the GMP web page: “*Low-level positive-integer, hard-to-use, very low overhead functions are found in the `mpn` category. No memory management is performed; the caller must ensure enough space is available for the results.*”

In this layer, a natural number is represented by an array of words, called *limbs*, each word corresponding to a digit in high radix ( $2^{32}$  or  $2^{64}$ ). The main GMP functions that will be useful for us are: the addition (resp. subtraction) of two  $N$ -limb numbers, with carry (resp. borrow) out; ditto between an  $N$ -limb number and a limb; left shift; right shift; negation with borrow out; complement.

Each MPFR floating-point number has its own precision in bits, starting at 2. All arithmetic operations are correctly rounded to the precision of the destination number in one of the 5 supported rounding modes: `MPFR_RNDN` (to nearest, with the even rounding rule), `MPFR_RNDD` (toward  $-\text{Inf}$ ), `MPFR_RNDU` (toward  $+\text{Inf}$ ), `MPFR_RNDZ` (toward zero), `MPFR_RNDA` (away from zero).

The MPFR numbers are represented with 3 fields: a sign, a significand (always normalized, with the leading bit 1

<sup>1</sup><http://www.mpfr.org/>

<sup>2</sup><https://gmplib.org/>

represented, and any trailing bit in the least significant limb being 0) interpreted as being in  $[1/2, 1[$ , and an exponent field, which contains either the exponent or a special value for singular datums: zero, infinity, and NaN (the significant contains garbage in such cases); contrary to IEEE 754, MPFR has only a single kind of NaN and does not have subnormals. An important point is that the exponent range can be very large: up to  $\llbracket 1 - 2^{62}, 2^{62} - 1 \rrbracket$  on 64-bit machines.

Moreover, most arithmetic operations return a ternary value, giving the sign of the rounding error.

In MPFR, exponents are stored in a signed integer type `mpfr_exp_t`. If this type has  $N$  value bits, i.e. the maximum value is  $2^N - 1$ , then the maximum exponent range is defined so that any valid exponent fits in  $N - 1$  bits (sign bit excluded), i.e. it is  $\llbracket 1 - 2^{N-1}, 2^{N-1} - 1 \rrbracket$ . This implies a huge gap between the minimum value `MPFR_EXP_MIN` of the type and the minimum valid exponent `MPFR_EMIN_MIN`; this allows us to make the following practical hypothesis.

We assume here that each represented bit of a regular MPFR number has an exponent value representable in the `mpfr_exp_t` type, i.e. a value that is larger or equal to `MPFR_EXP_MIN`. However, unless the precision of a MPFR number is huge, such a requirement is satisfied. In practice, it is always satisfied on 64-bit machines due to memory limits. On 32-bit machines, users would probably reach other system limits first (mainly due to the limit on the address space).

#### IV. THE OLD `MPFR_SUM` IMPLEMENTATION

The implementation of `mpfr_sum` from the current MPFR releases (up to version 3.1.x) is based on Demmel and Hida's accurate summation algorithm [3], which consists in adding the inputs one by one in decreasing magnitude. But here, this has several drawbacks:

- This is an algorithm using only high-level operations, mainly floating-point additions (in MPFR, `mpfr_add`). This is the right way to do to get an accurate sum in true IEEE 754 arithmetic implemented in hardware, but in MPFR, which uses integers as basic blocks, this introduces overheads, and more important problems mentioned just below.
- Due to the high-level operations, correct rounding had to be implemented with a Ziv loop: the working precision is increased until the rounding can be guaranteed [9]. In the case of summation, this gives a time and memory worst-case complexity exponential in the number of bits of the exponent field. In practice, this is very slow in some cases, and worse, since the exponent range can be large, this can yield a crash due to the lack of memory (and possible denial of service for other processes running on the machine).
- Demmel and Hida's algorithm is based on the fact that the precision is the same for all floating-point numbers, meaning that in the MPFR implementation, the maximum precision had to be chosen. An alternative would be to split the input numbers to numbers with the same precision, but this would introduce another overhead.

Moreover the sign of an exact zero result is not specified and the ternary value is valid only when it is zero (a non-zero return value provides no information).

#### V. SPECIFICATION OF THE SUMMATION OPERATION

The prototype of the `mpfr_sum` function is:

```
int mpfr_sum (mpfr_ptr sum,
              mpfr_ptr *const x,
              unsigned long n,
              mpfr_rnd_t rnd)
```

where `sum` will contain the correctly rounded sum, `x` is an array of pointers to the inputs, `n` is the length of this array, and `rnd` is the rounding mode. The return value of type `int` will be the usual ternary value. It is currently assumed that no input pointers point to `sum`, i.e. that the output number `sum` is not shared with an input.<sup>3</sup>

If `n = 0`, then the result is `+0`, whatever the rounding mode. This is equivalent to `mpfr_set_ui` and `mpfr_set_si` on the integer 0, which both assign a MPFR number from a mathematical zero (not signed), and this choice is consistent with the IEEE 754 sum operation of vector length 0.

Otherwise the result must be the same as the one that would have been obtained with

- if `n = 1`: a copy with rounding (`mpfr_set`);
- if `n > 1`: a succession of additions (`mpfr_add`) done in infinite precision, then rounded (the order of these additions does not matter).

This is equivalent to apply the following rules:

- 1) If an input is NaN, then the result is NaN.
- 2) If there are at least a `+Inf` and a `-Inf`, then the result is NaN.
- 3) If there is at least an infinity (in which case all the infinities have the same sign), then the result is this infinity.
- 4) If the result is an exact zero:
  - if all the inputs have the same sign (thus all `+0`'s or all `-0`'s), then the result has the same sign as the inputs;
  - otherwise, either because all inputs are zeros with at least a `+0` and a `-0`, or because some inputs are non-zero (but they globally cancel), the result is `+0`, except for the `MPFR_RNDD` rounding mode, where it is `-0`.
- 5) Otherwise the exact result is a non-zero real number, and the conventional rounding function is applied.

#### VI. NEW ALGORITHM AND IMPLEMENTATION

The new algorithm is carefully designed so that the time and memory complexity no longer depends on the value of the exponents of the inputs, i.e. the orders of magnitude of the inputs. Instead of being high level (based on `mpfr_add`), the algorithm/implementation is low level, based on integer

<sup>3</sup>This will be fixed later. This can be handled by using a temporary area for the output when sharing occurs, or in some optimized way.

operations, equivalently seen as fixed-point operations. Efficiency in case of cancellations and Table Maker’s Dilemma is regarded as important as for cases without such issues. To be as fast as possible, we will use the `mpn` layer of GMP.

As a bonus, this will also solve overflow, underflow and normalization issues, since everything is done in fixed point and the exponent of the result will be considered only at the end (early overflow detection could also be done, but this would probably not be very useful in practice).

The idea is the following. After handling special cases (NaN, infinities, only zeros, and fewer than 3 regular inputs), we apply the generic case, which more or less consists in a fixed-point accumulation by blocks: we take into account the bits of the inputs whose exponent is in some window  $[\text{minexp}, \text{maxexp}]$ , and if this is not sufficient due to cancellation, then we reiterate, using a new window with lower exponents. Once we have obtained an accurate sum, if one still cannot round correctly because the result is too close to a rounding boundary (i.e. a machine number or the middle of two consecutive machine numbers), which is the problem known as the Table Maker’s Dilemma (TMD), then this problem is solved by determining the sign of the remainder by using the same method in a low precision.

#### A. Preliminary Steps

We start by detecting the special cases. The `mpfr_sum` function does the following.

If  $n \leq 2$ , we can use existing MPFR functions and macros, mainly for better efficiency since the algorithm described below can work with any number of inputs (only minor changes would be needed):

- If  $n = 0$ : return  $+0$  (by using MPFR macros).
- If  $n = 1$ : use `mpfr_set` (which copies a number, with rounding to the target precision).
- If  $n = 2$ : use `mpfr_add` (which adds two numbers, with rounding to the target precision).

Now,  $n \geq 3$ . We iterate on the  $n$  input numbers to:

- (A) detect singular values (NaN, infinity, zero);
- (B) among the regular values, get the maximum exponent.

Such information can be retrieved very quickly and this does not need to look at the significand. Moreover, in the current internal number representation, the kind of a singular value is represented as special values of the exponent field, so that (B) does not need to fetch more data in memory after doing (A).

In detail, during this iteration, 4 variables will be set, but the loop will terminate earlier if one can determine that the result will be NaN, either because of a NaN input or because of infinity inputs of opposite signs:

- `maxexp`, which will contain the maximum exponent of the inputs. Thus it is initialized to `MPFR_EXP_MIN`.
- `rn`, which will contain the number of regular inputs, i.e. those which are non-zero finite numbers.
- `sign_inf`, which will contain the sign of an infinity result. It is initialized to 0, meaning no infinities yet. When the first infinity is encountered, this value is changed to

the sign of this infinity ( $+1$  or  $-1$ ). When a new infinity is encountered, either it has the same sign of `sign_inf`, in which case nothing changes, or it has the opposite sign, in which case the loop terminates immediately and a NaN result is returned.

- `sign_zero`, which will contain the sign of the zero result *in the case where all the inputs are zeros*. Thanks to the IEEE 754 rules, this can be tracked with this variable alone: There is a weak sign ( $-1$ , except for `MPFR_RNDD`, where it is  $+1$ ), which can be obtained only when all the inputs are zeros of this sign, and a strong sign ( $+1$ , except for `MPFR_RNDD`, where it is  $-1$ ), which is obtained in all the other cases, i.e. when there is at least a zero of this sign. One could have initialized the value of `sign_zero` to the weak sign. But we have chosen to initialize it to 0, which means that the sign is currently unknown, and do an additional test in the loop. In practice, one should not see the difference; this second solution was chosen just because it was implemented first, and on a test, it made the code slightly shorter.

When the loop has terminated “normally”, the result cannot be NaN. We do in the following order:

- 1) If `sign_inf`  $\neq 0$ , then the result is an infinity of this sign, and we return it.
- 2) If `rn` = 0, then all the inputs are zeros, so that we return the result zero whose sign is given by `sign_zero`.
- 3) If `rn`  $\leq 2$ , then one can use `mpfr_set` or `mpfr_add` as an optimization, similarly to what was done for  $n \leq 2$ . We reiterate in order to find the concerned input(s), call the function and return.
- 4) Otherwise we call a function `sum_aux`, which implements the generic case. In addition to the parameters of `mpfr_sum`, we pass to this function:
  - the maximum exponent;
  - the number `rn` of regular inputs, i.e. the number of non-zero inputs. This number will be used instead of `n` to determine bounds on the sum (to avoid internal overflows) and error bounds.

#### B. Introduction to the Generic Case

Let us define  $\text{logn} = \lceil \log_2(\text{rn}) \rceil$ .

The basic idea is to compute a truncated sum in the two’s complement representation, by using a fixed-point accumulator stored in a fixed memory area.

Two’s complement is preferred to the sign + magnitude representation because the signs of the temporary (and final) results are not known in advanced, and the computations (additions and subtractions) in two’s complement are more natural in this context. There will be a conversion to sign + magnitude (representation used by MPFR numbers) at the end, but this should not take much time compared to the other calculations.

The precision of the accumulator needs to be a bit larger than the output precision, denoted `sq`, for two reasons:

- We need some additional bits on the side of the most significant part due to the accumulation of `rn` values,

which can make the sum grow and overflow without enough extra bits. The absolute value of the sum is less than  $rn \cdot 2^{\maxexp}$ , taking up to  $\log n$  extra bits, and one needs one more bit to be able to determine the sign due to two's complement, so that a total of  $cq = \log n + 1$  extra bits will be necessary.

- We need some additional bits on the side of the least significant part to take into account the accumulation of the truncation errors. The choice of this number  $dq$  of bits is quite arbitrary: the larger this value is, the longer an iteration will take, but conversely, the less likely a costly new iteration (due to cancellations and/or the Table Maker's Dilemma) will be needed. In order to make the implementation simpler, the precision of the accumulator will be a multiple of the limb size `GMP_NUMB_BITS`. Moreover the algorithm will need at least 4 bits. The final choice should be done after testing various applications. In the current implementation, we chose the smallest value larger or equal to  $\log n + 2$  such that the precision of the accumulator is a multiple of `GMP_NUMB_BITS`. Since  $\log n \geq 2$ , we have  $dq \geq 4$  as wanted.

The precision of the accumulator is defined as:

$$wq = cq + \maxexp - \minexp = cq + sq + dq$$

and  $\minexp$  will always be the exponent of the least significant bit (LSB) of the accumulator. In the accumulation, the selected bits from the inputs will range from  $\minexp$  (included) to  $\maxexp$  (excluded), and the most significant  $cq$  bits can only be reached due to carry propagation.

When the Table Maker's Dilemma occurs, the needed precision for the truncated sum would grow. In particular, one could easily reach a huge precision with a few small-precision inputs: for instance, in directed rounding modes,  $\text{sum}(2^E, 2^F)$  with  $F$  much smaller than  $E$ . We want to avoid increasing the precision of the accumulator. This will be done by detecting the Table Maker's Dilemma, and when it occurs, solving it consists in determining the sign of some error term. This will be done by computing an approximation to the error term in low precision. The algorithm to compute this error term is the same as the one to compute an approximation to the sum, the main difference being that we just need a low precision here. Thus we will define a function `sum_raw`, used for both computations; it is described in the next section.

### C. The `sum_raw` Function

The `sum_raw` function will work in a fixed-point accumulator, having a fixed precision (a multiple of `GMP_NUMB_BITS` bits) and using a two's complement representation. An iteration will consist in accumulating the bits of the inputs whose exponents are in  $[\minexp, \maxexp]$ , where  $\maxexp - \minexp$  is less than the precision of the accumulator: as said above, we need some additional bits in order to avoid overflows during the accumulation. On the entry, the accumulator may already contain a value from previous computations (it is the caller that clears it if need be): in some cases, some bits will have to be kept between the two `sum_raw` invocations.

During the accumulation, the bits of the  $i$ -th input  $x[i]$  whose exponents are strictly less than  $\minexp$  form the *tail* of this input. When the tail of  $x[i]$  is not empty, its exponent  $e_i$  is defined as the minimum between  $\minexp$  and the exponent of  $x[i]$ . Thus the absolute value of this tail is strictly less than  $2^{e_i}$ . This will give an error bound on the computed sum at each iteration.

At the end of an iteration, we do the following. If the computed result is 0 (meaning full cancellation), set  $\maxexp$  to the maximum exponent of the tails, set  $\minexp$  so that it corresponds to the least significant bit of the accumulator, and reiterate. Otherwise, let  $e$  and  $err$  denote the exponents of the computed result and of the error bound respectively. While  $e - err$  is less than some given bound denoted `prec`, shift the accumulator (as detailed later), update  $\maxexp$  and  $\minexp$ , and reiterate. For the caller, this bound must be large enough in order to reach some wanted accuracy. However it cannot be too large since the accumulator has a limited precision: we will need to make sure that if a reiteration is needed, then the cause is a partial cancellation, so that the determined shift count is non-zero, otherwise the variable  $\minexp$  would not change and one would get an infinite loop. Details and formal definitions are given later.

Notes:

- The reiterations will end when there are no more tails, but in the code, this is detected only when needed.
- This definition of the tails allows one to skip potentially huge holes between inputs in case of full cancellation, e.g.  $1 + (-1) + r$  where  $r$  has a tiny exponent.
- We choose not to include  $\maxexp$  in the exponent interval in order to match the convention chosen to represent floating-point numbers in MPFR, where the significand is in  $[1/2, 1[$ , i.e. the exponent of a floating-point number is the one of the most significant bit + 1. Another advantage is that  $\minexp$  at some iteration will be  $\maxexp$  at the next iteration, unless there is a hole between the inputs (i.e. the exponent of each tail is less than  $\minexp$ ).

Now let us give the details about this `sum_raw` function. First, it takes the following arguments:

- `wp`: pointer to the accumulator (least significant limb first).
- `ws`: size of the accumulator (in limbs).
- `wq`: precision of the accumulator ( $ws \times \text{GMP\_NUMB\_BITS}$ ).
- `x`: array of the input numbers.
- `n`: size of this array (number of inputs, regular or not).
- `minexp`: exponent of the LSB of the first window.
- `maxexp`: exponent of the first window (i.e. exponent of its MSB + 1).
- `tp`: pointer to a temporary area (pre-allocated).
- `ts`: size of this temporary area.
- `logn`:  $\lceil \log_2(rn) \rceil$ ,  $rn$  being the number of regular inputs.
- `prec`: lower bound for  $e - err$  (as described above).
- `ep`: pointer to `mpfr_exp_t` (see below).
- `minexp`: pointer to `mpfr_exp_t` (see below).
- `maxexp`: pointer to `mpfr_exp_t` (see below).

We require as preconditions (explanations are given later):  $\text{prec} \geq 1$  and  $\text{wq} \geq \log n + \text{prec} + 2$ .

This function returns 0 if the accumulator is 0 (which implies that the exact sum for this `sum_raw` invocation is 0), otherwise the number of cancelled bits, defined as the number of identical bits on the most significant part of the accumulator. In the latter case, it also returns the following data in variables passed by reference (via pointers):

- for `ep`: the exponent  $e$  of the computed result;
- for `minexp`: the last value of the variable `minexp`;
- for `maxexp`: the new value of the variable `maxexp` (the one for a potential new iteration).

The temporary area must be large enough to hold a shifted input block, and the value of `ts` is used only when the full assertions are checked, in order to make sure that a buffer overflow does not occur.

Some notation used below:

- $E(v)$ : the exponent of a MPFR number  $v$ .
- $P(v)$ : the precision of a MPFR number  $v$ .
- $Q(v) = E(v) - P(v)$ : the exponent of the ulp of a MPFR number  $v$ .

A `maxexp2` variable will contain the maximum exponent of the tails. Thus it is initialized to the minimum value of the exponent type: `MPFR_EXP_MIN`; this choice means that at the end of the loop below, `maxexp2 = MPFR_EXP_MIN` if and only if there are no more tails (this case implies that the truncated sum is exact). If a new iteration is needed, then `maxexp2` will be assigned to the `maxexp` variable for this new iteration.

Then one has a loop over the inputs `x[i]`. Each input is processed with the following steps:

- 1) If the input is not regular (i.e. is zero), skip it. Note: if there are many zero inputs, it may be more efficient to have an array pointing to the regular inputs only, but such a case is assumed to be rare, and the number of iterations of this inner loop is also limited by the relatively small number of regular inputs.
- 2) If  $E(x[i]) \leq \text{minexp}$ , then no bits of `x[i]` need to be considered here. We set the `maxexp2` variable to  $\max(\text{maxexp2}, E(x[i]))$ , then go to the next input.
- 3) Now, we have:  $E(x[i]) > \text{minexp}$ . If the tail of `x[i]` is not empty, i.e. if  $Q(x[i]) < \text{minexp}$ , then we set the `maxexp2` variable to `minexp`.
- 4) We prepare the input for the accumulation. First, this means that if its significand is not aligned with the accumulator, then we need to align it by shifting a part of the significand (containing bits that will be accumulated at this iteration), storing the result to the temporary area at address `tp`. In both cases, we compute:
  - the address and size of the destination, corresponding either to the full accumulator or to some most significant part of it (which can occur if  $Q(x[i]) > \text{minexp}$ );
  - the address and size of the block (of the input or its shifted version) to accumulate to the destination;

- the number of the most significant bits of the block that need to be masked later, taking into account that if the input is not shifted, this block is regarded as a read-only area.

- 5) If `x[i]` is positive: an addition is done with `mpn_add_n`; if the most significant limb needs to be masked, then it is not taken into account in the addition, but the masked part is just added to the carry; carry propagation is done with `mpn_add_1` if the size of the destination is larger than the size of the block.

Note: There may be still be a carry out, but it is just ignored. This occurs when a negative value in the accumulator becomes non-negative, and this fact is part of the usual two's complement arithmetic.

If `x[i]` is negative, we do similar computations by using `mpn_sub_n` for the subtraction and `mpn_sub_1` to propagate borrows.

After the loop over the inputs, we need to see whether the accuracy of the truncated sum is sufficient. We first determine the number of cancelled bits, defined as the number of consecutive identical bits starting with the most significant one in the accumulator. At the same time, we can determine whether the truncated sum is 0 (all the bits are identical and their value is 0). If it is 0, we have two cases: if `maxexp2 = MPFR_EXP_MIN` (meaning no more tails), then we return 0, otherwise we reiterate at the beginning of `sum_raw` with `maxexp` set to `maxexp2` and `minexp` set to  $\text{cq} + \text{maxexp} - \text{wq}$ .

We can now assume that the truncated sum is not 0.

Let us note that our computation of the number `cancel` of cancelled bits was limited to the accumulator representation, while from a mathematical point of view, the binary expansion is unlimited and the bits of exponent less than `minexp` are regarded as 0's. So, we need to check that the value `cancel` matches this mathematical point of view:

- If the cancelled bits are 0's: the truncated sum is not 0, therefore the accumulator must contain at least a bit 1.
- If the cancelled bits are 1's: this sequence of 1's entirely fits in the accumulator, since the first non-represented bit is a 0.

The analysis below virtually maps the truncated sum to the destination without considering rounding yet. Let us denote:  $e = \text{maxexp} + \text{cq} - \text{cancel} = \text{minexp} + \text{wq} - \text{cancel}$  and  $\text{err} = \text{maxexp2} + \log n$ .

Then  $e$  is the exponent of the least significant cancelled bit, thus the absolute value of the truncated sum is in  $[2^{e-1}, 2^e]$  (binade closed on both ends due to two's complement). And since there are at most  $\log n$  regular inputs and the absolute value of each tail is strictly less than  $2^{\text{maxexp2}}$ , the absolute value of the error is strictly less than  $2^{\text{err}}$ . If `maxexp2 = MPFR_EXP_MIN` (meaning no more tails), then the error is 0.

We need  $\text{prec} \geq 1$  to be at least able to determine the sign of the result, hence this precondition.

If  $e - \text{err} \geq \text{prec}$ , then the `sum_raw` function returns as described above.

Otherwise, due to cancellation, we need to reiterate after shifting the value of the accumulator to the left and updating

the  $\text{minexp}$  and  $\text{maxexp}$  variables. Let  $\text{shiftq}$  denote the shift count, which must satisfy:  $0 < \text{shiftq} < \text{cancel}$ . The left inequality must be strict to ensure termination, and the right inequality ensures that the value of the accumulator will not change with the updated  $\text{minexp}$ :  $\text{shiftq}$  is subtracted from  $\text{minexp}$  at the same time. The reiteration is done with  $\text{maxexp}$  set to  $\text{maxexp2}$ , as said above.

We now need to determine the value of  $\text{shiftq}$ . We prefer it to be as large as possible: this is some form of normalization. Moreover, it must satisfy the above double inequality and be such that:

- (A) the new value of  $\text{minexp}$  is smaller than the new value of  $\text{maxexp}$ , i.e.  $\text{minexp} - \text{shiftq} < \text{maxexp2}$ , i.e.  $\text{shiftq} > \text{minexp} - \text{maxexp2}$ ;
- (B) overflows will still be impossible in the new iteration.

Note that since  $\text{maxexp2} \leq \text{minexp}$ , (A) will imply  $\text{shiftq} > 0$ . And (B) is an extension of  $\text{shiftq} < \text{cancel}$ . Thus the double inequality above is a weak form of what is actually required in (A) and (B).

Since we prefer  $\text{shiftq}$  to be maximum, we focus on (B) first. The absolute value of the accumulator at the end of the next iteration will be strictly bounded by:  $2^e + 2^{\text{err}} \leq 2^{e+1+\max(0, \text{err}-e)}$ . This means that if we do not shift the value in the accumulator, then at the end of the next iteration, the accumulator will contain at least  $\text{cancel} - 1 - \max(0, \text{err} - e)$  identical bits on its most significant part. Only the last of these bits is needed (which gives the sign) and the other ones are redundant. Therefore, in order to satisfy (B), we can choose:

$$\text{shiftq} = \text{cancel} - 2 - \max(0, \text{err} - e).$$

Now, let us prove that for this value, (A) is satisfied, using the fact that  $wq \geq \log n + \text{prec} + 2$  on input.

- If  $\text{err} - e \geq 0$ , then by using  $\text{err} = \text{maxexp2} + \log n$  and  $e = \text{minexp} + wq - \text{cancel}$ , we obtain:

$$\begin{aligned} \text{shiftq} &= \text{cancel} - 2 - (\text{err} - e) \\ &= \text{minexp} - \text{maxexp2} + (wq - \log n - 2) \\ &> \text{minexp} - \text{maxexp2} \end{aligned}$$

- If  $\text{err} - e < 0$ , then this is the case where the error can be potentially small: to be able to prove the inequality, we need to use the fact that the stop condition was not satisfied, i.e.  $e - \text{err} < \text{prec}$ . Thus  $(\text{minexp} + wq - \text{cancel}) - (\text{maxexp2} + \log n) < \text{prec}$ , and as a consequence:

$$\begin{aligned} \text{shiftq} - (\text{minexp} - \text{maxexp2}) &= \text{cancel} - 2 - (\text{minexp} - \text{maxexp2}) \\ &> wq - \log n - \text{prec} - 2 \geq 0 \end{aligned}$$

Note: It is expected in general that when a cancellation occurs so that a new iteration is needed, the cancellation is not very large (but this really depends on the problem), in which case the new additions will take place only in a small part of the accumulator, except in case of long carry propagation.

#### D. The Generic Case

Let us recall that the accumulator for the summation is decomposed into three parts:  $cq = \log n + 1$  bits to avoid overflows,  $sq$  bits corresponding to the target precision, and  $dq$  additional bits to take into account the truncation error and improve the accuracy ( $dq \geq \log n + 2$  in the current implementation). Thus  $wq = cq + sq + dq$ .

Memory is allocated both for the accumulator and for the temporary area needed by `sum_raw`. For performance reasons, the allocation is done in the stack if the size is small enough. No other memory allocation will be needed (except for automatic variables).

The accumulator is zeroed and `sum_raw` is invoked to compute an accurate approximation of the sum. Among its parameters,  $\text{maxexp}$  was computed during the preliminary steps,  $\text{minexp} = \text{maxexp} - (wq - cq)$ , and  $\text{prec} = sq + 3$ , which satisfies the  $wq \geq \log n + \text{prec} + 2$  precondition.

If `sum_raw` returns 0, then the exact sum is 0, so that we just set the target sum to 0 with the correct sign according to the IEEE 754 rules (positive, except for `MPFR_RNDD`, where it is negative), and return with ternary value 0.

Now, the accumulator contains the significand of a good approximation to the non-zero exact sum. The corresponding exponent is  $e$  and the sign is determined from one of the cancelled bits. The exponent of the ulp for the target precision is denoted  $u = e - sq$ . The absolute value of the error is strictly less than  $2^{-3}$  times the ulp of the computed value:  $2^{u-3}$ .

When  $\text{maxexp}$  (value returned by `sum_raw`) is different from `MPFR_EXP_MIN`, i.e. when some bits of the inputs have still not been considered, we will need to determine whether the TMD occurs. We compute  $d = u - \text{err}$ , which is larger or equal to 3 (see above) and can be very large if  $\text{maxexp}$  is very small; however  $d$  is representable in a `mpfr_exp_t` since  $u \leq \text{MPFR_EMAX_MAX}$ ,  $\text{err} \geq \text{MPFR_EMIN_MIN}$ , and this type can hold the sum of two valid exponents. The TMD occurs when the sum is close enough to a breakpoint, which is either a machine number (i.e. a number whose significand fits on  $sq$  bits) or a midpoint between two consecutive machine numbers, depending on the rounding mode:

Rounding mode	Breakpoint
to nearest	midpoint
to nearest directed	machine number
	machine number

(in the second case, the correctly rounded sum can be determined, but not the ternary value, and this is why the TMD occurs). More precisely, the TMD occurs when:

- in directed rounding modes: the  $d$  bits following the ulp bit are identical;
- in round-to-nearest mode: the  $d - 1$  bits following the rounding bit are identical.

Several things need to be considered for the significand, in arbitrary order:

- the copy to the destination (significand of `sum`);
- a shift (for the normalization), if the shift count is non-zero (this is the most probable case);

- a negation/complement if the value is negative (cancelled bits = 1), since the significand of MPFR numbers uses the conventional sign + absolute value representation;
- rounding (the TMD needs to be resolved first if it occurs).

It is more efficient to merge some of these operations, i.e. do them at the same time, and this possibility depends on the operations provided by the mpn layer of GMP. Ideally, all these operations should be merged together, but this is not possible with GMP.

For practical reasons, the shift should be done before the rounding, so that all the bits are represented for the rounding. The copy itself should be done together with the shift or the negation, because this is where most of the limbs are changed in general. We chose to do it with the shift as it is assumed that the proportion of non-zero shift counts is higher than the proportion of negations.

Moreover, for negative values, the difference between negation and complement is similar to the difference between rounding directions (these operations are identical on the real numbers, i.e. in infinite precision), so that negation/complement and rounding can naturally be merged.

We start by doing the first two operations at the same time: the bits of exponents in  $\llbracket \max(u, \text{minexp}), e \rrbracket$  are copied with a possible shift to the most significant part of the destination, and the least significant limbs (if any) are zeroed.

By definition of  $e$ , the most significant bit that is copied is the opposite of the value of the cancelled bits. A variable `pos` will contain its value, i.e. `pos = 1` for a positive number, `pos = 0` for a negative number.

The values of three variables are also determined at about the same time:

- `inex`: 0 if the final sum is *known* to be exact, else 1.
- `rbit`: the rounding bit (0 or 1) of the truncated sum, corrected to 0 for halfway cases that round downward if rounding is to nearest (so that this bit gives the rounding direction).
- `tmd`: three possible values: 0 if the TMD does not occur, 1 if the TMD occurs on a machine number, 2 if the TMD occurs on a midpoint.

All this is done by considering two cases:  $u > \text{minexp}$  and  $u \leq \text{minexp}$ . Details are not given in the paper.

Then we seek to determine how the value will be rounded, more precisely, what correction will be done to the significand that has been copied just above. We currently have a significand, a trailing term  $t$  in the accumulator (bits whose exponent is in  $\llbracket \text{minexp}, u \rrbracket$ ) such that  $0 \leq t < 1 \text{ ulp}$ , and an error on the trailing term bounded by some value  $t'$  in absolute value, so that the error  $\varepsilon$  on the significand satisfies  $-t' \leq \varepsilon < 1 \text{ ulp} + t'$ . Thus one has 4 correction cases, denoted by a `corr` integer value between  $-1$  and  $2$ , which depends on  $\varepsilon$ , the sign of the significand, `rbit`, and the rounding mode:

- 1: equivalent to `nextDown`;
- 0: no correction;
- +1: equivalent to `nextUp`;
- +2: equivalent to 2 consecutive `nextUp`.

At the same time, we will also determine the ternary value and store it in `inex`. This will be the ternary value *before* the check for overflow and underflow, which is done at the very end of `sum_aux` with the `mpfr_check_range` function (this check is common to almost all MPFR functions).

To determine `corr` and the ternary value, we distinguish two cases:

- `tmd = 0`. The TMD does not occur, so that  $t'$  has no influence on the rounding and the ternary value (one can assume  $t' = 0$ ). In the directed rounding modes, one currently has: `inex = 0` if and only if  $t = 0$  (from the various cases not detailed in the paper). Therefore `inex` is the absolute value of the ternary value, and we set `corr` as follows:

- for `MPFR_RNDD`, `corr = 0`;
- for `MPFR_RNDU`, `corr = inex`;
- for `MPFR_RNDZ`, `corr = inex && !pos`;
- for `MPFR_RNDA`, `corr = inex && pos`;
- for `MPFR_RNDN`, `corr = rbit`.

We now correct the sign of the ternary value: if `inex  $\neq$  0` (i.e. `inex = 1`) and `corr = 0`, we set `inex` to  $-1$ .

- `tmd  $\neq$  0`. The TMD occurs and will be resolved by determining the sign ( $-1$ ,  $0$  or  $+1$ ) of a secondary term thanks to a second `sum_raw` invocation with a low-precision accumulator. Since the truncated sum in the accumulator was already copied to the destination, we can reuse the memory of the accumulator, but its precision is now set to `cq + dq` rounded up to the next multiple of the limb size (`GMP_NUMB_BITS`). There may remain some bits following the truncated sum, but they will be taken into account as described below.

Let us recall that the  $d - 1$  bits from exponent  $u - 2$  to  $u - d$  ( $= \text{err}$ ) are identical. We distinguish two subcases:

- $\text{err} \geq \text{minexp}$ . We keep the last two over the  $d - 1$  identical bits and the following bits, i.e. the bits from  $\text{err} + 1$  to  $\text{minexp}$ .
- $\text{err} < \text{minexp}$ . Here at least one of the identical bits is not represented, meaning that it is 0 and all these bits are 0's. Thus the accumulator is set to 0. The new `minexp` value is determined from `maxexp`, with `cq` bits reserved to avoid overflows, just like in the main sum.

Once the accumulator is set up, `sum_raw` is called with `prec = 1`, satisfying the first `sum_raw` precondition. And we have:

$$wq \geq cq + dq \geq \log n + 3 = \log n + \text{prec} + 2,$$

corresponding to the second `sum_raw` precondition. This allows us to get the correction case `corr` and the ternary value `inex` (details are not given in the paper).

We now distinguish the two cases `pos = 1` (positive sum) and `pos = 0` (negative sum), to set the sign of the MPFR number here and update the significand field to its final contents: rounding based on the correction case `corr`, change of representation at the same time in the negative case, and



clear the trailing bits. One can show that in the positive case, this corresponds to an operation of the form  $x + \text{corr}$  on the significand field, and in the negative case,  $\bar{x} + (1 - \text{corr})$  where  $\bar{x}$  is the complement; GMP does not provide such a composed complement with addition/subtraction of a limb (or similar operation), but we can do this efficiently. The only correction to do in case of the change of binade<sup>4</sup> is to set the MSB of the significand to 1 and correct the variable  $e$ .

Finally, we set the exponent of the MPFR number to  $e$ , and check the range with `mpfr_check_range`.

## VII. TESTING

Different kinds of tests are done. First, there are usual generic random tests, with limited precisions and exponent range: the exact sum is computed with basic additions (`mpfr_add`) with enough precision, then rounded to the target precision, allowing us to check the result of `mpfr_sum`. Note that this test could be able to detect bugs in either `mpfr_add` or `mpfr_sum`; it is very unlikely to get a same wrong result for both computations, because completely different algorithms are used (when the array has at least 3 regular numbers).

As usual, cases involving singular values are also tested. In particular, tests are done with an array of 6 values and every combination of values among NaN, +Inf, -Inf, +0, -0, +1 and -1.

We have some specific tests to trigger particular cases in the implementation, the goal being to have a high code coverage. For instance, the sum of 4 numbers  $i \cdot 2^{46} + j \cdot 2^{45} + k \cdot 2^{44} + f \cdot 2^{-2}$  with  $-1 \leq i, j, k \leq 1$ ,  $i \neq 0$  and  $-3 \leq f \leq 3$  is tested with the target precision chosen to have the ulp of the exact sum equal to  $2^0$  or to  $2^{44}$  (all the cases satisfying these conditions are tested).

Code (not enabled by default) has been introduced in the `mpfr_sum` implementation to be able to check some combined parameter value coverage in the TMD cases, allowing us to make sure that all allowed combinations of rounding mode, `tmd` value (1 or 2), `rbit` value, sign of the secondary term and sign of the sum are tested.

We have generic random tests with cancellations. This is done by starting with some array of random numbers, then computing a correctly rounded sum with `mpfr_sum`, and appending the opposite value to the array, so that the next `mpfr_sum` call will have cancellations. We reiterate several times.

Finally, we also have tests with underflows and overflows.

## VIII. CONCLUSION

We have designed and implemented a new algorithm to compute the correctly rounded sum of several floating-point numbers in radix 2 in arbitrary precision for GNU MPFR, where each number (the inputs and the output) has its own precision. Together with the sum, the sign of the error is returned too.

The description in the paper gives a part of a proof. Since it is almost impossible to check that a proof like that covers everything, the quality of the test suite is important. Various kinds of tests are included in MPFR, and good coverage, in particular combined parameter value coverage in some cases, is checked. Since not all C implementations and not all value combinations can be tested, a formal proof would be useful, but it would have to be expressed in a very low level.

One of the main goals was to make sure that this algorithm is efficient in any corner case. This is particularly important to avoid denial of service in a client-server system. Contrary to the initial algorithm, the worst-case complexity is now polynomial. More about this will be said in future work (there are various ways to express the complexity here).

Future work will also consist in finding real applications to check whether we may want to modify some parameters. For instance, the precision of the accumulator may be increased if need be.

## REFERENCES

- [1] IEEE, "IEEE standard for floating-point arithmetic," *IEEE Std 754-2008*, 2008. [Online]. Available: <http://dx.doi.org/10.1109/IEEESTD.2008.4610935>
- [2] —, "IEEE standard for interval arithmetic," *IEEE Std 1788-2015*, 2015. [Online]. Available: <http://dx.doi.org/10.1109/IEEESTD.2015.7140721>
- [3] J. Demmel and Y. Hida, "Fast and accurate floating point summation with application to computational geometry," *Numerical Algorithms*, vol. 37, no. 1–4, pp. 101–112, Dec. 2004. [Online]. Available: <http://dx.doi.org/10.1023/B:NUMA.0000049458.99541.38>
- [4] S. M. Rump, "Ultimately fast accurate summation," *SIAM Journal on Scientific Computing*, vol. 31, no. 5, pp. 3466–3502, 2009. [Online]. Available: <http://dx.doi.org/10.1137/080738490>
- [5] P. Kornerup, V. Lefèvre, N. Louvet, and J.-M. Muller, "On the computation of correctly rounded sums," *IEEE Transactions on Computers*, vol. 61, no. 3, pp. 289–298, Mar. 2012. [Online]. Available: <http://dx.doi.org/10.1109/TC.2011.27>
- [6] S. M. Rump, T. Ogita, and S. Oishi, "Accurate floating-point summation Part II: Sign,  $K$ -fold faithful and rounding to nearest," *SIAM Journal on Scientific Computing*, vol. 31, no. 2, pp. 1269–1302, 2008. [Online]. Available: <http://dx.doi.org/10.1137/07068816X>
- [7] U. W. Kulisch, *Advanced Arithmetic for the Digital Computer: Design of Arithmetic Units*. Springer-Verlag, Berlin, 2002.
- [8] J.-M. Muller, N. Brisebarre, F. de Dinechin, C.-P. Jeannerod, V. Lefèvre, G. Melquiond, N. Revol, D. Stehlé, and S. Torrès, *Handbook of Floating-Point Arithmetic*, 1st ed. Birkhäuser Boston, 2010. [Online]. Available: <http://www.springer.com/birkhauser/mathematics/book/978-0-8176-4704-9>
- [9] L. Fousse, G. Hanrot, V. Lefèvre, P. Péliissier, and P. Zimmermann, "MPFR: A multiple-precision binary floating-point library with correct rounding," *ACM Transactions on Mathematical Software*, vol. 33, no. 2, Jun. 2007. [Online]. Available: <http://doi.acm.org/10.1145/1236463.1236468>

<sup>4</sup>This can be detected with a carry/borrow out of a GMP operation.